

MATH 587/589 COURSE NOTES

LOUIGI ADDARIO-BERRY

Abstract. **Course notes for Math 587/Math 589.**

These notes cover everything I teach in Math 587, and have substantial overlap with what I teach in Math 589, though the symmetric difference with the Math 589 content is nonempty.

Contents

1. Acknowledgements	3
2. Notation	3
3. Why measure theory?	3
4. Measure theory	5
4.1. Rings, fields and σ-fields	5
4.2. Building measures	5
4.3. Measures on \mathbb{R}	11
4.4. Independent events	13
5. Random variables	16
5.1. Generated σ-fields	18
5.2. Independence of random variables	19
5.3. Existence of random variables with given distributions	19
5.4. Kolmogorov's zero-one law	21
5.5. Almost sure convergence, convergence in probability and convergence in distribution	22
6. Integration and expectation	25
6.1. Expectation and independence	30
7. An interlude: the probabilistic method.	33
8. Densities and change of variables	35
8.1. Product measure and Fubini's theorem	37
9. Sums of independent random variables	43
9.1. Convolutions	43
10. Laws of large numbers	44
11. Convexity, inequalities, and L_p spaces	52
11.1. The geometric structure of L_2	55
12. Conditional expectation	59
12.1. Properties of conditional expectation	63
12.2. Conditional expectations, tightness and uniform integrability	67
13. Martingales	68
13.1. Martingale convergence theorems	72
Uniform integrability	75
13.2. Maximal inequalities and the L_p martingale convergence theorem.	78
13.3. Filtrations and changes of measure	80
14. Branching process limits	83

Date: September 4, 2023.

Copyright 2019; reproduction prohibited without permission of the author.

14.1.	Branching process recap	83
14.2.	Branching processes with immigration	86
14.3.	The Kesten-Stigum theorem	88
15.	Transforms 1: Moment-generating functions	90
15.1.	Introduction	90
15.2.	The moment generating function	91
15.3.	The moment generating function and uniqueness	94
16.	Transforms 2: Characteristic functions	95
16.1.	Characteristic functions: basic properties and examples	97
16.2.	The inversion and continuity theorems	99
16.3.	The central limit theorem	104
16.4.	Characteristic functions and moments	105
17.	Weak convergence	110
17.1.	Measures on metric spaces; the Portmanteau theorem	110
17.2.	Weights and measures	112
17.3.	Aside: the existence of non-tight probability measures	113
17.4.	Bounded Lipschitz functions and $\text{prob}(\mathcal{M})$	114
17.5.	Recursively partitioning Polish spaces	116
17.6.	Metrizing weak convergence	119
17.7.	Probability kernels and conditional probabilities	121
	List of notation and terminology	124
	References	125

1. Acknowledgements

Thanks to Huangchen Zhou for spotting many typos and making useful suggestions.

2. Notation

We write \mathcal{L}_X or μ_X (**Make this consistent.**) for the distribution of X . Given a σ -finite measure μ on \mathbb{R} and $p > 0$, write $|\mu|_p = \left(\int_{\mathbb{R}} |x|^p d\mu(x)\right)^{1/p}$. If μ is a probability distribution and X has law μ then $|\mu|_p = (\mathbf{E}[|X|^p])^{1/p}$. For a random variable X we'll write $\|X\|_p = (\mathbf{E}[|X|^p])^{1/p}$; if the space in question is ambiguous we may instead write $\|X\|_{L_p(\Omega, \mathcal{F}, \mathbf{P})}$.

3. Why measure theory?

This is a somewhat hard question to answer rigorously, but I recently (2020) learned about the connection between *hat problems*, the axiom of choice, and measurability, that I find helpful as a motivation. It works as follows.

First, here is a cooperative game for n players, labeled $1, \dots, n$. Each player is wearing either a white (0) or a black (1) hat, and can see the colour of all the hats aside from their own. In order from $1, \dots, n$, the players guess the colour of their own hat. (No communication is allowed aside from observing the guesses of the previous players.) The players win if everyone except for player 1 guesses correctly.

A winning strategy for the players is as follows. Player 1 reports the parity of the number of black hats worn by players $2, \dots, n$. Based on this, everyone can deduce the colour of their own hat (because player i is wearing a black hat iff player i sees a different parity of black hats among $2, \dots, n$ than that reported by player 1).

The infinite player version is similar to the finite-player version, except that there are infinitely many players (so the set of players is now $\mathbb{N} = \{1, 2, 3, \dots\}$). The serial (sequential) game is still interesting, but the parallel game is easier to analyze, so we focus on that. We thus insist that all players guess simultaneously and that no communication is allowed (equivalently, no one has access to the guesses of anyone else when making their own guess). For this infinite game, we will relax the winning condition - now say that the players win if all but a finite number of players guess correctly.

The winning strategy for the infinite game requires the axiom of choice.¹ The axiom of choice states that for any collection $(S_i, i \in I)$ of nonempty sets indexed by some set I , there exists a collection $s = (s_i, i \in I)$ with $s_i \in S_i$ for all $i \in I$; so s is a function “choosing” one element from each of the sets S_i .

We use the axiom of choice as follows. Each assignment of hats to players can be represented by an element of $\{0, 1\}^{\mathbb{N}}$. Say that two assignments ω, ω' are equivalent, and write $\omega \sim \omega'$, if ω and ω' differ in only finitely many coordinates. Then \sim is an equivalence relation, so defines a partition $(S_i, i \in I)$ of $\{0, 1\}^{\mathbb{N}}$, where ω, ω' lie in the same part of the partition if and only if $\omega \sim \omega'$.

Let $(s_i, i \in I)$ be such that $s_i \in S_i$ for all $i \in I$; this chooses for us a representative for each equivalence class. Write $s_i = (s_i(j), j \in \mathbb{N})$, so $s_i(j) = 1$ if player j has a black hat in assignment s_i , and $s_i(j) = 0$ otherwise. Now define a strategy as follows. Given $\omega \in \{0, 1\}^{\mathbb{N}}$, let $i = i(\omega) \in I$ be such that $\omega \sim s_i$; note that all players can deduce i from observing the hat assignment ω , since i is the unique index for which ω and s_i differ in only finitely many entries. Then player j guesses that their hat has colour $s_i(j)$. Since $\omega \sim s_i$, all but finitely many players will have $s_i(j) = \omega(j)$ and so will guess correctly; so this is a winning strategy.

It may seem counterintuitive that a winning strategy could exist, since no player can see their own hat, and moreover, changing the colour of any given hat doesn't change the strategy. And, indeed, we now show that any *measurable* strategy has probability at least 1/2 of failing, when the hat colours are independent and each hat is equally likely to be black or to be white. The measure

¹I use the word “required” advisedly; see Theorem 8 of *An introduction to infinite hat problems*, Hardin and Taylor, 2008; <https://link.springer.com/article/10.1007/BF03038092>

space modelling this is $(\{0, 1\}^{\mathbb{N}}, \mathcal{F}, \mathbb{P})$, where \mathcal{F} is the σ -field generated by the cylinder sets and \mathbb{P} is determined by the condition that for any $(x_1, \dots, x_n) \in \{0, 1\}^n$,

$$\mathbf{P} \left\{ \{\omega \in \{0, 1\}^{\mathbb{N}} : (\omega_1, \dots, \omega_n) = (x_1, \dots, x_n)\} \right\} = 2^{-n}.$$

A measurable strategy for player i is a measurable map $X_i : \{0, 1\}^{\mathbb{N}} \rightarrow \{0, 1\}$ such that the events $\{X_i = 1\}$ and $\{\omega_i = 1\}$ are independent; this condition reflects the fact that player i can not see the colour of their own hat.

So fix any measurable strategies $(X_i, i \in \mathbb{N})$ for the players. Then for all i , the above independence implies that

$$\mathbf{P} \{X_i \neq \omega_i\} = 1/2,$$

so by linearity of expectation, for all $n \in \mathbb{N}$,

$$\mathbf{E}|\{i \in [n] : X_i \neq \omega_i\}| = n/2.$$

Fixing $K \in \mathbb{N}$, the above identity and the inequality

$$\mathbf{E}|\{i \in [n] : X_i \neq \omega_i\}| < K + n\mathbf{P} \{|\{i \in [n] : X_i \neq \omega_i\}| \geq K\}$$

together imply that

$$\mathbf{P} \{|\{i \in [n] : X_i \neq \omega_i\}| \geq K\} > \frac{n - K}{2n - K}.$$

Since

$$\lim_{n \rightarrow \infty} \{|\{i \in [n] : X_i \neq \omega_i\}| \geq K\} = \{|\{i \in \mathbb{N} : X_i \neq \omega_i\}| \geq K\},$$

by the monotone convergence theorem we deduce that

$$\mathbf{P} \{|\{i \in \mathbb{N} : X_i \neq \omega_i\}| \geq K\} = \lim_{n \rightarrow \infty} \mathbf{P} \{|\{i \in [n] : X_i \neq \omega_i\}| \geq K\} \geq 1/2.$$

Next, the events $\{|\{i \in \mathbb{N} : X_i \neq \omega_i\}| \geq K\}$ are decreasing in K , and their limit is the event $\{|\{i \in \mathbb{N} : X_i \neq \omega_i\}| = \infty\}$, so by the dominated convergence theorem (or the monotone convergence theorem applied to the complementary events) we obtain that

$$\mathbf{P} \{|\{i \in \mathbb{N} : X_i \neq \omega_i\}| = \infty\} \geq 1/2,$$

as claimed.

On the one hand, we have shown that there exists a strategy which always succeeds. On the other hand, we have shown that any *measurable* strategy fails with positive probability. The probabilistic tools we use in the proof are very basic: independence, linearity of expectation, and the monotone convergence theorem. It is hard to see how to do much probability without those tools. So if we are to develop probability theory, either we need to find some substantially different approach, or measure theory is required.

Conjecture. Any measurable strategy fails with probability 1. More formally: consider the probability space $(\Omega, \mathcal{F}, \mathbf{P})$ with $\Omega = \{0, 1\}^{\mathbb{N}}$, \mathcal{F} the σ -field generated by the cylinder sets, and \mathbb{P} the measure under which

$$\mathbf{P} \{(\omega_1, \dots, \omega_n) = (x_1, \dots, x_n)\} = 2^{-n}$$

for all $n \geq 1$ and all $(x_1, \dots, x_n) \in \{0, 1\}^n$. For $i \geq 1$ let $B_i : \Omega \rightarrow \{0, 1\}$ be the i 'th coordinate map: $B_i(\omega) = \omega_i$.

Fix any measurable maps $(X_i, i \geq 1)$ with $X_i : \Omega \rightarrow \{0, 1\}$ such that

$$\mathbf{P} \{X_i = B_i\} = 1/2$$

for all $i \geq 1$. Then

$$\mathbf{P} \{X_i \neq B_i \text{ for infinitely many } i\} = 1.$$

I don't know how to prove this, though I haven't spent a long time trying. I'll give some bonus % to the course grade of the first person or group who solves this; since I don't yet know how hard the question is, I'll determine the amount of the bonus after seeing the solution.

4. Measure theory

Measure theory is the algebraic underpinning of probability theory. It can feel rather abstract; but it is worth setting things up clearly.

4.1. **Rings, fields and σ -fields.** Fix a set Ω and a set \mathcal{A} of subsets of Ω with $\emptyset \in \mathcal{A}$. We say \mathcal{A} is a *ring* if the following hold.

- (a) If $E, F \in \mathcal{A}$ then $E \cup F \in \mathcal{A}$.
- (b) If $E, F \in \mathcal{A}$ then $F \setminus E \in \mathcal{A}$.

We say \mathcal{A} is a π -system if the following holds.

- (c) If $E, F \in \mathcal{A}$ then $E \cap F \in \mathcal{A}$.

We say \mathcal{A} is a *field* if it is a ring and also the following holds.

- (d) If $E \in \mathcal{A}$ then $E^c \in \mathcal{A}$.

We say \mathcal{A} is a σ -field if it is a field and also the following holds.

- (a') For any sequence $(A_n, n \geq 1)$ of elements of \mathcal{A} , $\bigcup_{n \geq 1} A_n \in \mathcal{A}$.

In all the above cases, we refer to Ω as the *ground set*. Finally, for an arbitrary set \mathcal{A} of subsets of Ω , the σ -field generated by \mathcal{A} is

$$\sigma(\mathcal{A}) := \bigcap_{\{\mathcal{F} \supset \mathcal{A} : \mathcal{F} \text{ a } \sigma\text{-field}\}} \mathcal{F};$$

this is the smallest σ -field containing the set \mathcal{A} .

Exercise 4.1. (i) Show that properties (a) and (d) together imply properties (b) and (c).
 (ii) Show that a field which is closed under countable disjoint unions is a σ -field.

Exercise 4.2. Write $\mathbb{N} := \{1, 2, 3, \dots\}$. For $n \in \mathbb{N}$ we let $[n] := \{1, 2, \dots, n\}$. Say that $S \subset \mathbb{N}$ has an asymptotic density if

$$\mu(S) := \limsup_{n \rightarrow \infty} \frac{|S \cap [n]|}{n} = \liminf_{n \rightarrow \infty} \frac{|S \cap \{1, 2, \dots, n\}|}{n}.$$

Write \mathcal{A} for the set of subsets of \mathbb{N} which have an asymptotic density. Is \mathcal{A} a π -system? Is it a ring? A field? A σ -field?

4.2. **Building measures.** A *measurable space* is a pair (Ω, \mathcal{F}) , where \mathcal{F} is a σ -field over Ω . Given such a space, a *measure* μ on \mathcal{F} is a function $\mu : \mathcal{F} \rightarrow [0, \infty]$ such that $\mu(\emptyset) = 0$, and for any sequence $(A_n, n \geq 1)$ of disjoint elements of \mathcal{F} ,

$$\mu \left(\bigcup_{n \geq 1} A_n \right) = \sum_{n \geq 1} \mu(A_n).$$

We then call $(\Omega, \mathcal{F}, \mu)$ a *measure space*. You should think of a measure space as a model for a physical system involving randomness. Sometimes this can be quite concrete. For example, one might take $\Omega = [6]$, $\mathcal{F} := 2^\Omega$ is the power set of Ω , and $\mu(S) = |S|/6$, to model the roll of a fair die; here $\mu(S)$ is the probability that the roll yields a value in S . If we took $\Omega = [6]^{[2]} = \{(i, j) : i, j \in [6]\}$ and $\mu(S) = |S|/36$, we could view this as modelling two successive rolls of a fair die.

On the other hand, when doing probability it is often useful to leave the details of the measure space rather implicit. There are various tools which justify doing this (change of variables, existence theorems,...), which we'll see later.

Exercise 4.3. Let μ be a measure on a σ -field \mathcal{F} .

- (i) **[Monotone convergence/Continuity from below.]** Show that for any increasing sequence $(E_n, n \geq 1)$ of elements of \mathcal{F} , it holds that $\mu(\bigcup_{n \geq 1} E_n) = \lim_{n \rightarrow \infty} \mu(E_n)$.
- (ii) **[Dominated convergence/Continuity from above.]** Show that for any decreasing sequence $(E_n, n \geq 1)$ of elements of \mathcal{F} with $\mu(E_1) < \infty$, it holds that $\mu(\bigcap_{n \geq 1} E_n) = \lim_{n \rightarrow \infty} \mu(E_n)$.

Ring

π -system

Field

σ -field

$\sigma(\mathcal{A})$

Throughout these notes, "countable" means "finite or countably infinite".

Measurable space

Measure

Measure space

(iii) **[Subadditivity.]** Show that for any sequence $(E_n, n \geq 1)$ of elements of \mathcal{F} , it holds that $\mu(\bigcup_{n \geq 1} E_n) \leq \sum_{n \geq 1} \mu(E_n)$.

Exercise 4.4. Which of the following triples $(\Omega, \mathcal{F}, \mu)$ are measure spaces? Can you think of physical systems which they model?

- (a) $\Omega = \mathbb{N}$, \mathcal{F} the set of subsets of \mathbb{N} which have an asymptotic density, $\mu(S)$ the asymptotic density of S .
- (b) $\Omega = \{0, 1\}^{\mathbb{N}}$, \mathcal{F} the power set of Ω , $\mu(\{\omega\}) = p^{|\{i \in [n]: \omega_i = 1\}|} (1-p)^{|\{i \in [n]: \omega_i = 0\}|}$, where $p \in (0, 1)$ is fixed.
- (c) $\Omega = \{0, 1\}^{\mathbb{N}}$, \mathcal{F} the power set of Ω , $\mu(\omega) = p^{|\{i \in [n]: \omega_i = 1\}|}$, where $p \in [0, 1]$ is fixed.
- (d) $\Omega = [0, 1]$, \mathcal{F} the collection of sets $S \subset [0, 1]$ such that either S or $[0, 1] \setminus S$ is countable, and $\mu(S) = |S|$.

You have likely seen probability distributions described in terms of *cumulative distribution functions* (CDFs). For example, the standard exponential distribution has CDF $F(x) = (1 - e^{-x})\mathbf{1}_{[x \geq 0]}$, corresponding to the fact that for E is a standard exponential random variable, $\mathbf{P}\{E \leq x\} = (1 - e^{-x})\mathbf{1}_{[x \geq 0]}$.² What $F(x)$ lets us easily compute is probabilities of the form $\mathbf{P}\{E \in (a, b]\} = F(b) - F(a)$, or $\mathbf{P}\{E \in \bigcup_{i=1}^n (a_i, b_i]\} = \sum_{i=1}^n (F(b_i) - F(a_i))$, where $(a_1, b_1], \dots, (a_n, b_n]$ are disjoint intervals. On the other hand, it's not clear how we would use the above CDF to determine $\mathbf{P}\{E \in \mathbb{Q}\}$, for example, although we know the answer must be zero. If we are going to specify probability distributions in this way, we should really prove that probability measures are uniquely determined by their CDFs; this is a corollary of the coming development.

Indicator of a set

Fix a ring \mathcal{A} over a ground set Ω . A *pre-measure on \mathcal{A}* is a function $\mu : \mathcal{A} \rightarrow [0, \infty]$ with $\mu(\emptyset) = 0$ such that for any sequence $(A_n, n \geq 1)$ of disjoint elements of \mathcal{F} , if $\bigcup_{n \geq 1} A_n \in \mathcal{A}$ then

Pre-measure

$$\mu\left(\bigcup_{n \geq 1} A_n\right) = \sum_{n \geq 1} \mu(A_n).$$

We then say that $(\Omega, \mathcal{A}, \mu)$ is a pre-measure space.

Pre-measure space

Here is a key example of a pre-measure space. We hereafter write

 $\mathcal{A}(\mathbb{R})$

$$\mathcal{A}(\mathbb{R}) = \left\{ \bigcup_{i=1}^n (a_i, b_i] : n \geq 1, -\infty < a_1 \leq b_1 \leq a_2 \leq \dots \leq a_n \leq b_n < \infty \right\}.$$

Exercise 4.5. Prove that $\mathcal{A}(\mathbb{R})$ is a ring over \mathbb{R} .

We will see later that if F is a CDF then we can define a function μ on $\mathcal{A}(\mathbb{R})$ by setting $\mu(\bigcup_{i=1}^n (a_i, b_i]) = \sum_{i=1}^n (F(b_i) - F(a_i))$ when $((a_i, b_i], n \geq 1)$ are pairwise disjoint, and that the resulting triple $(\mathbb{R}, \mathcal{A}(\mathbb{R}), \mu)$ is a pre-measure space. The primordial³ existence theorem for measures is the following.

Theorem 4.1 (Carathéodory extension theorem). *Let $(\Omega, \mathcal{A}, \mu)$ be a pre-measure space. Then there exists a σ -field \mathcal{F} containing \mathcal{A} such that μ extends to a measure on \mathcal{F} .*

The previous theorem provides existence; the next theorem provides uniqueness.

Theorem 4.2 (Dynkin's theorem). *Let (Ω, \mathcal{F}) be a measurable space, and let $\mathcal{P} \subset \mathcal{F}$ be a π -system $\Omega \in \mathcal{P}$ and with $\sigma(\mathcal{P}) = \mathcal{F}$. Fix measures μ_1, μ_2 on \mathcal{F} , and suppose that (a) $\mu_1(E) = \mu_2(E)$ for all $E \in \mathcal{P}$ and (b) there exist sets $(\Omega_n, n \geq 1)$ in \mathcal{P} with $\Omega_n \uparrow \Omega$ as $n \rightarrow \infty$ and with $\mu_1(\Omega_n) < \infty$. Then $\mu_1 \equiv \mu_2$.*

²Here and throughout these notes, for a set S and a subset $T \subset S$, we write $I_T : S \rightarrow \{0, 1\}$ for the indicator of set T , so $I_T(x) = 1$ for $x \in T$ and $I_T(x) = 0$ otherwise.

³Primordial, adj. and n.: That constitutes the origin or starting point from which something else is derived or developed, or on which something else depends; fundamental, basic; elemental. —Oxford English Dictionary

The proof of the Carathéodory extension theorem consists of two parts. Starting from the pre-measure space $(\Omega, \mathcal{A}, \mu)$ provided by the hypothesis of the theorem, we first use the pre-measure μ provided by to produce an upper bound on any putative⁴ extension of μ . Next we show that the upper bound indeed yields a measure on a σ -field extending the ring \mathcal{A} .

Proposition 4.3. *Let $(\Omega, \mathcal{A}, \mu)$ be a pre-measure space. For $B \subset \Omega$ let*

$$\mu^*(B) := \inf \left(\sum_{n \geq 1} \mu(A_n) : A_n \in \mathcal{A}, n \geq 1; B \subset \bigcup_{n \geq 1} A_n \right).$$

Then μ^* is an outer measure:

- (i) $\mu^*(\emptyset) = 0$;
- (ii) If $E \subset F$ then $\mu^*(E) \leq \mu^*(F)$;
- (iii) If $(E_i, i \geq 1)$ are subsets of Ω then $\mu^*(\bigcup_{i \geq 1} E_i) \leq \sum_{i \geq 1} \mu^*(E_i)$.

Note: usually I try to avoid putting definitions within the statements of Theorems, Propositions, Lemmas and so forth; but this is almost the only place where outer measures will be used.

Lemma 4.4 (Carathéodory lemma). *Given an outer measure μ^* over a set Ω , say $A \subset \Omega$ is μ^* -additive if for all $B \subset \Omega$,*

$$\mu^*(B) = \mu^*(A \cap B) + \mu^*(A^c \cap B).$$

Let $\mathcal{F} = \{A \subset \Omega : A \text{ is } \mu^*\text{-additive}\}$, and define $\mu^+ : \mathcal{F} \rightarrow [0, \infty]$ by $\mu^+(A) := \mu^*(A)$. Then $(\Omega, \mathcal{F}, \mu^+)$ is a measure space.

I think of μ^* -additive sets as knives; they “sharply cut” any set $B \subset \Omega$ in two without any change of μ^* -measure. I’m not sure how useful this perspective is to others.

Proof of Proposition 4.3. Point (i) is obvious since the empty set is a cover of itself. Point (ii), monotonicity, is also obvious, since if $E \subset F$ then any cover of F is a cover of E , so $\mu^*(E)$ is an infimum over a larger set than $\mu^*(F)$.

Finally, fix subsets $(E_i, i \geq 1)$ of Ω and write $E = \bigcup_{i \geq 1} E_i$. Next fix $\epsilon > 0$, and for each $i \geq 1$, fix a cover $(A_n^i, n \geq 1)$ of E_i with

$$\sum_{n \geq 1} \mu(A_n^i) \leq \mu^*(E_i) + \frac{\epsilon}{2^i};$$

such a cover exists by the definition of $\mu^*(E_i)$. Then $(A_n^i, n, i \geq 1)$ is a cover of E , so

$$\begin{aligned} \mu^*(E) &\leq \sum_{n, i \geq 1} \mu(A_n^i) \\ &\leq \sum_{i \geq 1} \left(\mu^*(E_i) + \frac{\epsilon}{2^i} \right) \\ &= \sum_{i \geq 1} \mu^*(E_i) + \epsilon. \end{aligned}$$

Since $\epsilon > 0$ was arbitrary it follows that $\mu^*(E) \leq \sum_{i \geq 1} \mu^*(E_i)$. □

Proof of Lemma 4.4. The conclusion of the Carathéodory lemma is that \mathcal{F} is a σ -field over Ω and μ^+ is a measure on \mathcal{F} . We prove these in order.

First, for any $B \subset \Omega$ we have $\mu^*(\emptyset \cap B) + \mu^*(\Omega \cap B) = \mu^*(\emptyset) + \mu^*(B) = \mu^*(B)$, so $\emptyset \in \mathcal{F}$. Also, the definition of μ^* additive sets is invariant under complementation, so $A \in \mathcal{F}$ if and only if $A^c \in \mathcal{F}$.

⁴Putative, adj.: That is commonly believed to be such; reputed, supposed; imagined; postulated, hypothetical. –Oxford English Dictionary

Outer measure

Lemma 4.4 applies to any outer measure, not just μ^* ; change notation?

We next show \mathcal{F} is closed under intersections. Fix any sets $A_1, A_2 \in \mathcal{F}$. For any $B \subset \Omega$, we may write B as a disjoint union

$$\begin{aligned} B &= B_0 \cup B_1 \cup B_2 \cup B_{12}, \text{ where} \\ B_0 &= B \cap A_1^c \cap A_2^c, \\ B_1 &= B \cap A_1 \cap A_2^c, \\ B_2 &= B \cap A_1^c \cap A_2, \text{ and} \\ B_{12} &= B \cap A_1 \cap A_2. \end{aligned}$$

This “cuts B into four pieces”, according to its intersection with B_1 and B_2 . Since A_1 and A_2 are μ^* -additive, we have

$$\begin{aligned} \mu^*(B) &= \mu^*(B \cap A_1) + \mu^*(B \cap A_1^c) \\ &= \mu^*(B_{12}) + \mu^*(B_1) + \mu^*(B_2) + \mu^*(B_0) \end{aligned}$$

If we likewise cut $B \setminus B_{12}$ into four pieces, only the last three pieces will be non-empty, and we obtain

$$\mu^*(B \setminus B_{12}) = \mu^*(B_0) + \mu^*(B_2) + \mu^*(B_1).$$

Together the last two equations give that

$$\mu^*(B) = \mu^*(B_{12}) + \mu^*(B \setminus B_{12}) = \mu^*(B \cap A_1 \cap A_2) + \mu^*(B \cap (A_1 \cap A_2)^c).$$

Thus $A_1 \cap A_2 \in \mathcal{F}$.

At this point we know \mathcal{F} is a field, so to show it is a σ -field it suffices to establish that it is closed under countable disjoint unions. Fix a sequence $(A_i, n \geq 1)$ of disjoint sets in \mathcal{F} , and any set $B \subset \Omega$. Writing $A = \bigcup_{i \geq 1} A_i$, we must show that for all $B \subset \Omega$ we have $\mu^*(B) = \mu^*(A \cap B) + \mu^*(A^c \cap B)$. The fact that $\mu^*(B) \leq \mu^*(A \cap B) + \mu^*(A^c \cap B)$ is immediate by subadditivity of outer measure, so we only need to show the reverse inequality.

We will again “cut B into pieces” according to its intersection with the sets A_n . However, since the sets are disjoint, our task is now simpler; we may rewrite

$$\begin{aligned} \mu^*(B) &= \mu^*(B_{12}) + \mu^*(B_1) + \mu^*(B_2) + \mu^*(B_0) \\ &= \mu^*(B_1) + \mu^*(B_2) + \mu^*(B_0) \\ &= \mu^*(B \cap A_1) + \mu^*(B \cap A_2) + \mu^*(B \cap A_1^c \cap A_2^c). \end{aligned}$$

More generally, since $A_n \cap B \cap A_1^c \cap \dots \cap A_{n-1}^c = A_n \cap B$, by induction we have

$$\mu^*(B) = \mu^*(B \cap A_1^c \cap \dots \cap A_n^c) + \sum_{i=1}^n \mu^*(B \cap A_i)$$

for all n . Now, since $A^c \subset A_1^c \cap \dots \cap A_n^c$ we have $\mu^*(B \cap A_1^c \cap \dots \cap A_n^c) \geq \mu^*(B \cap A^c)$ by the monotonicity of outer measure, so

$$\mu^*(B) \geq \mu^*(B \cap A^c) + \sum_{i=1}^n \mu^*(B \cap A_i);$$

taking a limit in n now gives

$$\mu^*(B) \geq \mu^*(B \cap A^c) + \sum_{i=1}^{\infty} \mu^*(B \cap A_i).$$

Since $A = \bigcup_{i \geq 1} A_i$, by subadditivity of outer measure we have $\mu^*(B \cap A) \leq \sum_{i=1}^{\infty} \mu^*(B \cap A_i)$, which with the previous bound gives

$$\mu^*(B) \geq \mu^*(B \cap A^c) + \mu^*(B \cap A).$$

Thus $\mu^*(B) = \mu^*(B \cap A^c) + \mu^*(B \cap A)$, so $A \in \mathcal{F}$ and \mathcal{F} is indeed a σ -field.

Finally, note that in proving \mathcal{F} is a σ -field, we also established that the restriction μ^+ of μ^* to \mathcal{F} is countably additive on \mathcal{F} . Also, μ^+ is monotone and has $\mu^+(\emptyset) = 0$ by definition; so μ^+ is a measure on \mathcal{F} , as required. \square

Proof of Theorem 4.1. Let μ^* be as in Proposition 4.3; then μ^* is an outer measure. We first verify that μ^* agrees with μ on \mathcal{A} . Fix $A \in \mathcal{A}$ and any sequence $(A_i, i \geq 1)$ of elements of \mathcal{A} which cover A . Writing $B_n = A_n \setminus (A_1 \cup \dots \cup A_n)$, then $B_n \subset A_n$ and $(A \cap B_n, n \geq 1)$ is another cover of A with elements of \mathcal{A} . Since $A = \bigcup_{n \geq 1} A \cap B_n$, by countable additivity for pre-measures we have

$$\begin{aligned} \mu(A) &= \sum_{n \geq 1} \mu(A \cap B_n) \\ &\leq \sum_{n \geq 1} \mu(B_n) \\ &\leq \sum_{n \geq 1} \mu(A_n). \end{aligned}$$

Taking the infimum over covers $(A_n, n \geq 1)$ of A we obtain that $\mu(A) \leq \mu^*(A)$. Also, clearly $\mu(A) \geq \mu^*(A)$ since A covers itself; so $\mu(A) = \mu^*(A)$.

Next, let \mathcal{F} be the collection of μ^* -additive sets, and let μ be the restriction of μ^* to \mathcal{F} ; we are recycling notation here but this is OK since we already checked that μ and μ^* agree on their common domain of definition. By Lemma 4.4, $(\Omega, \mathcal{F}, \mu^*)$ is a measure space, and by the first paragraph we know that μ^* agrees with μ on \mathcal{A} . So, to complete the proof of the theorem it remains to show that $\mathcal{A} \subset \mathcal{F}$, or in other words that the sets in \mathcal{A} are μ^* -additive.

So fix any set $A \in \mathcal{A}$ and any set $B \subset \Omega$. By subadditivity of μ^* we have

$$\mu^*(B) \leq \mu^*(A \cap B) + \mu^*(A^c \cap B);$$

we need to prove the reverse inequality. If $\mu^*(B) = \infty$ then this is obvious, so we suppose $\mu^*(B) < \infty$. Fix $\epsilon > 0$; then we may find a cover $(A_n, n \geq 1)$ of B with elements of \mathcal{A} such that

$$\sum_{n \geq 1} \mu(A_n) \leq \mu^*(B) + \epsilon$$

Finally, $(A \cap A_n, n \geq 1)$ is a cover of $A \cap B$ with elements of \mathcal{A} , and $(A^c \cap A_n, n \geq 1)$ is a cover of $A^c \cap B$ with elements of \mathcal{A} , so from the definition of μ^* we have

$$\begin{aligned} \mu^*(A \cap B) + \mu^*(A^c \cap B) &\leq \sum_{n \geq 1} \mu(A \cap A_n) + \sum_{n \geq 1} \mu(A^c \cap A_n) \\ &= \sum_{n \geq 1} (\mu(A \cap A_n) + \mu(A^c \cap A_n)) \\ &= \sum_{n \geq 1} \mu(A_n) \\ &\leq \mu(B) + \epsilon. \end{aligned}$$

Since $\epsilon > 0$ was arbitrary, it follows that $\mu^*(A \cap B) + \mu^*(A^c \cap B) \leq \mu^*(B)$, as required. \square

The proof of Theorem 4.2 relies on one more algebraic/set theoretic closure property, which we now state. We say a collection \mathcal{A} of subsets of a ground set Ω is a λ -system if $\Omega \in \mathcal{A}$ and additionally the following both hold.

- (i) For all $E, F \in \mathcal{A}$ with $E \subset F$ we have $F \setminus E \in \mathcal{A}$.
- (ii) For any increasing sequence $(A_n, n \geq 1)$ of subsets of \mathcal{A} we have $\bigcup_{n \geq 1} A_n \in \mathcal{A}$.

By *increasing* we mean that $A_n \subset A_{n+1}$ for all $n \geq 1$.

Exercise 4.6. (i) If \mathcal{A} is a σ -field then it is a π -system.
 (ii) If \mathcal{A} is both a π -system and a λ -system then it is a σ -field.
 (iii) Fix any collection $\{\mathcal{A}_i, i \in I\}$ of λ -systems with a common ground set. Then $\bigcap_{i \in I} \mathcal{A}_i$ is a λ -system.

λ -system

Part (i) is similar to Exercise 2.1.

Lemma 4.5 (Dynkin's π -system lemma). *Let \mathcal{P} be a π -system over a ground set Ω . Then*

$$\bigcap_{\{\mathcal{F} \supset \mathcal{P}: \mathcal{F} \text{ a } \sigma\text{-field}\}} \mathcal{F} = \bigcap_{\{\mathcal{F} \supset \mathcal{P}: \mathcal{F} \text{ a } \lambda\text{-system}\}} \mathcal{F}$$

Proof of Lemma 4.5. The left-hand side is $\sigma(\mathcal{P})$ by definition; temporarily writing $\lambda(\mathcal{P})$ for the right-hand side, we aim to show that $\sigma(\mathcal{P}) = \lambda(\mathcal{P})$.

Since σ -fields are λ -systems we automatically have $\sigma(\mathcal{P}) \supset \lambda(\mathcal{P})$, so to prove the lemma it suffices to show that $\lambda(\mathcal{P})$ is a σ -field. Moreover $\lambda(\mathcal{P})$ is a λ -system by Exercise 4.6 (ii), so by part (iii) of the same exercise, to show it is a σ -field we just have to show it is closed under intersections.

We proceed in two steps. For $E \in \lambda(\mathcal{P})$, say E is *cooperative* if $E \cap F \in \lambda(\mathcal{P})$ for all $F \in \mathcal{P}$, and that E is *helpful* if $E \cap F \in \lambda(\mathcal{P})$ for all $F \in \lambda(\mathcal{P})$.

If $E \in \mathcal{P}$ then $E \cap F \in \mathcal{P}$ for all $F \in \mathcal{P}$ since \mathcal{P} is a π -system; so E is cooperative. Next, if E and E' are both cooperative and $E \subset E'$ then for all $F \in \mathcal{P}$ we have $E \cap F \in \lambda(\mathcal{P})$ and $E' \cap F$ in $\lambda(\mathcal{P})$. Since $\lambda(\mathcal{P})$ is a λ -system, it follows that

$$(E' \setminus E) \cap F = (E' \cap F) \setminus (E \cap F) \in \lambda(\mathcal{P}),$$

so $E' \setminus E$ is cooperative. Third, if $(E_n, n \geq 1)$ is an increasing sequence of cooperative sets then for all $F \in \mathcal{P}$ we have

$$F \cap \bigcup_{n \geq 1} E_n = \bigcup_{n \geq 1} F \cap E_n.$$

Each of the sets $F \cap E_n$ lies in $\lambda(\mathcal{P})$ since the E_n are cooperative. Since $(F \cap E_n, n \geq 1)$ is increasing and $\lambda(\mathcal{P})$ is a λ -system, it follows that $F \cap \bigcup_{n \geq 1} E_n \in \lambda(\mathcal{P})$, so $\bigcup_{n \geq 1} E_n$ is cooperative.

We've now showed that the cooperative sets in $\lambda(\mathcal{P})$ contain \mathcal{P} and are closed under monotone difference and monotone limits: they are a λ -system; so all sets in $\lambda(\mathcal{P})$ are cooperative.

We now bootstrap this argument. If $E \in \mathcal{P}$ then for any $F \in \lambda(\mathcal{P})$, since F is cooperative we have $E \cap F \in \lambda(\mathcal{P})$; so E is in fact helpful. Next, if E, E' are helpful and $E \subset E'$ then for all $F \in \lambda(\mathcal{P})$, $E' \cap F$ and $E \cap F$ both lie in $\lambda(\mathcal{P})$, so

$$(E' \setminus E) \cap F = (E' \cap F) \setminus (E \cap F) \in \lambda(\mathcal{P}).$$

Finally, if $(E_n, n \geq 1)$ is an increasing sequence of helpful sets then for all $F \in \lambda(\mathcal{P})$ and all $n \in \mathbb{N}$, $F \cap E_n \in \lambda(\mathcal{P})$, so

$$F \cap \bigcup_{n \geq 1} E_n = \bigcup_{n \geq 1} F \cap E_n \in \lambda(\mathcal{P}).$$

Thus $\bigcup_{n \geq 1} E_n$ is helpful. We've just showed that the helpful sets are a λ -system containing \mathcal{P} , so all sets in $\lambda(\mathcal{P})$ are helpful. But this means that means that $E \cap F \in \lambda(\mathcal{P})$ for all $E, F \in \lambda(\mathcal{P})$; so $\lambda(\mathcal{P})$ is closed under intersections, as required. \square

We now show that Lemma 4.5 easily yields Theorem 4.2.

Proof of Theorem 4.2. Let μ_1, μ_2 be as in the theorem's statement. Fix any set $G \in \mathcal{P}$ with $\mu_1(G) < \infty$, and write $\Lambda = \{E \in \mathcal{F} : \mu_1(E \cap G) = \mu_2(E \cap G)\}$; then Λ contains \mathcal{P} by definition, and in particular $\Omega \in \Lambda$.

Next, if $(E_n, n \geq 1)$ is an increasing sequence of sets in Λ then

$$\mu_1\left(\bigcup_{n \geq 1} E_n \cap G\right) = \lim_{n \geq 1} \mu_1(E_n \cap G) = \lim_{n \geq 1} \mu_2(E_n \cap G) = \mu_2\left(\bigcup_{n \geq 1} E_n \cap G\right)$$

where we've used countable additivity (as "continuity from below" in the form given in Exercise 4.3 (i)) for the first and third equalities. Thus $\bigcup_{n \geq 1} E_n \in \Lambda$. Also, if $E \subset F$ and $E, F \in \Lambda$ then

$$\mu_1(G \cap (F \setminus E)) = \mu_1(G \cap F) - \mu_1(G \cap E) = \mu_2(G \cap F) - \mu_2(G \cap E) = \mu_2(G \cap (F \setminus E)),$$

where we've used additivity of μ_1 and μ_2 , together with the fact that $\mu_1(G) < \infty$, for the first and third equalities. Thus $F \setminus E \in \Lambda$. It follows that Λ is a λ -system containing \mathcal{P} , so Λ contains $\mathcal{F} = \sigma(\mathcal{P})$ by Lemma 4.5. If $\mu_1(\Omega) < \infty$ then by taking $G = \Omega$ the result follows.

For the general case, let $(\Omega^n, n \geq 1)$ be elements of \mathcal{P} with $\Omega^n \uparrow \Omega$ and with $\mu_1(\Omega^n) < \infty$ for all n . Then for all $E \in \mathcal{F}$, since μ_1 and μ_2 are measures

$$\mu_1(E) = \lim_{n \rightarrow \infty} \mu_1(E \cap \Omega_n) = \lim_{n \rightarrow \infty} \mu_2(E \cap \Omega_n) = \mu_2(E),$$

where we have taken $G = \Omega_n$ to deduce that $\mu_1(E \cap \Omega_n) = \mu_2(E \cap \Omega_n)$. □

Remark. We say a measure μ on measurable space (Ω, \mathcal{F}) is σ -finite if there exists an increasing sequence $(\Omega_n, n \geq 1)$ of elements of \mathcal{F} with $\bigcup_{n \geq 1} \Omega_n = \Omega$ and with $\mu(\Omega_n) < \infty$ for all $n \geq 1$. Condition (b) in Dynkin's theorem is *stronger* than σ -finiteness, as it requires the approximating sets to in fact lie in \mathcal{P} .

One might think that the requirement (b) in the statement of Theorem 4.2 could be weakened to simply require σ -finiteness. The result of Exercise 4.9, below, shows that this is not the case. (We must briefly postpone stating the example, until we have defined Borel sets - but they are coming very shortly.)

4.3. Measures on \mathbb{R} . The above development meant to be in service of defining probability measures in particular. The most fundamental example driving the theory is that of measures on \mathbb{R} . We already discussed the specification of probability distributions on \mathbb{R} via their CDFs. Returning to this more formally and slightly more generally, we say $F : \mathbb{R} \rightarrow \mathbb{R}$ is a Stieltjes function if F is non-decreasing and right-continuous with left limits. If additionally $\lim_{x \rightarrow -\infty} F(x) = 0$ and $\lim_{x \rightarrow \infty} F(x) = 1$ then F is called a *cumulative distribution function*. Recall from above that $\mathcal{A}(\mathbb{R})$ is the set of finite unions of intervals of the form $\bigcup_{i=1}^n (a_i, b_i]$. We now define a function $\mu_F : \mathcal{A}(\mathbb{R}) \rightarrow [0, \infty]$ starting from the supposition that $\mu_F((a, b]) = F(b) - F(a)$. We are then forced by additivity to set $\mu_F(\bigcup_{i=1}^n (a_i, b_i]) = \sum_{i=1}^n F(b_i) - F(a_i)$ whenever $(a_i, b_i]$ are disjoint intervals.

Lemma 4.6. For any Stieltjes function F , μ_F is a pre-measure on $\mathcal{A}(\mathbb{R})$.

Proof. It is easy to verify that $\mathcal{A}(\mathbb{R})$ is a ring (this is Exercise 4.5). We show that μ_F is a pre-measure in three steps.

The first step is to check that μ_F is well-defined, i.e., that the expression in the definition of μ_F does not depend on how the elements of $\mathcal{A}(\mathbb{R})$ are expressed as finite disjoint unions. To see this, suppose that

$$L := \bigcup_{i=1}^n (a_i, b_i] = \bigcup_{j=1}^m (c_j, d_j]$$

are two ways of expressing the same element of \mathcal{A} as a disjoint union. For each $i \in [n]$ and $j \in [m]$, if $(a_i, b_i] \cap (c_j, d_j]$ is non-empty we denote the intersection by $(\ell_{ij}, r_{ij}]$. Then

$$\begin{aligned} \sum_{i=1}^n F(b_i) - F(a_i) &= \sum_{i=1}^n \sum_{\{j:(a_i,b_i] \cap (c_j,d_j] \neq \emptyset\}} F(r_{ij}) - F(\ell_{ij}) \\ &= \sum_{j=1}^m \sum_{\{i:(a_i,b_i] \cap (c_j,d_j] \neq \emptyset\}} F(r_{ij}) - F(\ell_{ij}) = \sum_{j=1}^m F(d_j) - F(c_j), \end{aligned}$$

so μ_F is indeed well-defined.

We next check that μ_F is additive. This is easy: if $\bigcup_{i=1}^n (a_i, b_i]$ and $\bigcup_{j=1}^m (c_j, d_j]$ are disjoint elements of $\mathcal{A}(\mathbb{R})$ then

$$\mu_F \left(\bigcup_{i=1}^n (a_i, b_i] \cup \bigcup_{j=1}^m (c_j, d_j] \right) = \sum_{i=1}^n (F(b_i) - F(a_i)) + \sum_{j=1}^m (F(d_j) - F(c_j)),$$

σ -finite

Stieltjes function

Cumulative distribution function

which is indeed the sum of the measures of the two elements of $\mathcal{A}(\mathbb{R})$.

Finally, we check that μ_F is a pre-measure. For this, suppose that

$$L := \bigcup_{i=1}^n (a_i, b_i] = \bigcup_{j=1}^{\infty} (c_j, d_j]$$

where the two unions are over disjoint intervals. Then for all $m \in \mathbb{N}$,

$$\bigcup_{i=1}^n (a_i, b_i] \supset \bigcup_{j=1}^m (c_j, d_j],$$

Thus, by monotonicity of μ_F ,

$$\mu_F\left(\bigcup_{i=1}^n (a_i, b_i]\right) \geq \sup_{m \geq 1} \mu_F\left(\bigcup_{j=1}^m (c_j, d_j]\right) = \sum_{j=1}^{\infty} \mu_F(c_j, d_j];$$

to complete the proof, we must show that in fact equality holds.

Suppose for a contradiction that $\mu_F(L) = \sum_{j=1}^{\infty} \mu_F(c_j, d_j] + 2\epsilon$, for some $\epsilon > 0$, and write $\Delta_m := L \setminus \bigcup_{i=1}^m (c_i, d_i]$. Note that $\Delta_m \in \mathcal{A}$ — it is a difference of finite unions of intervals — and $\Delta_m \downarrow \emptyset$ as $m \rightarrow \infty$. Also, since $L = \Delta_m \cup \bigcup_{i=1}^m (c_i, d_i]$ is a disjoint union, it follows that

$$\mu_F(\Delta_m) = \mu_F(L) - \sum_{i=1}^m \mu_F(c_i, d_i] \geq 2\epsilon$$

for all m .

Choose $D_m \in \mathcal{A}$ with $\overline{D_m} \subset \Delta_m$ and such that $\mu_F(\Delta_m \setminus D_m) \leq \epsilon/2^m$ for all m .⁵ Since

$$\Delta_m = \bigcap_{i=1}^m \Delta_i = \bigcap_{i=1}^m D_i \cup (\Delta_i \setminus D_i) \subseteq \bigcap_{i=1}^m D_i \cup \bigcup_{i=1}^m (\Delta_i \setminus D_i),$$

by monotonicity

$$2\epsilon \leq \mu_F(\Delta_m) \leq \mu_F\left(\bigcap_{i=1}^m D_i\right) + \sum_{i=1}^m \mu_F(\Delta_i \setminus D_i) \leq \mu_F\left(\bigcap_{i=1}^m D_i\right) + \epsilon.$$

Thus $\mu_F\left(\bigcap_{i=1}^m D_i\right) \geq \epsilon$ for all m , so $\bigcap_{i=1}^m \overline{D_i} \neq \emptyset$ for all m , so $\bigcap_{i=1}^{\infty} \Delta_i \supset \bigcap_{i=1}^{\infty} \overline{D_i} \neq \emptyset$, contradicting the fact that $\Delta_m \downarrow \emptyset$ as $m \rightarrow \infty$. \square

The σ -field generated by $\mathcal{A}(\mathbb{R})$ is called the *Borel σ -field*, and denoted $\mathcal{B}(\mathbb{R})$; its elements are called *Borel sets* of \mathbb{R} . The next exercise asks you to show that $\mathcal{B}(\mathbb{R})$ is the smallest σ -field containing all open sets in \mathbb{R} . $\mathcal{B}(\mathbb{R})$.

Exercise 4.7. Show that $\sigma(\mathcal{A}(\mathbb{R})) = \sigma(\{U \subset \mathbb{R} : U \text{ open}\})$.

More generally, given a topological space M , the Borel σ -field over M is defined to be the σ -field generated by the open sets, $\mathcal{B}(M) := \sigma(\{U \subset M : U \text{ open}\})$. $\mathcal{B}(M)$.

With Lemma 4.6 under our belt, it now follows easily that Stieltjes functions \mathbb{R} uniquely determine measures on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$.

Theorem 4.7. Let F be a Stieltjes function. Then there exists a unique measure μ on $\mathcal{B}(\mathbb{R})$ such that $\mu(a, b] = F(b) - F(a)$ for all $-\infty < a \leq b < \infty$.

Proof. Write $\mathcal{P} = \{(a, b] : -\infty < a \leq b < \infty\}$. By Lemma 4.6, there exists a pre-measure μ on $\mathcal{A}(\mathbb{R})$ such that $\mu(a, b] = F(b) - F(a)$ for all $(a, b] \in \mathcal{P}$. By the Carathéodory Extension Theorem, Theorem 4.1, μ extends to a measure $\mu^+ : \mathcal{F} \rightarrow [0, \infty]$ for some σ -field \mathcal{F} containing

⁵Not hard to see this is possible - add a proof?

$\mathcal{A}(\mathbb{R})$. Since $\mathcal{B}(\mathbb{R})$ is the smallest σ -field containing $\mathcal{A}(\mathbb{R})$, the restriction of μ^+ to $\mathcal{B}(\mathbb{R})$ is well-defined, is a measure on $\mathcal{B}(\mathbb{R})$ and has $\mu^+(a, b] = \mu(a, b] = F(b) - F(a)$ for all $(a, b] \in \mathcal{P}$. This proves existence.

Now suppose that μ_1 and μ_2 are measures on $\mathcal{B}(\mathbb{R})$ satisfying the hypotheses of the theorem. Then μ_1 and μ_2 agree on \mathcal{P} . But \mathcal{P} is a π -system. Clearly $\sigma(\mathcal{P})$ contains $\mathcal{A}(\mathbb{R})$, so so we must have $\sigma(\mathcal{P}) = \sigma(\mathcal{A}(\mathbb{R})) = \mathcal{B}(\mathbb{R})$. It follows by Dynkin's theorem, Theorem 4.2, that $\mu_1 \equiv \mu_2$. This proves uniqueness. \square

The above proof refers to “some σ -field \mathcal{F} containing $\mathcal{A}(\mathbb{R})$ ”. Looking back at the statement of of the Carathéodory lemma reveals that the σ -field \mathcal{F} consists precisely of the μ^* -additive sets.

The next exercise reveals more information about the collection of μ^* -additive sets, and its relation to the Borel σ -fields. The exercise after that provides an example which shows that condition (b) in Dynkin's theorem can not be replaced by σ -finiteness. The following definition features in first of the two exercises: we say a measure space $(\Omega, \mathcal{F}', \mu')$ extends another measure space $(\Omega, \mathcal{F}, \mu)$ if $\mathcal{F} \subseteq \mathcal{F}'$ and $\mu'|_{\mathcal{F}} \equiv \mu$.

Exercise 4.8. Let $(\Omega, \mathcal{F}, \mu)$ be a measure space. Say $N \in \mathcal{F}$ is a null set if $\mu(N) = 0$. Say that $(\Omega, \mathcal{F}, \mu)$ is complete if for any null set N , for all $M \subset N$ we have $M \in \mathcal{F}$.

(a) Write $\overline{\mathcal{F}} := \bigcap_{\{(\Omega, \mathcal{F}', \mu') \text{ extending } (\Omega, \mathcal{F}, \mu): (\Omega, \mathcal{F}', \mu') \text{ complete}\}} \mathcal{F}'$. Prove that $\overline{\mathcal{F}} = \{E \cup M : E \in \mathcal{F}, M \subset N \text{ for some null set } N \in \mathcal{F}\}$.

(b) Let $\mu^* : 2^\Omega \rightarrow [0, \infty]$ be an outer measure on some ground set Ω , and let $\mathcal{F} = \{A \subset \Omega : A \text{ is } \mu^*\text{-additive}\}$. Show that $(\Omega, \mathcal{F}, \mu^*|_{\mathcal{F}})$ is complete.

(c) Let μ be the Lebesgue pre-measure on $\mathcal{A}(\mathbb{R})$, i.e., with $\mu(a, b] = b - a$ for bounded intervals $(a, b] \subset \mathbb{R}$. Let μ^* be the corresponding outer measure on \mathbb{R} , and let $\mathcal{L}(\mathbb{R}) = \{S \subset \mathbb{R} : S \text{ is } \mu^*\text{-additive}\}$. Show that $\mathcal{L}(\mathbb{R})$ is the completion of $\mathcal{B}(\mathbb{R})$.

NB: For (c) you will need the σ -finiteness of Lebesgue measure.

The set $\mathcal{L}(\mathbb{R})$ is known as the Lebesgue σ -field (actually, I have only ever seen it called the Lebesgue σ -algebra, but I decided to call them σ -fields, and I'm sticking to it).

Exercise 4.9. Consider the measures μ_1, μ_2 on $\mathcal{B}(\mathbb{R})$ defined by $\mu_1(B) = |B \cap \mathbb{Q}|, \mu_2(B) = 2|B \cap \mathbb{Q}|$.

- (a) Show that μ_1 and μ_2 are σ -finite measures.
- (b) Show that $\mu_1(A) = \mu_2(A)$ for all $A \in \mathcal{A}(\mathbb{R})$.

Given a measurable space (Ω, \mathcal{F}) a set S and a function $f : \Omega \rightarrow S$, the push-forward of \mathcal{F} to S is the set $f^*(\mathcal{F}) = \{B \subset S : f^{-1}(B) \in \mathcal{F}\}$.

Exercise 4.10. Show that the push-forward $f^*(\mathcal{F})$ is a σ -field.

4.4. **Independent events.** A probability space is a measure space $(\Omega, \mathcal{F}, \mathbf{P})$ with $\mathbf{P}(\Omega) = 1$. Elements of \mathcal{F} are called events; elements of Ω are called elementary events.⁶

We say that events $(E_i, i \in I)$ are mutually independent if for all $J \subset I$ finite,

$$\mathbf{P} \left\{ \bigcap_{j \in J} E_j \right\} = \prod_{j \in J} \mathbf{P} \{E_j\}. \tag{4.1}$$

(Often the word “mutually” is omitted.) For $k \geq 1$, we say the events $(E_i, i \in I)$ are k -wise independent if (4.1) holds for all $J \subset I$ with $|J| \leq k$. In particular, they are pairwise independent if $\mathbf{P} \{E_i \cap E_j\} = \mathbf{P} \{E_i\} \mathbf{P} \{E_j\}$ for any distinct $i, j \in I$.

Exercise 4.11. In this exercise we say that an event E in a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ is non-trivial if $\mathbf{P} \{E\} \in (0, 1)$.

- (a) Let $k \in \mathbb{N}$, let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space and let (E_1, \dots, E_k) be nontrivial, independent events in \mathcal{F} . Prove that $|\Omega| \geq 2^k$.

⁶An unfortunate aspect of this terminology: elementary events need not be events! But it is what it is.

Probability space

Independent events

- (b) Construct a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ with $|\Omega| = 2^{k-1}$ and nontrivial events (E_1, \dots, E_k) , such that for any $1 \leq i \leq k$, the events $(E_j, j \in [k] \setminus \{i\})$ are mutually independent, but (E_1, \dots, E_k) are not mutually independent.

The following example is further developed in the homework (and inspired by the use of Rademacher random variables in James Norris's "Probability and measure" lecture notes). It is a hands-on way to model an infinite sequence of independent fair coin tosses. Let $\Omega = [0, 1]$, $\mathcal{F} = \mathcal{B}([0, 1]) = \mathcal{B}(\mathbb{R})|_{[0,1]}$, and let \mathbf{P} be Lebesgue measure on $[0, 1]$, which is often called the *uniform* probability measure in this context; then $(\Omega, \mathcal{F}, \mathbf{P})$ is a probability space. For $k \geq 1$ define the event

$$A_k = \bigcup_{\substack{0 \leq i < 2^k \\ i \text{ even}}} \left(\frac{i}{2^k}, \frac{i+1}{2^k} \right]. \quad (4.2)$$

So $A_1 = (0, 1/2]$, $A_2 = (0, 1/4] \cup (1/2, 3/4]$, and so on.

Exercise 4.12. Show that $(A_k, k \geq 1)$ are mutually independent.

Note that A_i may be thought of as the set of $x \in (0, 1]$ for which the i 'th bit in the binary expansion is zero (provided we adopt the convention that we never use infinite strings of zeros in our binary representation).

The *Borel–Cantelli lemmas* are basic and important workhorses of probability theory; stating them will additionally help us away from the language of sets and toward probabilistic terminology. Given a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and events $(E_n, n \geq 1)$, we define

$$\limsup_{n \rightarrow \infty} E_n := \bigcap_{n \geq 1} \bigcup_{m \geq n} E_m = \{\omega \in \Omega : \omega \in E_n \text{ for infinitely many } n\}.$$

(In fact, this definition makes sense for any sequence of sets $(E_n, n \geq 1)$ over a common ground set Ω .) Thinking probabilistically, if $\omega \in \limsup_{n \rightarrow \infty} E_n$ then infinitely many of the events E_n occur; we therefore introduce $\{E_n \text{ occurs infinitely often}\}$ or simply $\{E_n \text{ i.o.}\}$ as alternative notation for the set $\limsup_{n \rightarrow \infty} E_n$.

Similarly, we define

$$\liminf_{n \rightarrow \infty} E_n := \bigcup_{n \geq 1} \bigcap_{m \geq n} E_m = \{\omega \in \Omega : \omega \in E_n \text{ for all but finitely many } n\}.$$

Note that $(\limsup_{n \rightarrow \infty} E_n)^c = \liminf_{n \rightarrow \infty} (E_n^c)$.

As an example, for the events $(A_k, k \geq 1)$ described above, we have

$$\begin{aligned} \limsup_{n \rightarrow \infty} A_n &= \{x \in [0, 1] : \text{there are infinitely many zeros in any binary expansion of } x\}, \text{ and} \\ \liminf_{n \rightarrow \infty} A_n &= \{x \in [0, 1] : x = k/2^n, \text{ for some integers } n, k \text{ with } n \geq 1, 0 \leq k \leq 2^n\}. \end{aligned}$$

Lemma 4.8 (First Borel–Cantelli Lemma). Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space and let $(E_n, n \geq 1)$ be events in \mathcal{F} . If $\sum_{n \geq 1} \mathbf{P}\{E_n\} < \infty$ then $\mathbf{P}\{E_n \text{ i.o.}\} = 0$.

Proof. Fix $\epsilon > 0$. Then there exists n_0 such that $\sum_{m \geq n_0} \mathbf{P}\{E_m\} < \epsilon$, so by monotonicity and subadditivity of measures,

$$\mathbf{P}\{E_n \text{ i.o.}\} \leq \mathbf{P}\left\{ \bigcap_{n \geq n_0} \bigcup_{m \geq n} E_m \right\} \leq \mathbf{P}\left\{ \bigcup_{m \geq n_0} E_m \right\} \leq \sum_{m \geq n_0} \mathbf{P}\{E_m\} < \epsilon.$$

Since $\epsilon > 0$ was arbitrary, the result follows. \square

Lemma 4.9 (Second Borel–Cantelli Lemma). *Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space and let $(E_n, n \geq 1)$ be mutually independent events in \mathcal{F} . If $\sum_{n \geq 1} \mathbf{P}\{E_n\} = \infty$ then $\mathbf{P}\{E_n \text{ i.o.}\} = 1$.*

Proof. Note that by definition,

$$\{E_n \text{ i.o.}\}^c = \liminf_{n \rightarrow \infty} (E_n^c) = \bigcup_{n \geq 1} \bigcap_{m \geq n} E_m^c,$$

so by subadditivity

$$\mathbf{P}\{\{E_n \text{ i.o.}\}^c\} \leq \sum_{n \geq 1} \mathbf{P}\left\{\bigcap_{m \geq n} E_m^c\right\}.$$

To prove the lemma we'll show that the summands on the right are all zero.

Writing $p_n = \mathbf{P}\{E_n\}$, for all $1 \leq n \leq N$, by monotonicity and independence we have

$$\mathbf{P}\left\{\bigcap_{m \geq n} E_m^c\right\} \leq \mathbf{P}\left\{\bigcap_{m=n}^N E_m^c\right\} = \prod_{m=n}^N (1 - p_m).$$

Since this holds for all N , and since $1 - p_n \leq e^{-p_n}$, it follows that

$$\mathbf{P}\left\{\bigcap_{m \geq n} E_m^c\right\} \leq \prod_{m=n}^{\infty} (1 - p_m) \leq e^{-\sum_{m=n}^{\infty} p_m} = 0,$$

as required. □

The two Borel–Cantelli lemmas together show that if $(E_n, n \geq 1)$ is any sequence of independent events, then $\mathbf{P}\{E_n \text{ i.o.}\} \in \{0, 1\}$. This is a first instance of a *zero-one law*, and a special case of Kolmogorov's zero-one law, which you will meet quite shortly.

We conclude the section by defining independence of σ -fields and establishing a sufficient condition for such independence. Given a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and a collection $(\mathcal{G}_i, i \in I)$ of subsets of \mathcal{F} , we say that $(\mathcal{G}_i, i \in I)$ are independent if for all $J \subset I$ finite and any events $(E_j, j \in J)$ with each $E_j \in \mathcal{G}_j$, we have

$$\mathbf{P}\left\{\bigcap_{j \in J} E_j\right\} = \prod_{j \in J} \mathbf{P}\{E_j\}.$$

Proposition 4.10. *Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and let \mathcal{P}, \mathcal{Q} be π -systems in \mathcal{F} . If $\mathbf{P}\{A \cap B\} = \mathbf{P}\{A\}\mathbf{P}\{B\}$ for all $A \in \mathcal{P}, B \in \mathcal{Q}$, then $\sigma(\mathcal{P})$ and $\sigma(\mathcal{Q})$ are independent.*

Proof. First fix $A \in \mathcal{P}$ and define measures μ_A, \mathbf{P}_A on $\sigma(\mathcal{Q})$ by

$$\mu_A(B) = \mathbf{P}\{A\}\mathbf{P}\{B\} \text{ and } \mathbf{P}_A(B) = \mathbf{P}\{A \cap B\}.$$

Then $\mu_A(B) = \mathbf{P}_A(B)$ for all $B \in \mathcal{Q}$, and $\mu_A(\Omega) = \mathbf{P}\{A\} = \mathbf{P}_A(\Omega)$, so $\mu_A = \mathbf{P}_A$ by Dynkin's theorem. Thus

$$\mathbf{P}\{A \cap B\} = \mathbf{P}\{A\}\mathbf{P}\{B\}$$

for all $A \in \mathcal{P}, B \in \sigma(\mathcal{Q})$.

Next, fix $B \in \sigma(\mathcal{Q})$ and define measures ν^B, \mathbf{P}^B on $\sigma(\mathcal{P})$ by

$$\nu^B(A) = \mathbf{P}\{A\}\mathbf{P}\{B\} \text{ and } \mathbf{P}^B(A) = \mathbf{P}\{A \cap B\}.$$

Then $\nu^B(A) = \mathbf{P}^B(A)$ for $A \in \mathcal{P}$, and $\nu^B(\Omega) = \mathbf{P}\{B\} = \mathbf{P}^B(\Omega)$, so by Dynkin's theorem we have $\nu^B = \mathbf{P}^B$. Thus $\mathbf{P}\{A \cap B\} = \mathbf{P}\{A\}\mathbf{P}\{B\}$ for all $A \in \sigma(\mathcal{P})$ and $B \in \sigma(\mathcal{Q})$, i.e., $\sigma(\mathcal{P})$ and $\sigma(\mathcal{Q})$ are independent. □

Exercise 4.13. Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and π -systems $(\mathcal{P}_i, i \in I)$ which are subsets of \mathcal{F} . Then the σ -fields $(\sigma(\mathcal{P}_i), i \in I)$ are independent if and only if $\mathbf{P} \left\{ \bigcap_{j \in J} E_j \right\} = \prod_{j \in J} \mathbf{P} \{E_j\}$ for all $J \subset I$ finite and any events $E_j \in \mathcal{P}_j$.

5. Random variables

Much of the richness of probability theory arises from the interaction of independence with random variables, but to explore that, we need to define random variables first!

To begin, given measurable spaces (R, \mathcal{R}) and (S, \mathcal{S}) , a $(\mathcal{R}/\mathcal{S})$ -measurable map is a function $f : R \rightarrow S$ such that $f^{-1}(E) \in \mathcal{R}$ for all $E \in \mathcal{S}$.⁷ If R and S are topological spaces and \mathcal{R}, \mathcal{S} are the Borel σ -algebras, then f is also called a *Borel function*.

Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$. A (real) *random variable* is a $(\mathcal{F}/\mathcal{B}(\mathbb{R}))$ -measurable function $X : \Omega \rightarrow \mathbb{R}$. In other words, random variables are just measurable maps but where the domain happens to be the ground set of a probability space. Real random variables and extended real random variables are the bread and butter of the course. The laws of large numbers are the jam. Basic measure theory is the plate. Enough with that metaphor. (We write $\mathbb{R}^* = \mathbb{R} \cup \{\pm\infty\}$ for the extended real line; its open sets are generated by those of \mathbb{R} together with sets of the form $(x, \infty]$ and $[-\infty, x)$ for $x \in \mathbb{R}$; an extended real random variable is a $(\mathcal{F}/\mathcal{B}(\mathbb{R}^*))$ -measurable map $X : \Omega \rightarrow \mathbb{R}^*$.) Oh, and random variables taking values in more general spaces are the *croissants au beurre*. If M is a topological space and $X : \Omega \rightarrow M$ is $(\mathcal{F}/\mathcal{B}(M))$ -measurable then we call X an M -valued random variable; if (S, \mathcal{S}) is a measurable space and $X : \Omega \rightarrow S$ is $(\mathcal{F}/\mathcal{S})$ -measurable then we call X an S -valued random variable.

For a function $X : \Omega \rightarrow \mathbb{R}$ or $X : \Omega \rightarrow \mathbb{R}^*$, it's very useful to introduce the notation $\{X \leq r\} := \{\omega \in \Omega : X(\omega) \leq r\}$ and to think of this set as “the event that $X \leq r$ ”. More generally for a function $X : R \rightarrow S$ and $U \subset S$ we write $\{X \in U\} := X^{-1}(U)$.

Before diving into the theory, it's worth motivating ourselves (and honing our intuition) by considering an example. We revisit the events A_k defined in (4.2), above, and define $R_k : [0, 1] \rightarrow \mathbb{R}$ by $R_k = \mathbf{1}_{[A_k]}$, so $R_k(x) = 1$ if and only if $x \in A_k$.

Exercise 5.1. Show that R_k is $\mathcal{B}([0, 1])/\mathcal{B}(\mathbb{R})$ -measurable.

Under the uniform probability measure on $[0, 1]$, we have

$$\mathbf{P} \{R_k = 1\} := \mathbf{P} \{\{x : R_k(x) = 1\}\} = \mathbf{P} \{A_k\} = \frac{1}{2}.$$

This agrees with the intuition that for a uniformly random point in $[0, 1]$, each bit of the binary expansion is equally likely to be zero or one. Moreover, intuition suggests that these bits should be independent, and that the asymptotic proportion of ones in the sequence $(R_n, n \geq 1)$ should be $1/2$. More precisely, we expect that

$$\mathbf{P} \left\{ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n R_n = \frac{1}{2} \right\} = 1. \quad (5.1)$$

To make rigorous sense of this assertion, we first need to know that

$$\left\{ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n R_n = \frac{1}{2} \right\} = \left\{ x \in [0, 1] : \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n R_n(x) = \frac{1}{2} \right\} \quad (5.2)$$

is a measurable set (otherwise its probability is not defined). Fortunately, this is not hard to see; the closure properties of σ -fields allow us to perform essentially any operations we please with random variables and obtain other random variables, provided we perform at most countably many operations in total. The next theorem provides a useful time-saving device for proving measurability of random variables; the subsequent exercise shows that many of the basic operations of arithmetic and analysis preserve measurability, and in particular implies that the set in (5.2) is measurable.

⁷Notice the similarity to the definition of continuous functions between topological spaces.

\mathbb{R}^*

Extended real random variable

Theorem 5.1. Let (R, \mathcal{R}) and (S, \mathcal{S}) be measurable spaces and let $f : R \rightarrow S$. Suppose that there is $\mathcal{A} \subset \mathcal{S}$ with $\sigma(\mathcal{A}) = \mathcal{S}$ such that $f^{-1}(A) \in \mathcal{R}$ for all $A \in \mathcal{A}$. Then f is $(\mathcal{R}/\mathcal{S})$ -measurable.

Proof. Let $\mathcal{S}_0 = \{E \in \mathcal{S} : f^{-1}(E) \in \mathcal{R}\}$. Then $\mathcal{A} \subset \mathcal{S}_0$ by assumption. Also, if $E \in \mathcal{S}_0$ then

$$f^{-1}(E^c) = \{r \in R : f(r) \in E^c\} = R \setminus \{r \in R : f(r) \in E\} = (f^{-1}(E))^c \in \mathcal{R},$$

so $E^c \in \mathcal{S}_0$. Similarly, if $(E_n, n \geq 1)$ are in \mathcal{S}_0 and $E_n \uparrow E_\infty$ then

$$f^{-1}(E_\infty) = \{r \in R : f(r) \in E_\infty\} = \bigcup_{n \geq 1} \{r \in R : f(r) \in E_n\} = \bigcup_{n \geq 1} f^{-1}(E_n) \in \mathcal{R},$$

so $E_\infty \in \mathcal{S}_0$. Thus \mathcal{S}_0 is a σ -field, so equals \mathcal{S} . □

Here are some examples of how the theorem is useful. Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$.

- If $X : \Omega \rightarrow \mathbb{R}$ satisfies that $\{X \leq r\} = X^{-1}(-\infty, r] \in \mathcal{F}$ for all $r \in \mathbb{R}$, then X is a real random variable (it is $(\mathcal{F}/\mathcal{B}(\mathbb{R}))$ -measurable).
- If $X : \Omega \rightarrow \mathbb{R}$ is a real random variable and $f : \mathbb{R} \rightarrow \mathbb{R}$ is continuous then $f(X)$ is another random variable. (Since if $U \subset \mathbb{R}$ is open, then $f^{-1}(U)$ is open, so $\{f(X) \in U\} = (f \circ X)^{-1}(U) = X^{-1}(f^{-1}(U)) \in \mathcal{B}(\mathbb{R})$; and the open sets are a π -system generating $\mathcal{B}(\mathbb{R})$.)
- If $X_J = (X_j, j \in J)$ is a finite collection of random variables then X_J may be viewed as a function from Ω to \mathbb{R}^J , sending ω to $(X_j(\omega), j \in J)$. The collection

$$\mathcal{P}_J := \left\{ \prod_{j \in J} (-\infty, b_j] : b_j \in \mathbb{R} \text{ for } j \in J \right\} \tag{5.3}$$

is a π -system generating the Borel sets $\mathcal{B}(\mathbb{R}^J)$. For any element $\prod_{j \in J} (-\infty, b_j]$ of \mathcal{P}_J , we have

$$X_J^{-1} \left(\prod_{j \in J} (-\infty, b_j] \right) = \bigcap_{j \in J} X_j^{-1}(b_j) \in \mathcal{F},$$

since \mathcal{F} is closed under finite intersections. Thus X is an \mathbb{R}^J -valued random variable.

- If R and S are topological spaces, and $h : R \rightarrow S$ is such that $h^{-1}(U) \in \mathcal{B}(R)$ for all open $U \subset S$, then h is a Borel function.

The next exercise asks you to check various closure properties of the collection of real-valued measurable maps. Some of these require enlarging the target space from the real numbers to the *extended* real numbers $\mathbb{R}^* := \mathbb{R} \cup \{-\infty, \infty\}$. The open sets of \mathbb{R}^* are generated by $\mathcal{A} = \{(a, b), a, b \in \mathbb{R}\} \cup \{(a, \infty], a \in \mathbb{R}\} \cup \{[-\infty, b), b \in \mathbb{R}\}$, so \mathcal{A} also generates the Borel sets of \mathbb{R}^* : that is, $\mathcal{B}(\mathbb{R}^*) = \sigma(\mathcal{A})$.

Exercise 5.2. Let (Ω, \mathcal{F}) be a measurable space and let X, Y , and $(X_n, n \geq 1)$ be $(\mathcal{F}/\mathcal{B}(\mathbb{R}))$ -measurable maps from Ω to \mathbb{R} .

- (a) Prove that $\mathbf{1}_{[X \geq 0]}, X + Y, XY, (X/Y)\mathbf{1}_{[Y \neq 0]}$ are all $(\mathcal{F}/\mathcal{B}(\mathbb{R}))$ -measurable.
- (b) Prove that $\sup_{n \geq 1} X_n, \inf_{n \geq 1} X_n, \limsup_{n \geq 1} X_n$ and $\liminf_{n \geq 1} X_n$ are all $(\mathcal{F}/\mathcal{B}(\mathbb{R}^*))$ -measurable.
- (c) Prove that if Z is any of the four expressions from part (b), then $Z\mathbf{1}_{[Z \in \mathbb{R}]}$ is $(\mathcal{F}/\mathcal{B}(\mathbb{R}))$ -measurable.
- (d) Prove that if $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is $(\mathcal{B}(\mathbb{R}^n)/\mathcal{B}(\mathbb{R}))$ -measurable then $f(X_1, \dots, X_n)$ is $(\mathcal{F}/\mathcal{B}(\mathbb{R}))$ -measurable.

Given a sequence $(a_n, n \geq 1)$ of real numbers, we say that $\lim_{n \rightarrow \infty} a_n$ exists if either there is $a \in \mathbb{R}$ such that $\lim_{n \rightarrow \infty} a_n = a$, or if $\lim_{n \rightarrow \infty} a_n = \infty$ or $\lim_{n \rightarrow \infty} a_n = -\infty$.

Proposition 5.2. If $(X_n, n \geq 1)$ is a sequence of random variables on probability space $(\Omega, \mathcal{F}, \mathbf{P})$, then

$$E := \left\{ \lim_{n \rightarrow \infty} X_n \text{ exists} \right\} = \left\{ \omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) \text{ exists} \right\}$$

is an element of \mathcal{F} .

Remove \mathbf{P} since it is not needed in the proposition? But trying to encourage readers to think probabilistically...

Proof. By Exercise 5.2 (b), $\overline{X} := \limsup_{n \geq 1} X_n$ and $\underline{X} := \liminf_{n \geq 1} X_n$ are extended real-valued random variables, so

$$E_\infty := \left\{ \lim_{n \rightarrow \infty} X_n = \infty \right\} = \{ \underline{X} = \infty \}$$

is an event, and

$$E_{-\infty} := \left\{ \lim_{n \rightarrow \infty} X_n = -\infty \right\} = \{ \overline{X} = -\infty \}$$

is an event. Also,

$$E_{\text{bd}} := \{(X_n, n \geq 1) \text{ is a bounded sequence}\} = \{-\infty < \underline{X}\} \cap \{\overline{X} < \infty\}$$

is an event, so

$$E_{\text{fin}} = \left\{ \lim_{n \rightarrow \infty} X_n \text{ exists and is finite} \right\} = E_{\text{bd}} \cap \bigcap_{m \in \mathbb{N}} \{ \overline{X} - \underline{X} < 1/m \}$$

is an event. Since $E = E_\infty \cup E_{-\infty} \cup E_{\text{fin}}$, this completes the proof. \square

Exercise 5.3. Let (Ω, \mathcal{F}) be a measurable space and let $(X_n, n \geq 1)$ be $(\mathcal{F}/\mathcal{B}(\mathbb{R}))$ -measurable maps from Ω to \mathbb{R} . Write $E := \{\omega : \lim_{n \rightarrow \infty} X_n(\omega) \text{ exists}\}$. Prove that, defining $X_\infty : \Omega \rightarrow \mathbb{R}^*$ by

$$X_\infty(\omega) = \begin{cases} \lim_{n \rightarrow \infty} X_n(\omega) & \text{if } \omega \in E \\ 0 & \text{otherwise,} \end{cases}$$

then X_∞ is $(\mathcal{F}/\mathcal{B}(\mathbb{R}^*))$ -measurable.

5.1. Generated σ -fields. Fix a set R and a measurable space (S, \mathcal{S}) . Given a collection $(X_i, i \in I)$ of functions from R to S , we define

$$\sigma(X_i, i \in I) := \sigma(\{X_i^{-1}(E) : i \in I, E \in \mathcal{S}\}) = \bigcap_{\substack{\mathcal{F} \text{ a } \sigma\text{-field over } R \\ \forall i \in I, X_i \text{ is } (\mathcal{F}/\mathcal{S})\text{-measurable}}} \mathcal{F}.$$

In words, $\sigma(X_i, i \in I)$ is the smallest σ -field over R to yield measurability of all the maps $(X_i, i \in I)$. If (R, \mathcal{R}) is a measurable space and the functions $(X_i, i \in I)$ are all $(\mathcal{R}/\mathcal{S})$ -measurable, then $\sigma(X_i, i \in I) \subset \mathcal{R}$ by definition.

The most important example is that of a collection of real random variables $(X_i, i \in I)$ over a common probability space. For $i \in I$ we have $\sigma(X_i) = \{\{X_i \in B\}, B \in \mathcal{B}(\mathbb{R})\} = \sigma(\{X_i \leq b\}, b \in \mathbb{R})$, so it follows that

$$\sigma(X_i, i \in I) = \sigma\left(\bigcup_{i \in I} \{X_i \leq b\} : b \in \mathbb{R}\right)$$

For any $J \subset I$ finite and $(b_j, j \in J) \in \mathbb{R}^J$, it follows that

$$\{X_j \leq b_j, j \in J\} = \bigcap_{j \in J} \{X_j \leq b_j\} \in \sigma(X_i, i \in I),$$

so we may also write

$$\sigma(X_i, i \in I) = \sigma(\{X_j \leq b_j, j \in J\} : J \subset I \text{ finite}, (b_j, j \in J) \in \mathbb{R}^J).$$

Exercise 5.4. Let $(X_i, i \in I)$ be random variables defined on a common probability space. Show that

$$\sigma(X_i, i \in I) = \bigcup_{J \subset I, J \text{ countable}} \sigma(X_j, j \in J).$$

Exercise 5.5 (Doob-Dynkin Lemma). Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and let $X, Y : \Omega \rightarrow \mathbb{R}$ be random variables. Show that Y is $(\sigma(X)/\mathcal{B}(\mathbb{R}))$ -measurable if and only if there exists a Borel function $f : \mathbb{R} \rightarrow \mathbb{R}$ such that $Y = f(X)$.

5.2. Independence of random variables. We say a collection of random variables $(X_i, i \in I)$ over a common probability space $(\Omega, \mathcal{F}, \mathbf{P})$ are *mutually independent* if the σ -fields $(\sigma(X_i), i \in I)$ are mutually independent. (We'll often drop the word "mutually".) In other words, $(X_i, i \in I)$ are independent if for any $J \subset I$ finite, and any Borel sets $(B_j, j \in J)$, we have

$$\mathbf{P} \{X_j \in B_j, j \in J\} = \prod_{j \in J} \mathbf{P} \{X_j \in B_j\}.$$

Proposition 5.3. *Real random variables $(X_i, i \in I)$ defined on a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ are mutually independent if and only if for all $J \subset I$ finite, for any real numbers $(b_j, j \in J) \in \mathbb{R}^J$,*

$$\mathbf{P} \{X_j \leq b_j \text{ for all } j \in J\} = \prod_{j \in J} \mathbf{P} \{X_j \leq b_j\}.$$

Proof. By definition, $(X_i, i \in I)$ are mutually independent if and only if the σ -fields $(\sigma(X_i), i \in I)$ are independent. For $i \in I$, set $\mathcal{P}_i = \{\{X_i \leq r\}, r \in \mathbb{R}\}$. Then \mathcal{P}_i is a π -system with $\sigma(\mathcal{P}_i) = \sigma(X_i)$, so by Exercise 4.13, the σ -fields in $(\sigma(X_i), i \in I)$ are independent if and only if for all $J \subset I$ finite, and any choice of events $E_j \in \mathcal{P}_j$ for $j \in J$, it holds that $\mathbf{P} \left\{ \bigcap_{j \in J} E_j \right\} = \prod_{j \in J} \mathbf{P} \{E_j\}$. This is equivalent to the condition in the proposition. \square

Note that this proposition implies that the Rademacher random variables $(R_k, k \geq 1)$ defined earlier are independent, since for any $n \in \mathbb{N}$ and $b_1, \dots, b_n \in \mathbb{R}$,

$$\mathbf{P} \{R_k \leq b_k \text{ for all } k \in [n]\} = \left(\frac{1}{2}\right)^{\#\{k \in [n]: b_k \in [0,1]\}} = \prod_{k \in [n]} \mathbf{P} \{R_k \leq b_k\}.$$

5.3. Existence of random variables with given distributions. You already met cumulative distribution functions of random variables in passing in Section 4.2. Given a real random variable X on a probability space $(\Omega, \mathcal{F}, \mathbf{P})$, its cumulative distribution function $F_X : \mathbb{R} \rightarrow [0, 1]$ is given by $F_X(r) = \mathbf{P} \{X \leq r\}$. Its *distribution* is the measure μ_X on $\mathcal{B}(\mathbb{R})$ given by $\mu_X(B) = \mathbf{P} \{X \in B\}$ for $B \in \mathcal{B}(\mathbb{R})$. In other words, μ_X is the push-forward of the measure \mathbf{P} by X .

It's easy to see that F_X is a Stieltjes function, and that the Borel measure corresponding to F_X — which by Theorem 4.7 is unique — is μ_X . The next proposition says that, in turn, any cumulative distribution function (CDF) is the CDF of some random variable.

Proposition 5.4. *Let F be any CDF. Then there exists a random variable $X : [0, 1] \rightarrow \mathbb{R}$ on the probability space $([0, 1], \mathcal{B}([0, 1]), \text{Leb}_{[0,1]})$ such that $F_X = F$.*

Proof. It's both efficient and pedagogically useful to first treat a special case. Suppose F is the Uniform $[0, 1]$ CDF; that is,

$$F(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } 0 \leq x \leq 1 \\ 1 & \text{if } x \geq 1. \end{cases}$$

We claim that $U := \sum_{k \geq 1} 2^{-k} R_k$ has $F_U = F$. First, note that $U = \sup_{\ell \geq 1} \sum_{k=1}^{\ell} 2^{-k} R_k$. Each of the terms in the supremum is a finite sum of random variables, so is a random variable; thus U is a random variable by Exercise 5.2.

To see that $F_U = F$, note that for any $n \geq 1$ and $0 \leq m < 2^n$, if we write $m/2^n$ in binary as $m/2^n = 0.b_1 b_2 \dots b_n$ then

$$\mathbf{P} \left\{ U \in \left(\frac{m}{2^n}, \frac{m+1}{2^n} \right] \right\} = \mathbf{P} \{R_1 = b_1, \dots, R_n = b_n\} = \frac{1}{2^n}.$$

It follows that $\mathbf{P} \{U \leq m/2^n\} = m/2^n$.

μ_X
Note this is a different use of the term "push-forward" from earlier.

Writing $D = \bigcup_{n \geq 1} \{m/2^n, 0 \leq m \leq 2^n\}$ for the dyadic fractions in $[0, 1]$, for any $x \in (0, 1]$ we may thus find an increasing sequence $(x_k, k \geq 1)$ of elements of D with $x_k \rightarrow x$ as $k \rightarrow \infty$. For monotone sequences of events we may interchange limit and probability, so

$$\mathbf{P}\{U < x\} = \lim_{k \rightarrow \infty} \mathbf{P}\{U \leq x_k\} = \lim_{k \rightarrow \infty} x_k = x.$$

We also have $\mathbf{P}\{U \leq x\} \leq \inf\{\mathbf{P}\{U \leq y\} : y \in D, y \geq x\} = x$, so in fact we must have $\mathbf{P}\{U \leq x\} = x$. Thus F is indeed the CDF of U .

For the general case, fix any CDF $F : \mathbb{R} \rightarrow [0, 1]$, and let $G : [0, 1] \rightarrow \mathbb{R}^*$ be defined by

$$G(p) := \inf\{x : F(x) \geq p\}.$$

The function G is sometimes called the “right inverse” of F . It is straightforward to check that G is Borel measurable.

Note that for $q \in [0, 1]$ and $r \in \mathbb{R}$, if $F(r) \geq q$ then $\{x \in \mathbb{R} : F(x) \geq q\} \subset [r, \infty)$, so

$$G(q) = \inf\{x : F(x) \geq q\} \geq \inf[r, \infty) = r.$$

Conversely, if $F(r) < q$ then, by right-continuity of F , there exists $s > r$ such that $F(s) < q$. For such s we have $\{x \in \mathbb{R} : F(x) \geq q\} \subset (s, \infty)$, so $G(q) \geq s > r$.

The preceding paragraph establishes that $F(r) \geq q$ if and only if $r \geq G(q)$. Now let $X = G(U)$. Then X is a random variable since G is Borel, and for $r \in \mathbb{R}$,

$$\mathbf{P}\{X \leq r\} = \mathbf{P}\{G(U) \leq r\} = \mathbf{P}\{U \leq F(r)\} = F(r).$$

□

There is a simpler way to construct a Uniform $[0, 1]$ random variable on the probability space $([0, 1], \mathcal{B}([0, 1]), \text{Leb}_{[0,1]})$. Simply let $X : [0, 1] \rightarrow \mathbb{R}$ be the identity function, $X(\omega) = \omega$. Then for $x \in \mathbb{R}$,

$$\mathbf{P}\{X \leq x\} = \mathbf{P}\{\omega \in [0, 1] : \omega \leq x\} = \text{Leb}_{[0,1]}\{\omega \in [0, 1] : \omega \leq x\} = \begin{cases} 0 & x \leq 0 \\ x & x \in (0, 1] \\ 1 & x > 1, \end{cases}$$

so X is Uniform $[0, 1]$. Moreover, the function U defined in the course of the proof is essentially just the identity function ([expand on this](#)), which may make the proof seem unnecessarily complicated. However, by building a Uniform $[0, 1]$ random variable in this way, the argument can be more easily bootstrapped to yield not just a single random variable, but sequences of independent random variables with arbitrary prescribed CDFs.

Theorem 5.5. *Fix any sequence $(F_n, n \geq 1)$ of cumulative distribution functions. Then there exists a sequence of independent random variables $(X_n, n \geq 1)$ such that X_n has CDF F_n .*

Proof. In the previous proof, we constructed a random variable with a given CDF by an appropriate transformation of a uniform random variable. We want to do the same thing but using an independent uniform for each term in the sequence. For this, we begin by splitting the sequence $(R_n, n \geq 1)$ of Rademacher random variables into infinitely many independent groups; there is no canonical way to do this so we just pick one.

List the prime numbers as $(p_i, i \geq 1)$, so $p_1 = 2, p_2 = 3$ and so forth. Then for $j, k \geq 1$ set $Q_{j,k} = R_{p_j^k}$. Then for any $i, j \geq 1$ with $i \neq j$, the sequences $(Q_{i,j}, k \geq 1)$ and $(Q_{j,i}, k \geq 1)$ contain no common terms.

Next, for $i \geq 1$ let $U_i = \sum_{k \geq 1} 2^{-k} Q_{i,k}$. The random variables $(U_i, i \geq 1)$ are each Uniform $[0, 1]$ by the same reasoning as in the proof of Proposition 5.4. Moreover, they are independent since $\sigma(U_i) \subseteq \sigma(Q_{i,k}, k \geq 1)$, and the σ -fields $(\sigma(Q_{i,k}, k \geq 1), i \geq 1)$ are independent.

Now, for $n \geq 1$ let $G_n : [0, 1] \rightarrow \mathbb{R}$ be defined by $G_n(p) = \inf\{x : F_n(x) \geq p\}$, and set $X_n = G_n(U_n)$. Then X_n has CDF F_n by the argument from the proof of Proposition 5.4, and $(X_n, n \geq 1) = (G_n(U_n), n \geq 1)$ are independent since $(U_n, n \geq 1)$ are independent. □

The independence of the random variables $(X_n, n \geq 1)$ constructed in the above proof is a special case of the result of the following exercise.

Exercise 5.6. *If $(Y_i, i \in I)$ are mutually independent random variables, $(I_n, n \geq 1)$ partitions I , and for each n , $g_n : \mathbb{R}^{I_n} \rightarrow \mathbb{R}$ is $(\sigma(Y_i, i \in I_n)/\mathcal{B}(\mathbb{R}))$ -measurable, then with $X_n = g_n(Y_i, i \in I_n)$, the random variables $(X_n, n \geq 1)$ are independent.*

5.4. Kolmogorov’s zero-one law. Fix a countable collection $X = (X_n, n \in \mathbb{N})$ of random variables over a probability space $(\Omega, \mathcal{F}, \mathbf{P})$. For $M \subset \mathbb{N}$, write $\mathcal{T}_M = \mathcal{T}_M(X) := \sigma(X_m, m \in \mathbb{N} \setminus M)$. The tail σ -field is $\mathcal{T}(X) := \bigcap_{M \subset \mathbb{N}: |M| < \infty} \mathcal{T}_M$. Informally, it contains all information about the sequence $(X_n, n \geq 1)$ that can be obtained while ignoring any given finite set of the random variables. The term “tail” comes from the (standard) setting when $\mathbb{N} = \mathbb{N}$, in which case $\mathcal{T}(X) = \bigcap_{n \geq 1} \sigma(X_m, m > n)$, and from thinking of \mathbb{N} as arranged on the number line.

At first blush, it might seem that if the entries of $(X_n, n \in \mathbb{N})$ are independent then \mathcal{T} ought to be the trivial σ -field $\{\emptyset, \Omega\}$; after all, for any fixed $n \in \mathbb{N}$, it appears not to contain any information about X_n . However, that’s not quite the case. For example, the event that $\lim_{n \rightarrow \infty} X_n$ exists is a tail event, as is any event of the form

$$\{X_n \in B_n \text{ infinitely often}\} = \bigcap_{n \geq 1} \bigcup_{m \geq n} \{X_m \in B_m\},$$

where $(B_n, n \geq 1)$ are Borel sets in \mathbb{R} .

Exercise 5.7. *Prove carefully that the two preceding examples are indeed examples of tail events.*

Kolmogorov’s zero-one law says that \mathcal{T} is at least trivial in a somewhat weaker sense.

Theorem 5.6 (Kolmogorov’s 0-1 law). *Let $X = (X_n, n \in \mathbb{N})$ be a countable collection of independent random variables on a common probability space $(\Omega, \mathcal{F}, \mathbf{P})$. Then $\mathbf{P}\{E\} \in \{0, 1\}$ for all $E \in \mathcal{T}(X)$.*

Proof. Fix $E \in \mathcal{T}$. For any $n \in \mathbb{N}$ and $F \in \sigma(X_n)$, since $\mathcal{T} \subset \sigma(X_m, m \in \mathbb{N} \setminus \{n\})$, the events E and F are independent. Let

$$\mathcal{G} = \sigma(X_n, n \in \mathbb{N}) = \sigma\left(\bigcup_{n \in \mathbb{N}} \sigma(X_n)\right).$$

Note that E is independent of all events in $\bigcup_{n \in \mathbb{N}} \sigma(X_n)$ so is independent of \mathcal{G} ; that is, for all $F \in \mathcal{G}$, $\mathbf{P}\{E \cap F\} = \mathbf{P}\{E\}\mathbf{P}\{F\}$. However, $\mathcal{T} \subset \mathcal{G}$ so $E \in \mathcal{G}$, so

$$\mathbf{P}\{E \cap E\} = \mathbf{P}\{E\}^2;$$

this is only possible if $\mathbf{P}\{E\} \in \{0, 1\}$. □

One appealing thing about this result is that there are applications that can be described without having developed the theory of integration (expectation). We sketch two.

First, let $(X_n, n \geq 1)$ be independent random variables, and let \mathcal{T} be their tail σ -field. Write $S_n = \sum_{i=1}^n X_i$ and $M^+ := \limsup_{n \rightarrow \infty} S_n/n$, $M^- := \liminf_{n \rightarrow \infty} S_n/n$.

For all $x \in \mathbb{R}$, we have $\{M^+ \geq x\} \in \mathcal{T}$, so $\mathbf{P}\{M^+ \geq x\} \in \{0, 1\}$. Letting $x^+ = \sup\{x : \mathbf{P}\{M^+ \geq x\} = 1\}$, then for $y > x^+$ we have $\mathbf{P}\{M^+ \geq y\} < 1$ so $\mathbf{P}\{M^+ \geq 0\} = 0$. Thus $\mathbf{P}\{M^+ = x^+\} = 1$, and likewise $\mathbf{P}\{M^- = x^-\} = 1$. Moreover, $\mathbf{P}\{\lim_{n \rightarrow \infty} S_n/n \text{ exists}\} \in \{0, 1\}$. The strong law of large numbers gives a necessary and sufficient condition for the last probability to equal 1, provided the entries of $(X_n, n \geq 1)$ are identically distributed.

The second example is that of *percolation*, one of the most active areas of modern probability theory. Let \mathbb{Z}^d be the d -dimensional integer lattice; in this context the elements of \mathbb{Z}^d are called *sites*. Fix $p \in [0, 1]$ and let $B = (B_v, v \in \mathbb{Z}^d)$ be independent Bernoulli(p) random variables on a common probability space $(\Omega, \mathcal{F}, \mathbf{P})$. (By Bernoulli(p) we mean that $\mathbf{P}\{B_v = 1\} = p = 1 - \mathbf{P}\{B_v = 0\}$ for all $v \in \mathbb{Z}^d$.) Let \mathcal{T} be the tail σ -field of B .

We use B to define *site percolation clusters* as follows. Write $\mathbb{Z}^d(B) = \{v \in \mathbb{Z}^d : B_v = 1\}$. For $x, y \in \mathbb{Z}^d$ say that x is *connected to* y in $\mathbb{Z}^d(B)$, and write $x \xrightarrow{B} y$, if there is a nearest-neighbour path from x to y containing only elements of $\mathbb{Z}^d(B)$. Then for $x \in \mathbb{Z}^d$ define

$$\mathcal{C}(x) := \{y \in \mathbb{Z}^d : x \xrightarrow{B} y\}.$$

Note that if $y \in \mathcal{C}(x)$ then $\mathcal{C}(x) = \mathcal{C}(y)$.

Now let

$$E = \{\exists x \in \mathbb{Z}^d; |\mathcal{C}(x)| = \infty\} = \{\mathbb{Z}^d(B) \text{ contains an infinite connected component}\}.$$

An infinite connected component can not be created or destroyed by adding or removing finitely many sites, so E is a tail event; therefore $x(p, d) := \mathbf{P}\{E\} \in \{0, 1\}$ by Kolmogorov's 0-1 law. Which of these values is correct depends on the parameter p of the Bernoulli random variables and on the dimension d .

The *critical probability* for site percolation on \mathbb{Z}^d is

$$p_c(\mathbb{Z}^d) := \sup\{p : x(p, d) = 0\}.$$

We necessarily have $x(p, d) = 1$ for all $p > p_c$, but unlike for $\limsup S_n/n$ the first example, this doesn't imply that $x(p_c, d) = 1$. In fact, it is conjectured that $x(p_c, d) = 0$, or in words that there is "no percolation at criticality", in any dimension. This is probably the most famous open question in probability.

The next exercise should take care of any measurability concerns in the definition of percolation. Recall that $2^{\mathbb{Z}^d}$ is the set of all subsets of \mathbb{Z}^d . So a set $S \subset 2^{\mathbb{Z}^d}$ is a set of subsets of \mathbb{Z}^d ; we say such S is a *cylinder set* if

$$S = \{V \subset \mathbb{Z}^d : A \subset V, B \subset V^c\},$$

for some finite sets A, B .

Exercise 5.8. Let $\mathcal{G} = \sigma(B_v, v \in \mathbb{Z}^d)$ and let $\mathcal{G}^* = B^*(\mathcal{G})$ be the push-forward of \mathcal{G} under the map

$$\omega \mapsto \{B_v(\omega), v \in \mathbb{Z}^d\} \in \{0, 1\}^{\mathbb{Z}^d}.$$

Show that $\mathcal{G}^* = \sigma(\text{Cylinder sets in } 2^{\mathbb{Z}^d})$.

Exercise 5.9. Show carefully that the event $\{\exists x \in \mathbb{Z}^d; |\mathcal{C}(x)| = \infty\}$ is in $\mathcal{T} \subset \mathcal{G}$.

5.5. Almost sure convergence, convergence in probability and convergence in distribution. Let $(X_n, 1 \leq n \leq \infty)$ be a sequence of random variables defined on a common space $(\Omega, \mathcal{F}, \mathbf{P})$. We say X_n *converges almost surely* to X_∞ , and write $X_n \xrightarrow{\text{a.s.}} X_\infty$, if

$$\mathbf{P}\left\{\lim_{n \rightarrow \infty} X_n = X_\infty\right\} = 1.$$

We say X_n *converges in probability* to X_∞ if for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbf{P}\{|X_n - X_\infty| > \epsilon\} = 0.$$

Next, given random variables $(X_n, 1 \leq n \leq \infty)$, with $X_n : \Omega_n \rightarrow \mathbb{R}$ for some probability space $(\Omega_n, \mathcal{F}_n, \mathbf{P}_n)$, we say X_n *converges in distribution* to X_∞ , and write $X_n \xrightarrow{d} X_\infty$, if

$$\lim_{n \rightarrow \infty} \mathbf{P}_n\{X_n \leq x\} = \mathbf{P}_\infty\{X_\infty \leq x\}$$

for all x with $\mathbf{P}_\infty\{X_\infty = x\} = 0$. This may seem complicated compared with the previous definitions; the reason for this is that convergence in distribution is really a property of the *distributions* of the random variables (or, equivalently, of their CDFs), and is insensitive to the specific spaces on which they are defined.

Exercise 5.10. (a) Check that $X_n \xrightarrow{d} X_\infty$ iff $F_{X_n}(x) \rightarrow F_{X_\infty}(x)$ for all continuity points x of F_{X_∞} .
 (b) Show that if $X_n \xrightarrow{d} X$ as $X \rightarrow \infty$ and $X_n \xrightarrow{d} Y$ as $n \rightarrow \infty$ then $F_X = F_Y$ and so $\mu_X = \mu_Y$.

I haven't checked these exercises carefully, proceed at your own risk

Almost sure convergence

Convergence in probability

Convergence in distribution

One warning, which partially explains the restriction to $x \in \mathbb{R}$ with $\mathbf{P}\{X_\infty = x\} = 0$ above, is in order. Write $U_n = \sum_{k=1}^n 2^{-k} R_k$, where $(R_n, n \geq 1)$ are the Rademacher random variables defined earlier, and let $U_\infty = \sum_{k \geq 1} 2^{-k} U_k$. Then $U_n \rightarrow U_\infty$ almost surely, since $|U_n - U_\infty| \leq \sum_{k > n} 2^{-k} = 2^{-n}$. However, $\mathbf{P}\{U_n \in \mathbb{Q}\} = 1$ and $\mathbf{P}\{U_\infty \in \mathbb{Q}\} = 0$. This shows that $X_n \xrightarrow{\text{a.s.}} X_\infty$ does not in general imply that

$$\mathbf{P}\{X_n \in A\} \rightarrow \mathbf{P}\{X_\infty \in A\}$$

for all $A \in \mathcal{B}(\mathbb{R})$; more care is needed.

An easy example also shows that convergence in probability does not imply almost sure convergence. Let $(B_n, n \geq 1)$ be independent with B_n a Bernoulli($1/n$) random variable, which is to say $\mathbf{P}\{B_n = 1\} = 1/n = 1 - \mathbf{P}\{B_n = 0\}$. Then for all $\epsilon \in (0, 1)$,

$$\mathbf{P}\{|B_n - 0| > \epsilon\} = \mathbf{P}\{B_n = 1\} = \frac{1}{n} \rightarrow 0$$

as $n \rightarrow \infty$, so $B_n \xrightarrow{\text{P}} 0$. However, $\sum_{n \geq 1} \mathbf{P}\{B_n = 1\} = \sum_{n \geq 1} 1/n = \infty$, so by the second Borel-Cantelli lemma, $\mathbf{P}\{B_n = 1 \text{ i.o.}\} = 1$. It follows that

$$\mathbf{P}\left\{\lim_{n \rightarrow \infty} B_n = 0\right\} = \mathbf{P}\left\{\{B_n = 1 \text{ i.o.}\}^c\right\} = 1 - \mathbf{P}\{B_n = 1 \text{ i.o.}\} = 0.$$

Thus B_n does not converge to 0 almost surely.

We now turn from warning examples to positive results.

Proposition 5.7. *Let $(X_n, n \geq 1)$ be a sequence of random variables defined on a common probability space $(\Omega, \mathcal{F}, \mathbf{P})$. If $X_n \xrightarrow{\text{a.s.}} X_\infty$ then $X_n \xrightarrow{\text{P}} X_\infty$.*

Proof. Fix $\epsilon > 0$. Then we have

$$\mathbf{P}\left\{\lim_{n \rightarrow \infty} X_n = X_\infty\right\} \leq \mathbf{P}\left\{\limsup_{n \rightarrow \infty} |X_n - X_\infty| \leq \epsilon\right\} = \mathbf{P}\left\{\exists n \in \mathbb{N} : \sup_{m \geq n} |X_m - X_\infty| \leq \epsilon\right\}.$$

The sequence of events $\{\sup_{m \geq n} |X_m - X_\infty| \leq \epsilon\}$ is increasing in n , and its limit is the event

$$\left\{\limsup_{n \rightarrow \infty} |X_n - X_\infty| \leq \epsilon\right\},$$

so

$$\mathbf{P}\left\{\limsup_{n \rightarrow \infty} |X_n - X_\infty| \leq \epsilon\right\} = \lim_{n \rightarrow \infty} \mathbf{P}\left\{\sup_{m \geq n} |X_m - X_\infty| \leq \epsilon\right\} \leq \lim_{n \rightarrow \infty} \mathbf{P}\{|X_n - X_\infty| \leq \epsilon\}.$$

It follows that if $\mathbf{P}\{\lim_{n \rightarrow \infty} X_n = X_\infty\} = 1$ then $\lim_{n \rightarrow \infty} \mathbf{P}\{|X_n - X_\infty| \leq \epsilon\} = 1$. \square

Proposition 5.8. *Let $(X_n, n \geq 1)$ be a sequence of random variables defined on a common probability space $(\Omega, \mathcal{F}, \mathbf{P})$. If $X_n \xrightarrow{\text{P}} X_\infty$ then there exists a subsequence $(n_k, k \geq 1)$ such that $X_{n_k} \xrightarrow{\text{a.s.}} X_\infty$ as $k \rightarrow \infty$.*

Proof. Suppose that $X_n \xrightarrow{\text{P}} X_\infty$. Then for each $k \in \mathbb{N}$, we may choose $n_k \in \mathbb{N}$ large enough that $\mathbf{P}\{|X_m - X_\infty| > 1/k\} < 1/2^k$ for all $m \geq n_k$. The n_k can clearly be chosen to be increasing, so that $(n_k, k \geq 1)$ is indeed a subsequence of \mathbb{N} . Then

$$\sum_{k \geq 1} \mathbf{P}\{|X_{n_k} - X_\infty| > 1/m\} \leq m + \sum_{k \geq m} \frac{1}{2^k} < \infty,$$

so by the first Borel-Cantelli lemma, $\mathbf{P}\{|X_{n_k} - X_\infty| > 1/m \text{ i.o.}\} = 0$. Thus

$$\begin{aligned} \mathbf{P}\left\{\lim_{k \rightarrow \infty} X_{n_k} \neq X_\infty\right\} &= \mathbf{P}\left\{\exists m \in \mathbb{N} : \limsup_{k \rightarrow \infty} |X_{n_k} - X_\infty| > 1/m\right\} \\ &\leq \sum_{m \in \mathbb{N}} \mathbf{P}\left\{\limsup_{k \rightarrow \infty} |X_{n_k} - X_\infty| > 1/m \text{ i.o.}\right\} \\ &= 0. \end{aligned} \quad \square$$

Proposition 5.9. *Let $(X_n, n \geq 1)$ be a sequence of random variables defined on a common probability space $(\Omega, \mathcal{F}, \mathbf{P})$. If $X_n \xrightarrow{p} X_\infty$ then $X_n \xrightarrow{d} X_\infty$.*

Proof. First, note that for any random variable X , for each $x \in \mathbb{R}$ with $\mathbf{P}\{X = x\} > 0$ the interval

$$(\mathbf{P}\{X < x\}, \mathbf{P}\{X \leq x\})$$

is non-empty, and these intervals are pairwise disjoint for different points $x, y \in \mathbb{R}$. Thus, if for each $x \in \mathbb{R}$ with $\mathbf{P}\{X = x\} > 0$ we choose a point $q(x) \in (\mathbf{P}\{X < x\}, \mathbf{P}\{X = x\}) \cap \mathbb{Q}$, then the values $q(x)$ are distinct rational numbers. We have thus defined an injective map from $\{x \in \mathbb{R} : \mathbf{P}\{X = x\} > 0\}$ to \mathbb{Q} , so $\{x \in \mathbb{R} : \mathbf{P}\{X = x\} > 0\}$ is at countable.

Now fix $x \in \mathbb{R}$ with $\mathbf{P}\{X_\infty = x\} = 0$. Then since $\{X_\infty < x\}$ is the increasing limit of the events $\{X_\infty \leq x - \delta\}$ as $\delta \downarrow 0$, we have

$$\mathbf{P}\{X_\infty \leq x\} = \mathbf{P}\{X_\infty < x\} = \lim_{\delta \downarrow 0} \mathbf{P}\{X_\infty \leq x - \delta\}.$$

Also, by continuity from above, $\mathbf{P}\{X_\infty \leq x\} = \lim_{\delta \downarrow 0} \mathbf{P}\{X_\infty \leq x + \delta\}$. Thus, for all $\epsilon > 0$ there is $\delta > 0$ such that

$$\mathbf{P}\{X_\infty \leq x\} - \epsilon < \mathbf{P}\{X_\infty \leq x - \delta\} \leq \mathbf{P}\{X_\infty \leq x + \delta\} < \mathbf{P}\{X_\infty \leq x\} + \epsilon.$$

Now, if $X_n \leq x$ then either $X_\infty \leq x + \delta$ or $|X_n - X_\infty| > \delta$, so

$$\begin{aligned} \limsup_{n \rightarrow \infty} \mathbf{P}\{X_n \leq x\} &\leq \limsup_{n \rightarrow \infty} (\mathbf{P}\{X_\infty \leq x + \delta\} + \mathbf{P}\{|X_n - X_\infty| > \delta\}) \\ &= \mathbf{P}\{X_\infty \leq x - \delta\} \leq \mathbf{P}\{X_\infty \leq x\} + \epsilon. \end{aligned}$$

Likewise, if $X_\infty \leq x - \delta$ then either $X_n \leq x$ or $|X_n - X_\infty| > \delta$, so

$$\begin{aligned} \liminf_{n \rightarrow \infty} \mathbf{P}\{X_n \leq x\} &\geq \liminf_{n \rightarrow \infty} (\mathbf{P}\{X_\infty \leq x - \delta\} - \mathbf{P}\{|X_n - X_\infty| > \delta\}) \\ &= \mathbf{P}\{X_\infty \leq x - \delta\} \geq \mathbf{P}\{X_\infty \leq x\} - \epsilon. \end{aligned}$$

Since $\epsilon > 0$ was arbitrary, this completes the proof. \square

For the last, and perhaps most interesting, implication between different modes of convergence, we require an additional definition. Fix a collection of measures $(\mu_i, i \in I)$. A *coupling* of $(\mu_i, i \in I)$ is a collection $(Y_i, i \in I)$ of random variables defined on a *common* probability space $(\Omega, \mathcal{F}, \mathbf{P})$ such that $\mu_{Y_i} = \mu_i$ for all $i \in I$. If $(X_i, i \in I)$ is a collection of random variables, possibly defined on different probability spaces, with $\mu_{X_i} = \mu_i$, we might also refer to $(Y_i, i \in I)$ as a coupling of $(X_i, i \in I)$.

For example, suppose that μ_1 and μ_2 are both the uniform measure on the set $[6] = \{1, 2, 3, 4, 5, 6\}$. Then with $\Omega = [6]$, $\mathcal{F} = 2^{[6]}$ and \mathbf{P} the uniform measure on Ω , setting $Y_1(\omega) = \omega$ and $Y_2(\omega) = 7 - \omega$ gives a coupling of μ_1 and μ_2 .⁸ Alternately, with $\Omega = [6]^2 = \{(i, j), 1 \leq i, j \leq 6\}$, $\mathcal{F} = 2^\Omega$, and \mathbf{P} the uniform measure on Ω , setting $Y_1(i, j) = i$ and $Y_2(i, j) = j$ gives another coupling of μ_1 and μ_2 ; this is an “independent coupling” since Y_1 and Y_2 are independent. By Theorem 5.5, if $(\mu_i, i \in I)$ is a countable collection of probability measures then a coupling of $(\mu_i, i \in I)$ always exists.

Theorem 5.10 (Skorohod representation theorem). *Fix random variables $(X_n, 1 \leq n \leq \infty)$, with $X_n : \Omega_n \rightarrow \mathbb{R}$ for some probability space $(\Omega_n, \mathcal{F}_n, \mathbf{P}_n)$. If $X_n \xrightarrow{d} X_\infty$ then there exists a coupling $(Y_n, 1 \leq n \leq \infty)$ of $(X_n, 1 \leq n \leq \infty)$ such that $Y_n \xrightarrow{a.s.} Y_\infty$.*

Proof. We write $F_n = F_{X_n}$. Our coupling lives on the probability space

$$(\Omega, \mathcal{F}, \mathbf{P}) = ([0, 1], \mathcal{B}([0, 1]), \text{Leb}_{[0,1]}).$$

⁸This is the “glass table” coupling of a die roll: the value that comes up and the value seen by someone lying under the table.

For $1 \leq n \leq \infty$, let $Y_n : \Omega \rightarrow \mathbb{R}$ be defined by

$$Y_n(p) = \inf\{x : F_n(x) \geq p\}.$$

Then by the same argument as in the proof of Proposition 5.4, we have $F_{Y_n} = F_n$ for all n , so $(Y_n, 1 \leq n \leq \infty)$ is indeed a coupling of $(X_n, 1 \leq n \leq \infty)$. The bulk of the proof consists in showing that $Y_n \xrightarrow{\text{a.s.}} Y_\infty$.

Note that for all $1 \leq n \leq \infty$, $Y_n(p)$ is increasing in p , so has at most countably many points of discontinuity (reprising the argument from the start of Proposition 5.9 gives an injective map from the discontinuity points into \mathbb{Q}). Thus to prove that $Y_n \xrightarrow{\text{a.s.}} Y_\infty$ it is sufficient to prove that $Y_n(p) \rightarrow Y_\infty(p)$ whenever Y_∞ is continuous at p .

So fix $p \in [0, 1]$ a continuity point of Y_∞ , and write $y = Y_\infty(p) = \inf\{x \in \mathbb{R} : F_\infty(x) \geq p\}$. Then $F_\infty(x) < p$ for $x < y$. Writing $p' = F_\infty(y)$, by right-continuity of F_∞ we must have $p' \geq p$. Moreover, since p is a continuity point of F_∞ we must have $F_\infty(z) > p$ for all $z > y$. (If $F_\infty(z) = p$ for some $z > y$ then for all $q > p$ we have $Y_\infty(q) \geq z$, contradicting that p is a continuity point.)

Now fix $\epsilon > 0$, and choose $x < y < z$ with x, z continuity points of F_∞ and such that

$$y - \epsilon < x < y < z < y + \epsilon.$$

Then $F_\infty(x) < p$ and $F_\infty(z) > p$. Since z is a continuity point of F_∞ and $X_n \xrightarrow{d} X_\infty$, it follows that

$$F_n(z) \rightarrow F_\infty(z) > p$$

so $F_n(z) > p$ for all n sufficiently large. Thus $Y_n(p) \leq z < y + \epsilon$ for n large. Likewise, $F_n(x) \rightarrow F_\infty(x) < p$, so $F_n(x) < p$ for n large. Thus for $Y_n(p) \geq x > y - \epsilon$ for n large. Since $\epsilon > 0$ was arbitrary, it follows that $Y_n(p) \rightarrow Y_\infty(p)$, as required. \square

6. Integration and expectation

Let $(\Omega, \mathcal{F}, \mu)$ be a σ -finite measure space. In this section, unless otherwise specified, when we refer to a *measurable function* f , we mean a $\mathcal{F}/\mathcal{B}(\mathbb{R})$ -measurable function $f : \Omega \rightarrow \mathbb{R}$.⁹ We say that an event $E \in \mathcal{F}$ occurs μ -almost everywhere, or μ -a.e., if $\mu(E^c) = 0$.

Our aim is to define the (definite) integral

$$\int f d\mu \equiv \int_\Omega f d\mu \equiv \int_\Omega f(x)\mu(dx) \equiv \mu(f)$$

for as rich a class of measurable functions as possible. The preceding display lists four different bits of notation for this integral; **these notes use at least the first three.**

The way the integral is defined is by starting from functions taking only finitely many values, where the correct definition of the integral is obvious, then taking limits. We say a measurable function f is *simple* if it takes only finitely many values. Thus, f is simple if for some $n \in \mathbb{N}$ there are sets $E_1, \dots, E_n \in \mathcal{F}$ and constants $c_1, \dots, c_n \in \mathbb{R}$ such that $f = \sum_{i=1}^n c_i \mathbf{1}_{E_i}$.

Exercise 6.1. For any simple function $f : \Omega \rightarrow \mathbb{R}$, there is a unique choice of pairwise disjoint sets $D_1, \dots, D_\ell \in \mathcal{F}$ and of distinct constants $b_1, \dots, b_\ell \in \mathbb{R}$ such that $f = \sum_{i=1}^\ell b_i \mathbf{1}_{D_i}$.

Let $f = \sum_{i=1}^\ell b_i \mathbf{1}_{D_i}$ be a simple function from Ω to \mathbb{R} , with (D_1, \dots, D_ℓ) pairwise disjoint and (b_1, \dots, b_ℓ) distinct. We say f is *integrable* if $\mu(D_i) < \infty$ for all $1 \leq i \leq \ell$; if this holds then we define

$$\int_\Omega f d\mu = \sum_{i=1}^\ell c_i \mu(E_i).$$

If $\mu(\Omega) < \infty$ then every simple function is integrable. The next exercise says that for a simple integrable function, the definition of the integral doesn't depend on the representation of f as a sum of indicators of sets of bounded measure.

⁹Most of what follows also works if $f : \Omega \rightarrow \mathbb{R}^*$ is $\mathcal{F}/\mathcal{B}(\mathbb{R}^*)$ -measurable, provided one takes appropriate care around situations where $\infty - \infty$ might show up.

Simple function.

Exercise 6.2. Suppose that $\sum_{i=1}^n c_i \mathbf{1}_{[E_i]} = \sum_{i=1}^m d_i \mathbf{1}_{[F_i]}$ define the same function (where $E_1, \dots, E_n \in \mathcal{F}$ and $F_1, \dots, F_m \in \mathcal{F}$ all have finite measure, and $c_1, \dots, c_n, d_1, \dots, d_m \in \mathbb{R}$). Then $\sum_{i=1}^n c_i \mu(E_i) = \sum_{i=1}^m d_i \mu(F_i)$.

The next proposition states some basic properties of the integral for simple integrable functions.

Proposition 6.1. Let $(\Omega, \mathcal{F}, \mu)$ be a σ -finite measure space and let $f, g : \Omega \rightarrow \mathbb{R}$ be simple integrable functions.

- If $f \geq 0$ μ -a.e. then $\int f d\mu \geq 0$.
- If $a \in \mathbb{R}$ then $\int af + gd\mu = a \int f d\mu + \int gd\mu$.
- If $f \leq g$ μ -a.e. then $\int f d\mu \leq \int g d\mu$.

Proof. Write $f = \sum_{i=1}^n c_i \mathbf{1}_{[E_i]}$ with $E_1, \dots, E_n \in \mathcal{F}$ disjoint. If some $c_i < 0$ then since $f \geq 0$ μ -almost everywhere we must have $\mu(E_i) = 0$. Thus

$$\int f d\mu = \sum_{i:c_i>0} c_i \mu(E_i) \geq 0,$$

proving (a). Next, write $g = \sum_{j=1}^m d_j \mathbf{1}_{[F_j]}$. Then $af + g = a \sum_{i=1}^n c_i \mathbf{1}_{[E_i]} + \sum_{j=1}^m d_j \mathbf{1}_{[F_j]}$ is simple so by definition

$$\int af + gd\mu = a \sum_{i=1}^n c_i \mu(E_i) + \sum_{j=1}^m d_j \mu(F_j) = a \int f d\mu + \int g d\mu,$$

proving (b). Finally, if $f \leq g$ μ -a.e. then $g - f \geq 0$ μ -a.e. so by (a) and (b),

$$0 \leq \int g - f d\mu = \int g d\mu - \int f d\mu,$$

proving (c). □

In what follows we'll sometimes write “ f s.i.” to mean that f is simple and integrable. We extend the definition from simple functions first to non-negative functions, then to general functions. For f a non-negative measurable function, define

$$\int f d\mu = \sup_{\substack{g \leq f \\ g \text{ s.i.}}} \int g d\mu.$$

Note that if f is itself simple then for $g \leq f$ simple we have $\int g d\mu \leq \int f d\mu$ by the previous proposition; it follows that this new definition agrees with the previous definition when f is simple.

One may think of this definition as a “horizontal” definition via lower approximations, whereas the Riemann integral uses a “vertical” approximation. Alternatively, one may say that the Riemann approximation to the integral decomposes the domain, whereas the above definition (one might call it a “Lebesgue approximation”) decomposes the range.

Finally, for a general measurable function f , write $f^+ = \max(f, 0)$ and $f^- = -\min(f, 0)$. If either $\int f^+ d\mu < \infty$ or $\int f^- d\mu < \infty$ then we set

$$\int f d\mu = \int f^+ d\mu - \int f^- d\mu,$$

and say the integral of f is defined. Note that if $f \geq 0$ then $f = f^+$ and $f^- = 0$, so this definition agrees with the definition for non-negative functions.

Having extended the definition of the integral from simple functions to this more general class, we now need to check again that the basic properties of the integral all hold.

Proposition 6.2. Let $(\Omega, \mathcal{F}, \mu)$ be a σ -finite measure space and let f, g and $(f_n, n \geq 1)$ be measurable functions.

- **Weak monotonicity.** If $f \leq g$ and $\int f d\mu$ and $\int g d\mu$ are defined then $\int f d\mu \leq \int g d\mu$.
- **Weak monotone convergence theorem.** If $f_n \geq 0$ and $f_n \uparrow f$, then $\int f_n d\mu \uparrow \int f d\mu$.

f s.i.: simple integrable

Definition of f^+ and f^- for a function f ; note that we use this notation differently earlier in the notes.

• **Linearity of expectation.** If $f, g \geq 0$ and $a \geq 0$ then $\int af + g d\mu = a \int f d\mu + \int g d\mu$.

To prove linearity of expectation, we need the following lemma.

Lemma 6.3. Let $f \geq 0$ be measurable. Then there exist non-negative simple functions $(f_n, n \geq 1)$ such that $f_n \uparrow f$ as $n \rightarrow \infty$.

s

Proof. For $0 \leq k < n \cdot 2^n$ let $B_{n,k} = \{k/2^n \leq f < (k+1)/2^n\}$. Then set

$$f_n = \sum_{k=0}^{n2^n-1} \frac{k}{2^n} \mathbf{1}_{[B_{n,k}]}$$

Then $0 \leq f_n \leq f$, and $f \mathbf{1}_{[f \leq n]} \leq f_n + 1/2^n$ so $\liminf_{n \rightarrow \infty} f_n \geq \liminf_{n \rightarrow \infty} (f \mathbf{1}_{[f \leq n]} - 2^{-n}) = f$. \square

For later use, we remark that the functions f_n constructed in the course of proving the above theorem are all $(\sigma(f)/\mathcal{B}(\mathbb{R}))$ -measurable. This means that if $(\Omega, \mathcal{F}, \mathbf{P})$ is a probability space and $X : \Omega \rightarrow \mathbb{R}$ and $Y : \Omega \rightarrow \mathbb{R}$ are non-negative independent random variables, then there exist $(\sigma(X)/\mathcal{B}(\mathbb{R}))$ -measurable random variables $(X_n, n \geq 1)$ and $(\sigma(Y)/\mathcal{B}(\mathbb{R}))$ -measurable random variables $(Y_n, n \geq 1)$ such that $X_n \uparrow X$ and $Y_n \uparrow Y$ as $n \rightarrow \infty$. The collections $(X_n, n \geq 1)$ and $(Y_n, n \geq 1)$ of random variables then independent due to the independence of X and Y . This extends to more than two variables in an obvious way.

Proof of Proposition 6.2. If $f \geq 0$ then this is obvious because the supremum in the definition of $\int g d\mu$ is over a larger set than in the definition of $\int f d\mu$. For general f , since $f \leq g$ we have $f^+ \leq g^+$ and $f^- \geq g^-$ so

$$\int f d\mu = \int f^+ d\mu - \int f^- d\mu \leq \int g^+ d\mu - \int g^- d\mu = \int g d\mu.$$

This proves the first assertion.

Next, suppose $0 \leq f_n \uparrow f$. Then for each n by monotonicity we have $f_n \leq f$ so $\int f_n d\mu \leq \int f d\mu$, so

$$\lim_{n \rightarrow \infty} \int f_n d\mu = \sup_{n \in \mathbb{N}} \int f_n d\mu \leq \int f d\mu.$$

To prove the reverse inequality, fix any simple function $g = \sum_{i=1}^m c_i \mathbf{1}_{[E_i]}$ with $0 \leq g \leq f$. We may assume that (E_1, \dots, E_m) are disjoint and that $c_i > 0$ for all $1 \leq i \leq m$.

First suppose $\int f d\mu = \infty$. For $n \geq 1$ let

$$E_{i,n} = E_i \cap \{f_n > c_i/2\} = \{\omega \in E_i : f_n(\omega) \geq c_i/2\}.$$

Then $E_{i,n} \uparrow E_i$ as $n \rightarrow \infty$, so $\mu(E_{i,n}) \uparrow \mu(E_i)$. Since also $f_n \geq \sum_{i=1}^m (c_i/2) \mathbf{1}_{[E_{i,n}]}$, it follows that

$$\begin{aligned} \liminf_{n \rightarrow \infty} \int f_n d\mu &\geq \liminf_{n \rightarrow \infty} \int \sum_{i=1}^m (c_i/2) \mathbf{1}_{[E_{i,n}]} d\mu \\ &= \liminf_{n \rightarrow \infty} \frac{1}{2} \sum_{i=1}^m c_i \mu(E_{i,n}) \\ &= \liminf_{n \rightarrow \infty} \frac{1}{2} \sum_{i=1}^m c_i \mu(E_i) \\ &= \frac{1}{2} \int g d\mu. \end{aligned}$$

Thus

$$\liminf_{n \rightarrow \infty} \int f_n d\mu \geq \frac{1}{2} \sup_{\substack{g \leq f \\ g \text{ simple}}} \int g d\mu = \infty.$$

Next suppose $\int f d\mu < \infty$. Fix $\epsilon > 0$, let $\delta = \epsilon / \int g d\mu$, and for $n \geq 1$ let

$$E_{i,n} = E_i \cap \{f_n > c_i - \epsilon\}.$$

Then again $E_{i,n} \uparrow E_i$ as $n \rightarrow \infty$, so $\mu(E_{i,n}) \uparrow \mu(E_i)$. Moreover,

$$f_n \geq \sum_{i=1}^m (c_i - \delta) \mathbf{1}_{[E_{i,n}]} \geq \sum_{i=1}^m c_i \mathbf{1}_{[E_{i,n}]} - \delta \sum_{i=1}^m \mathbf{1}_{[E_i]},$$

so

$$\begin{aligned} \liminf_{n \rightarrow \infty} \int f_n d\mu &\geq \liminf_{n \rightarrow \infty} \left(\int \sum_{i=1}^m c_i \mathbf{1}_{[E_{i,n}]} d\mu - \int \delta \sum_{i=1}^m \mathbf{1}_{[E_i]} d\mu \right) \\ &= \liminf_{n \rightarrow \infty} \left(\int \sum_{i=1}^m c_i \mathbf{1}_{[E_{i,n}]} d\mu \right) - \delta \sum_{i=1}^m c_i \mu(E_i) \\ &= \liminf_{n \rightarrow \infty} \sum_{i=1}^m c_i \mu(E_{i,n}) - \epsilon \\ &= \sum_{i=1}^m c_i \mu(E_i) - \epsilon. \end{aligned}$$

Since $\epsilon > 0$ and $g \leq f$ were arbitrary, it follows that

$$\liminf_{n \rightarrow \infty} \int f_n d\mu \geq \int f d\mu$$

as before. This proves the second assertion.

Finally, fix non-negative measurable functions f, g and constant $a \geq 0$. Then let $(f_n, n \geq 1)$ and $(g_n, n \geq 1)$ be simple functions with $0 \leq f_n \uparrow f$ and $0 \leq g_n \uparrow g$. Then $a f_n + g_n \uparrow a f + g$, so

$$\int c f + g d\mu = \lim_{n \rightarrow \infty} \int c f_n + g_n d\mu = \lim_{n \rightarrow \infty} c \int f_n d\mu + \int g_n d\mu = c \int f d\mu + \int g d\mu,$$

where we have used monotone convergence, plus linearity of integration for simple functions, in the above string of identities. This completes the proof. \square

Notice that linearity of integration for non-negative functions implies that $\int |f| d\mu = \int f^+ d\mu + \int f^- d\mu$, since $|f| = f^+ + f^-$. If $\int |f| d\mu < \infty$ we say that f is μ -integrable and write $f \in L_1(\mu)$.

$L_1(\mu)$, μ -integrable

Exercise 6.3. (a) Show that if f, g are μ -integrable and $a \in \mathbb{R}$ then $\int a f + g d\mu = a \int f d\mu + \int g d\mu$.
(b) Let $f \geq 0$ be measurable. Show that $\int f d\mu = 0$ if and only if $f = 0$ μ -almost everywhere.

Proposition 6.4 (Monotonicity of integrals). If $f \leq g$ μ -almost everywhere and both integrals are defined, then $\int f d\mu \leq \int g d\mu$.

Proof. Write $\hat{g} = g + (f - g) \mathbf{1}_{[f > g]}$. Then $\hat{g} \geq f$ so

$$\int \hat{g} d\mu \geq \int f d\mu.$$

But $\hat{g} - g = (f - g) \mathbf{1}_{[f > g]}$ is non-negative and μ -a.e. equals zero, so

$$\int g d\mu = \int \hat{g} d\mu - \int (\hat{g} - g) d\mu = \int \hat{g} d\mu. \quad \square$$

Note that Proposition 6.4 implies that if $f \stackrel{\mu\text{-a.e.}}{=} g$ and $\int f d\mu$ is defined, then $\int g d\mu$ is defined and $\int f d\mu = \int g d\mu$.

We now state and prove the fundamental convergence theorems for sequences of functions. In all of them, $(f_n, n \geq 1)$, f , and g are measurable functions defined on a σ -finite measure space $(\Omega, \mathcal{F}, \mu)$. The first result, the (strong) monotone convergence theorem, is really a corollary of the weak monotone convergence theorem combined with the previous proposition.

Theorem 6.5 (Monotone convergence theorem). *If $(f_n, n \geq 1)$ and f are measurable functions and $0 \leq f_n \uparrow f$ holds μ -almost everywhere then*

$$\int f_n d\mu \rightarrow \int f d\mu,$$

as $n \rightarrow \infty$.

Proof. Let

$$E = \{\omega \in \Omega : f_n(\omega) \uparrow f(\omega) \text{ as } n \rightarrow \infty\}.$$

Then $\mu(E^c) = 0$ by assumption. Writing $f'_n = f_n \mathbf{1}_{[E]}$ and $f' = f \mathbf{1}_{[E]}$, then $0 \leq f'_n \uparrow f'$ so by the weak monotone convergence theorem $\int f'_n d\mu \rightarrow \int f' d\mu$. But $f'_n \stackrel{\mu\text{-a.e.}}{=} f_n$ and $f' \stackrel{\mu\text{-a.e.}}{=} f$, so Proposition 6.4 we have

$$\int f_n d\mu = \int f'_n d\mu \quad \text{and} \quad \int f d\mu = \int f' d\mu$$

and the result follows. □

Theorem 6.6 (Fatou's lemma). *If $f_n \geq 0$ for all n then*

$$\int \liminf_{n \rightarrow \infty} f_n d\mu \leq \liminf_{n \rightarrow \infty} \int f_n d\mu.$$

Proof. Note that $\inf_{k \geq n} f_k$ is increasing in n , and its limit is $\liminf_{n \rightarrow \infty} f_n$, so by the monotone convergence theorem

$$\int \liminf_{n \rightarrow \infty} f_n d\mu = \lim_{n \rightarrow \infty} \int \inf_{k \geq n} f_k d\mu.$$

But for each $k \geq n$, $\int \inf_{k \geq n} f_k d\mu \leq \int f_k d\mu$, so

$$\lim_{n \rightarrow \infty} \int \inf_{k \geq n} f_k d\mu \leq \lim_{n \rightarrow \infty} \inf_{k \geq n} \int f_k d\mu = \liminf_{n \rightarrow \infty} \int f_n d\mu. \quad \square$$

Theorem 6.7 (Dominated convergence theorem). *Suppose that $f_n \rightarrow f$ μ -almost everywhere. If there exists $g \in L_1(\mu)$ such that $|f_n| \leq g$ μ -almost everywhere then*

$$\int f_n d\mu \rightarrow \int f d\mu.$$

Proof. We now know that changing a function on a set of measure zero doesn't change its integral, so we can assume that $f_n \rightarrow f$ and $|f_n| \leq g$ for all n . It follows that $|f| \leq g$ so $f \in L_1(\mu)$ as well.

Now apply Fatou's lemma to both $g + f_n$ and $g - f_n$; since $\liminf_{n \rightarrow \infty} g + f_n = g + f$ and $\liminf_{n \rightarrow \infty} g - f_n = g - f$ we obtain

$$\int g + f d\mu = \int \liminf_{n \rightarrow \infty} (g + f_n) d\mu \leq \liminf_{n \rightarrow \infty} \int (g + f_n) d\mu = \int g d\mu + \liminf_{n \rightarrow \infty} \int f_n d\mu,$$

and

$$\int g - f d\mu = \int \liminf_{n \rightarrow \infty} (g - f_n) d\mu \leq \liminf_{n \rightarrow \infty} \int (g - f_n) d\mu = \int g d\mu - \limsup_{n \rightarrow \infty} \int f_n d\mu.$$

Subtracting $\int g d\mu$ from both equations, this gives

$$\int f d\mu \leq \liminf_{n \rightarrow \infty} \int f_n d\mu \leq \limsup_{n \rightarrow \infty} \int f_n d\mu \leq \int f d\mu,$$

so the limit of $\int f_n d\mu$ must exist and equal $\int f d\mu$. \square

Corollary 6.8. *Suppose $\mu(\Omega) < \infty$. If $f_n \rightarrow f$ μ -almost everywhere and there is $M > 0$ such that $|f_n| \leq M$ for all $n \geq 1$, then*

$$\int f_n d\mu \rightarrow \int f d\mu.$$

Proof. In this case the constant function $g \equiv M$ satisfies $\int g d\mu = M\mu(\Omega) < \infty$, and $|f_n| \leq g$ for all $n \geq 1$. \square

Exercise 6.4. (a) *Fix $g \in L_1(\mu)$. Suppose that $\sum_{n \geq 1} f_n$ converges μ -almost everywhere and that $|\sum_{k=1}^n f_k| \leq g$ μ -almost everywhere, for all $n \geq 1$. Show that $f_n \in L_1(\mu)$ for all $n \geq 1$, that $\sum_{n \geq 1} f_n \in L_1(\mu)$ and that*

$$\int \sum_{n \geq 1} f_n d\mu = \sum_{n \geq 1} \int f_n d\mu.$$

(b) *Suppose that $\sum_{n \geq 1} \int |f_n| d\mu < \infty$. Prove that $\sum_{n \geq 1} f_n$ is μ -a.e. absolutely convergent and that $\sum_{n \geq 1} \int f_n d\mu = \int \sum_{n \geq 1} f_n d\mu$.*

6.1. Expectation and independence. All the theorems of the preceding section can be applied to real random variables defined over a probability space $(\Omega, \mathcal{F}, \mathbf{P})$. In this setting, for a random variable $X : \Omega \rightarrow \mathbb{R}$ we have one additional way to write the integral: $\mathbf{E}X := \int X d\mathbf{P}$; in this setting the integral is called the *expected value* of X .

So the theorems of the preceding section imply, for example, that if $0 \leq X_n \uparrow X$ almost surely then $\mathbf{E}X_n \rightarrow \mathbf{E}X$ as $n \rightarrow \infty$; and if $|X_n| \leq M$ for all n and $X_n \xrightarrow{\text{a.s.}} X$ then $|X| \leq M$ almost surely and $\mathbf{E}X_n \rightarrow \mathbf{E}X$.

The main goal of the current section is to exhibit the strong connection between independence of random variables and factorization of expectations into product form.

Theorem 6.9 (Independence means multiply). *Let $(X_i, 1 \leq i \leq n)$ be random variables defined over a probability space $(\Omega, \mathcal{F}, \mathbf{P})$. Then $(X_i, 1 \leq i \leq n)$ are independent if and only if*

$$\mathbf{E} \left[\prod_{k=1}^n f_k(X_k) \right] = \prod_{k=1}^n \mathbf{E}[f_k(X_k)] \quad (6.1)$$

for any bounded Borel measurable functions $f_k : \mathbb{R} \rightarrow \mathbb{R}$.

This theorem gives us the chance to introduce one of the last “simplifying techniques” of the notes: the *monotone class theorem*.

Theorem 6.10 (Monotone class theorem). *Let (Ω, \mathcal{F}) be a measurable space, and let $\mathcal{P} \subset \mathcal{F}$ be a π -system over Ω with $\Omega \in \mathcal{P}$ and $\sigma(\mathcal{P}) = \mathcal{F}$. Let \mathcal{S} be a collection of functions $f : \Omega \rightarrow \mathbb{R}$ with the following properties.*

- (a) *For all $P \in \mathcal{P}$, $\mathbf{1}_{[P]} \in \mathcal{S}$.*
- (b) *For all $f, g \in \mathcal{S}$ and $c \in \mathbb{R}$, $cf + g \in \mathcal{S}$.*
- (c) *If $(f_n, n \geq 1)$ are elements of \mathcal{S} and $0 \leq f_n \uparrow f$ for a bounded function f , then $f \in \mathcal{S}$.*

Then \mathcal{S} contains all bounded $\mathcal{F}/\mathcal{B}(\mathbb{R})$ -measurable functions.

Proof. Let $\Lambda = \{F \in \mathcal{F} : \mathbf{1}_{[F]} \in \mathcal{S}\}$. Then $\mathcal{P} \subset \Lambda$ by definition. Moreover, if $E, F \in \Lambda$ and $E \subset F$ then $\mathbf{1}_{[F \setminus E]} = \mathbf{1}_{[F]} - \mathbf{1}_{[E]} \in \mathcal{S}$ by (b) and so $F \setminus E \in \Lambda$. Also, if $F_n \uparrow F$ then $\mathbf{1}_{[F_n]} \uparrow \mathbf{1}_{[F]}$ and $\mathbf{1}_{[F]}$ is bounded so lies in \mathcal{S} ; thus $F \in \Lambda$. This means that Λ is a λ -system, containing \mathcal{P} , so contains \mathcal{F} by Dynkin’s π -system lemma (Lemma 4.5).

We now know that $\mathbf{1}_{[F]} \in \mathcal{S}$ for all $F \in \mathcal{F}$. Since by (b), the collection \mathcal{S} is closed under linear combinations, it follows that \mathcal{S} contains all simple functions. Any bounded non-negative function is a monotone limit of simple functions by Lemma 6.3, so by (c) it follows that \mathcal{S} contains all non-negative bounded measurable functions. Finally, for any bounded measurable function f , we may write $f = f^+ - f^-$ as a difference of bounded measurable functions, so another application of (b) shows that $f \in \mathcal{S}$. \square

First proof of Theorem 6.9. First suppose that (6.1) holds for any bounded measurable functions f_1, \dots, f_n . Fix any events $E_1, \dots, E_n \in \mathcal{F}$ with $E_k \in \sigma(X_k)$. Since $\sigma(X_k) = \{X_k^{-1}(B), B \in \mathcal{B}(\mathbb{R})\}$, we may write $E_k = \{X_k \in B_k\}$ for some $B_k \in \mathcal{B}(\mathbb{R})$. It follows that

$$\begin{aligned} \mathbf{P} \left\{ \bigcap_{k=1}^n E_k \right\} &= \mathbf{P} \left\{ \bigcap_{k=1}^n \{X_k \in B_k\} \right\} = \mathbf{E} \left[\prod_{k=1}^n \mathbf{1}_{[B_k]}(X_k) \right] \\ &= \prod_{k=1}^n \mathbf{E} [\mathbf{1}_{[B_k]}(X_k)] = \prod_{k=1}^n \mathbf{P} \{X_k \in B_k\} = \prod_{k=1}^n \mathbf{P} \{E_k\}. \end{aligned}$$

Thus X_1, \dots, X_n are independent.

Conversely, suppose X_1, \dots, X_n are independent. Let

$$\mathcal{S}_n := \left\{ f : \mathbb{R} \rightarrow \mathbb{R} \text{ bounded, Borel} : \forall B_1, \dots, B_{n-1} \in \mathcal{B}(\mathbb{R}), \right.$$

$$\left. \mathbf{E} \left[f(X_n) \cdot \prod_{k=1}^{n-1} \mathbf{1}_{[X_k \in B_k]} \right] = \mathbf{E} [f(X_n)] \cdot \prod_{k=1}^{n-1} \mathbf{P} \{X_k \in B_k\} \right\}.$$

Then by assumption, \mathcal{S}_n contains the indicator functions $\{\mathbf{1}_{[B]} : B \in \mathcal{B}(\mathbb{R})\}$. Moreover, if $f, g \in \mathcal{S}$ and $c \in \mathbb{R}$ then for any $B_1, \dots, B_{n-1} \in \mathcal{B}(\mathbb{R})$, by linearity of expectation,

$$\begin{aligned} \mathbf{E} \left[(cf + g)(X_n) \cdot \prod_{k=1}^{n-1} \mathbf{1}_{[X_k \in B_k]} \right] &= c \mathbf{E} \left[f(X_n) \cdot \prod_{k=1}^{n-1} \mathbf{1}_{[X_k \in B_k]} \right] + \mathbf{E} \left[g(X_n) \cdot \prod_{k=1}^{n-1} \mathbf{1}_{[X_k \in B_k]} \right] \\ &= c \mathbf{E} [f(X_n)] \cdot \prod_{k=1}^{n-1} \mathbf{P} \{X_k \in B_k\} + \mathbf{E} [g(X_n)] \cdot \prod_{k=1}^{n-1} \mathbf{P} \{X_k \in B_k\} \\ &= \mathbf{E} [(cf + g)(X_n)] \cdot \prod_{k=1}^{n-1} \mathbf{P} \{X_k \in B_k\}, \end{aligned}$$

so $cf + g \in \mathcal{S}_n$. Also, if $0 \leq f_m \uparrow f$ with f bounded and $f_m \in \mathcal{S}_n$ for all $m \geq 1$, then for any $B_1, \dots, B_{n-1} \in \mathcal{B}(\mathbb{R})$, by the monotone convergence theorem

$$\begin{aligned} \mathbf{E} \left[f(X_n) \cdot \prod_{k=1}^{n-1} \mathbf{1}_{[X_k \in B_k]} \right] &= \lim_{m \rightarrow \infty} \mathbf{E} \left[f_m(X_n) \cdot \prod_{k=1}^{n-1} \mathbf{1}_{[X_k \in B_k]} \right] \\ &= \lim_{m \rightarrow \infty} \mathbf{E} [f_m(X_n)] \cdot \prod_{k=1}^{n-1} \mathbf{P} \{X_k \in B_k\} \\ &= \mathbf{E} [f(X_n)] \cdot \prod_{k=1}^{n-1} \mathbf{P} \{X_k \in B_k\}. \end{aligned}$$

so $f \in \mathcal{S}_n$. Thus \mathcal{S}_n contains all bounded measurable functions. Next let

$$\mathcal{S}_{n-1} := \left\{ f : \mathbb{R} \rightarrow \mathbb{R} \text{ Borel} : \forall B_1, \dots, B_{n-2} \in \mathcal{B}(\mathbb{R}), \forall g : \mathbb{R} \rightarrow [0, \infty) \text{ Borel}, \right.$$

$$\left. \mathbf{E} \left[f(X_{n-1})g(X_n) \cdot \prod_{k=1}^{n-2} \mathbf{1}_{[X_k \in B_k]} \right] = \mathbf{E} f(X_{n-1}) \cdot \mathbf{E} g(X_n) \prod_{k=1}^{n-2} \mathbf{P} \{X_k \in B_k\} \right\}.$$

By repeating the same arguments as for \mathcal{S}_n , we see that \mathcal{S}_{n-1} contains all bounded measurable functions (the monotone convergence theorem can be used since we took g non-negative). Repeating this argument (i.e. by induction), we obtain that

$$\mathbf{E} \left[\prod_{k=1}^n f_k(X_k) \right] = \prod_{k=1}^n \mathbf{E} [f_k(X_k)]$$

for any non-negative bounded Borel functions f_1, \dots, f_n . Using linearity of expectation once more, it follows that this identity indeed holds for any bounded Borel functions. \square

Second proof of Theorem 6.9. This proof replaces the use of the monotone class theorem with a direct argument (which has a similar flavour). We refer to (6.1) as “the factorization formula”. The proof that if the factorization formula holds then X_1, \dots, X_n are independent is the same as in the first proof.

Now suppose that X_1, \dots, X_n are independent. Then for all $B_1, \dots, B_n \in \mathcal{B}_n$, the events $(\{X_k \in B_k\}, 1 \leq k \leq n)$ are independent, so for any constants $c_1, \dots, c_n \in \mathbb{R}$, writing $c = \prod_{i=1}^n c_i$, we have

$$\mathbf{E} \prod_{k=1}^n c_k \mathbf{1}_{[B_k]}(X_k) = c \mathbf{P} \left\{ \bigcap_{k=1}^n \{X_k \in B_k\} \right\} = c \prod_{k=1}^n \mathbf{P} \{X_k \in B_k\} = \prod_{k=1}^n \mathbf{E} c_k \mathbf{1}_{[B_k]}(X_k),$$

proving the factorization formula for (multiples of) indicator functions.

Now let f_1, \dots, f_n be simple Borel functions. Then we may write $f_k = \sum_{\ell=1}^m c_{k,\ell} \mathbf{1}_{[B_{k,\ell}]}$ for some real constants $(c_{k,\ell}, k \in [n], \ell \in [m])$ and Borel sets $(B_{k,\ell}, k \in [n], \ell \in [m])$. (We can always “pad” some of the sums so that they all have the same number of terms.) Then using linearity of expectation and the factorization formula for indicator functions,

$$\begin{aligned} \mathbf{E} \prod_{k=1}^n f_k(X_k) &= \mathbf{E} \prod_{k=1}^n \sum_{\ell=1}^m c_{k,\ell} \mathbf{1}_{[B_{k,\ell}]}(X_k) \\ &= \sum_{\ell_1, \dots, \ell_n=1}^m \mathbf{E} \prod_{k=1}^n c_{k,\ell_k} \mathbf{1}_{[B_{k,\ell_k}]}(X_k) \\ &= \sum_{\ell_1, \dots, \ell_n=1}^m \prod_{k=1}^n c_{k,\ell_k} \mathbf{E} \mathbf{1}_{[B_{k,\ell_k}]}(X_k) \\ &= \prod_{k=1}^n \mathbf{E} \sum_{\ell=1}^m c_{k,\ell} \mathbf{1}_{[B_{k,\ell}]}(X_k) \\ &= \prod_{k=1}^n \mathbf{E} f_k(X_k), \end{aligned}$$

so the factorization formula holds for simple functions.

Now suppose f_1, \dots, f_n are non-negative Borel functions, and write $Y_k = f_k(X_k)$. Then Y_k is $\sigma(X_k)/\mathcal{B}(\mathbb{R})$ -measurable, so Y_1, \dots, Y_n are independent. By Lemma 6.3, for each $1 \leq k \leq n$ we may find simple functions $(Y_{k,m}, m \geq 1)$ such that $0 \leq Y_{k,m} \uparrow Y_k$ and such that $Y_{k,m}$ is $\sigma(X_k)/\mathcal{B}(\mathbb{R})$ -measurable for all $m \in \mathbb{N}$. Then $(Y_{1,m}, \dots, Y_{n,m})$ are independent for all m , and $\prod_{k=1}^n Y_{k,m} \uparrow \prod_{k=1}^n Y_k$ as $m \rightarrow \infty$, so by the monotone convergence theorem and the factorization formula for simple functions,

$$\mathbf{E} \prod_{k=1}^n f_k(X_k) = \mathbf{E} \prod_{k=1}^n Y_k = \lim_{m \rightarrow \infty} \mathbf{E} \prod_{k=1}^n Y_{k,m} = \lim_{m \rightarrow \infty} \prod_{k=1}^n \mathbf{E} Y_{k,m} = \prod_{k=1}^n \mathbf{E} Y_k,$$

proving the factorization formula for non-negative functions.

Finally, if f_1, \dots, f_n are bounded and Borel measurable then we can again use linearity of expectation to write $f(X_k) =: Y_k = Y_k^+ - Y_k^-$, and we then have

$$\begin{aligned} \mathbf{E} \prod_{k=1}^n f_k(X_k) &= \mathbf{E} \prod_{k=1}^n (Y_k^+ - Y_k^-) \\ &= \sum_{(\sigma_1, \dots, \sigma_n) \in \{-, +\}^n} (-1)^{\#\{k \in [n] : \sigma_k = -\}} \mathbf{E} \prod_{k=1}^n Y_k^{\sigma_k} \\ &= \sum_{(\sigma_1, \dots, \sigma_n) \in \{-, +\}^n} (-1)^{\#\{k \in [n] : \sigma_k = -\}} \prod_{k=1}^n \mathbf{E} Y_k^{\sigma_k} = \prod_{k=1}^n \mathbf{E} [Y_k^+ - Y_k^-] = \prod_{k=1}^n \mathbf{E} f_k(X_k) \end{aligned}$$

so the factorization formula holds in general. □

An extremely important corollary of Theorem 6.1 is that the factorization formula holds when the functions f_1, \dots, f_n are simply the identity (this is not an immediate consequence of the theorem as the identity function is unbounded).

Corollary 6.11. *Suppose that X_1, \dots, X_n are independent and either (a) $X_k \geq 0$ for $1 \leq k \leq n$ or (b) $X_k \in L_1(\mathbf{P})$ for $1 \leq k \leq n$. If (b) holds then $\prod_{k=1}^n X_k \in L_1(\mathbf{P})$ and if either (a) or (b) hold then $\mathbf{E} \prod_{k=1}^n X_k = \prod_{k=1}^n \mathbf{E} X_k$.*

In the proof (and later in the notes?), we use the following notation: for a function $f : \Omega \rightarrow \mathbb{R}$ and $r > 0$ we write $f_{\leq r} := f \mathbf{1}_{\{|f| \leq r\}}$.

Proof. It suffices to prove the corollary when $n = 2$; the result then follows by induction. So suppose X, Y are independent. If X and Y are non-negative then by the monotone convergence theorem and by (6.1),

$$\mathbf{E}[XY] = \lim_{n \rightarrow \infty} \mathbf{E}[X_{\leq n} Y_{\leq n}] = \lim_{n \rightarrow \infty} \mathbf{E} X_{\leq n} \mathbf{E} Y_{\leq n} = \mathbf{E} X \mathbf{E} Y,$$

so if (a) holds then the factorization formula holds.

Next, if $X, Y \in L_1(\mathbf{P})$ then write $X = X^+ - X^-$ and $Y = Y^+ - Y^-$. Then $|XY| = (X^+ + X^-)(Y^+ + Y^-)$, so by linearity of expectation and the conclusion of the previous paragraph,

$$\begin{aligned} \mathbf{E}|XY| &= \mathbf{E}[X^+Y^+] + \mathbf{E}[X^+Y^-] + \mathbf{E}[X^-Y^+] + \mathbf{E}[X^-Y^-] \\ &= \mathbf{E}X^+ \mathbf{E}Y^+ + \mathbf{E}X^+ \mathbf{E}Y^- + \mathbf{E}X^- \mathbf{E}Y^+ + \mathbf{E}X^- \mathbf{E}Y^- \\ &= \mathbf{E}|X| \mathbf{E}|Y| < \infty, \end{aligned}$$

so $XY \in L_1(\mathbf{P})$. We may then again use linearity of expectation to deduce that

$$\begin{aligned} \mathbf{E}XY &= \mathbf{E}[X^+Y^+] - \mathbf{E}[X^+Y^-] - \mathbf{E}[X^-Y^+] + \mathbf{E}[X^-Y^-] \\ &= \mathbf{E}X^+ \mathbf{E}Y^+ - \mathbf{E}X^+ \mathbf{E}Y^- - \mathbf{E}X^- \mathbf{E}Y^+ + \mathbf{E}X^- \mathbf{E}Y^- \\ &= \mathbf{E}X \mathbf{E}Y. \end{aligned}$$

□

Exercise 6.5. *Fix random variables $X, Y \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ and let $\mathcal{P} \subset \mathcal{F}$ be a π -system with $\sigma(\mathcal{P}) = \mathcal{F}$. Show that if $\mathbf{E}[X \mathbf{1}_{[A]}] = \mathbf{E}[Y \mathbf{1}_{[A]}]$ for all $A \in \mathcal{P}$ then $X \stackrel{\text{a.s.}}{=} Y$.*

7. An interlude: the probabilistic method.

One of the challenges of teaching a first rigorous probability course is the amount of setup that's required before one gets to "the real stuff". Billingsley's textbook avoids this issue by focussing exclusively on simple functions in the early chapters. This makes the book more engaging at the outset; the cost is that many of the most important random variables (Gaussian, exponential, Gamma, Beta, ...) are excluded from consideration.

My approach this course has been to bite the bullet and do the necessary setup, while doing my best to motivate its development. I've also postponed a few things that in most courses would

Notation $f_{\leq r}$

have already been introduced or would be next on the menu: Fubini's theorem, convergence in L_p , Hölder, Minkowski and Cauchy-Schwartz inequalities, to name a few.

Even so, I know that the first third to half of the course can feel like a bit of a slog. To liven things up a bit, I've decided to describe one of the ways in which probability has contributed to other branches of mathematics: the *probabilistic method*. In a nutshell, the idea of the probabilistic method is this. One wishes to show the existence of a mathematical object m with some property P . Rather constructing m directly, one instead constructs a *random* object M and shows that

$$\mathbf{P} \{M \text{ has property } P\} > 0.$$

This immediately shows that there must be at least one object m with property P , proving existence.

Example 1: existence of continuous nowhere differentiable functions. Let $D_n := \{i/2^n, 0 \leq i \leq 2^n\}$, so that $D := \bigcup_{n \geq 0} D_n$ are the dyadic rationals in $[0, 1]$. Note that $D_{n-1} \subset D_n$ for each $n \geq 1$. Let $(N_x, x \in D)$ be IID $N(0, 1)$ random variables. Define a sequence of random functions B_n from $[0, 1]$ to \mathbb{R} as follows.

D_n : n 'th level dyadic rationals in $[0, 1]$.

For $x \in [0, 1]$ let $B_1(x) = xN_1$. Actually, a more wordy but equivalent definition of B_1 presages¹⁰ the subsequent construction more effectively. Let $B_1(0) = 0$ and let $B_1(1) = N_1$; then for $x \in (0, 1)$ define $B_1(x)$ by linear interpolation between points of $D_0 = \{0, 1\}$.

Inductively, given B_{n-1} , let

$$B_n(x) = \begin{cases} B_{n-1}(x) & \text{if } x \in D_{n-1} \\ B_{n-1}(x) + \frac{N_x}{\sqrt{2^n}} & \text{if } x \in D_n \setminus D_{n-1} \\ pB_n(i/2^n) + (1-p)B_n((i+1)/2^n) & \text{if } x = \frac{pi+(1-p)i+1}{2^n}, p \in (0, 1). \end{cases}$$

Then it is possible to show that¹¹

$$\mathbf{P} \{(B_n, n \geq 0) \text{ is a uniformly convergent sequence of functions}\} = 1,$$

so one may define a random function B_∞ as the almost sure limit of the sequence B_n . The limit B_∞ is *Brownian motion* on the interval $[0, 1]$. The fact that B_∞ is a.s. a uniform limit of continuous functions implies that B_∞ is a.s. continuous. However, it turns out that

a.s.: almost surely

$$\mathbf{P} \{B_\infty \text{ is nowhere differentiable}\} = 1.$$

Thus, Brownian motion provides an example of a continuous, nowhere differentiable function. In fact, Brownian motion is (in a sense which can be made precise) a *uniformly random continuous function*, so the above statement can be interpreted as stating that *almost all continuous functions are nowhere differentiable*.

Example 2: small-norm signings of vectors.

This example is a bit more down-to-earth, and we may actually prove all our statements with the machinery we currently have available to us. It is drawn from Alon and Spencer's book "The probabilistic method".

Proposition 7.1. *Let v_1, \dots, v_n be vectors in \mathbb{R}^m (for some m) with $|v_i| = 1$ for all i . Then there exist $\sigma_1, \dots, \sigma_n \in \{-1, 1\}$ such that*

$$|\sigma_1 v_1 + \dots + \sigma_n v_n| \leq \sqrt{n}.$$

Proof. Let $\sigma_1, \dots, \sigma_n$ be independent and uniform on $\{-1, 1\}$. Set

$$X = |\sigma_1 v_1 + \dots + \sigma_n v_n|^2 = (\sigma_1 v_1 + \dots + \sigma_n v_n) \cdot (\sigma_1 v_1 + \dots + \sigma_n v_n) = \sum_{i,j=1}^n \sigma_i \sigma_j v_i \cdot v_j.$$

¹⁰presage, v.: 1. transitive. a. To constitute a supernatural sign of (a future event); to be an omen of, to portend. b. To be indicative or suggestive of; to be a natural precursor of, to give warning of. –Oxford English Dictionary

¹¹discuss measurability issues?

Then by linearity of expectation,

$$\mathbf{E}X = \mathbf{E} \sum_{i,j=1}^n \sigma_i \sigma_j v_i \cdot v_j = \sum_{i,j=1}^n \mathbf{E} [\sigma_i \sigma_j] v_i \cdot v_j.$$

If $i \neq j$ then σ_i and σ_j are independent so $\mathbf{E} [\sigma_i \sigma_j] = \mathbf{E} \sigma_i \mathbf{E} \sigma_j = 0$; so the above identity simplifies to

$$\mathbf{E}X = \sum_{i=1}^n \mathbf{E} [\sigma_i^2] v_i \cdot v_i = \sum_{i=1}^n 1 = n$$

since $v_i \cdot v_i = |v_i|^2 = 1$ and $\sigma_i^2 \equiv 1$.

But if $\mathbf{E}X = n$ then $\mathbf{P} \{X \leq n\} > 0$ by monotonicity of expectations; so there must be some choice of $\sigma_1, \dots, \sigma_n$ which makes $X \leq n$, and for this choice

$$|\sigma_1 v_1 + \dots + \sigma_n v_n| = X^{1/2} \leq \sqrt{n}. \quad \square$$

8. Densities and change of variables

This section is about how to actually do computations with random variables and their expectations.

If X is a random variable taking values in \mathbb{N} , then $X = \lim_{n \rightarrow \infty} X_{\leq n}$, and this is an increasing limit. For each n , $X_{\leq n}$ is a simple function so by the definition of the integral of simple functions, we have

$$\mathbf{E}X = \lim_{n \rightarrow \infty} \mathbf{E}X_{\leq n} = \lim_{n \rightarrow \infty} \sum_{k=0}^n k \mathbf{P} \{X_{\leq n} = k\} = \sum_{k \geq 0} k \mathbf{P} \{X = k\}.$$

More generally, if X is non-negative and takes values in a countable set N , then the same sort of argument gives that $\mathbf{E}X = \sum_{n \in N} n \mathbf{P} \{X = n\}$. This allows us to do computations with discrete random variables. For example, if P is Poisson(λ) then

$$\mathbf{E}P = \sum_{k \geq 0} k \cdot \frac{\lambda^k e^{-\lambda}}{k!} = \lambda \cdot \sum_{k \geq 1} \frac{\lambda^{k-1} e^{-\lambda}}{(k-1)!} = \lambda.$$

But it's not yet clear how to tackle computations involving non-discrete random variables. For example, suppose that $(N_i, i \geq 1)$ are independent standard Gaussian random variables, independent of P . How should we compute (or even approximate)

$$\mathbf{P} \left\{ \sum_{i=1}^P N_i \geq 1 \right\} = \mathbf{E} \left[\mathbf{1}_{[\sum_{i=1}^P N_i \geq 1]} \right] = \int_{\Omega} \mathbf{1}_{[\sum_{i=1}^P N_i \geq 1]}(\omega) d\mathbf{P} ?$$

Using the tools we now develop.

Definition 8.1. Given a measure space $(\Omega, \mathcal{F}, \mu)$ and $f : \Omega \rightarrow \mathbb{R}$ non-negative and $\mathcal{F}/\mathcal{B}(\mathbb{R})$ -measurable, define a new measure μf on (Ω, \mathcal{F}) by setting

$$\mu f(A) = \int_A f d\mu := \int f \mathbf{1}_{[A]} d\mu.$$

If $\nu = \mu f$ then we say ν has density f with respect to μ .

Exercise 8.1. The function μf defined above is a measure on (Ω, \mathcal{F}) .

Exercise 8.2. If $f' \stackrel{\mu\text{-a.e.}}{=} f$ then $\mu f = \mu f'$.

We also say a real random variable X has density f with respect to Lebesgue measure if its distribution μ_X has density f with respect to Lebesgue measure, or in other words if

$$\mu_X(B) = \int_B f(x) dx$$

for any Borel $B \subset \mathbb{R}$. In this case we also say that f is the *probability density function* of X . These definitions are justified by the following two results.

Proposition 8.2. *Fix a σ -finite measure space $(\Omega, \mathcal{F}, \mu)$ and measurable $f : \Omega \rightarrow [0, \infty)$, and suppose ν has density f with respect to μ . Then for measurable $g : \Omega \rightarrow \mathbb{R}$,*

$$\int g d\nu = \int g f d\mu$$

provided that either $g \geq 0$ or $g \in L_1(\nu)$. Moreover, $g \in L_1(\nu)$ if and only if $gf \in L_1(\mu)$.

Proof. If $g = \mathbf{1}_{[A]}$ for some $A \in \mathcal{F}$ then the equality of the two integrals holds by definition. The theorem then follows straightforwardly using the monotone class theorem and the monotone convergence theorem. \square

Proposition 8.3 (Change of variables formula). *Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space, and let $X : \Omega \rightarrow \mathbb{R}$ be a random variable. Then for all measurable $g : \mathbb{R} \rightarrow \mathbb{R}$ with $\mathbf{E}|g(X)| < \infty$,*

$$\mathbf{E}g(X) = \int_{\mathbb{R}} g(x) \mu_X(dx). \quad (8.1)$$

Moreover, if X has density f with respect to Lebesgue measure then also $\mathbf{E}g(X) = \int_{\mathbb{R}} g(x) f(x) dx$.

Proof. Again, the assertions are true if $g = \mathbf{1}_{[B]}$ for Borel $B \subset \mathbb{R}$, and the rest of the proof follows using the monotone class theorem and the monotone convergence theorem. \square

More generally, if $(\Omega, \mathcal{F}, \mu)$ is a σ -finite measure space and (S, \mathcal{S}) is another measurable space, then for an $(\mathcal{F}/\mathcal{S})$ -measurable function $f : \Omega \rightarrow S$ we may define

$$\nu(E) = \mu(f^{-1}(E))$$

for $E \in \mathcal{S}$. Then for all non-negative $(\mathcal{S}/\mathcal{B}(\mathbb{R}))$ -measurable functions $g : S \rightarrow \mathbb{R}$, we have

$$\int g d\nu = \int g \circ f d\mu.$$

The change of variables formula (8.1) is a special case of this, but this also tells us, for example, that if $\mathbf{X} = (X_1, \dots, X_n)$ are random variables defined on a common space, then

$$\mathbf{E}g(X_1, \dots, X_n) = \int_{\mathbb{R}^n} g(\vec{x}) \mu_{\mathbf{X}}(\vec{x}).$$

Making this a useful computational tool, even for independent random variables, will require Fubini's theorem.

Exercise 8.3. *Prove that if $\varphi : [a, b] \rightarrow \mathbb{R}$ is C_1 (continuously differentiable) and strictly increasing then for any Borel function $g : [\varphi(a), \varphi(b)] \rightarrow [0, \infty)$,*

$$\int_{\varphi(a)}^{\varphi(b)} g(y) dy = \int_a^b g(\varphi(y)) \varphi'(y) dy.$$

Example: size-biasing the Poisson and folded normal distributions. Let X be a non-negative random variable on a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ with $0 < \mathbf{E}X < \infty$. Then $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mu_X)$ is another probability space. The *size-biased* distribution of X is $\hat{\mu}_X := (\mu_X \cdot X / \mathbf{E}X)$. In other words, for Borel $B \subset \mathbb{R}$,

$$\hat{\mu}_X(A) = (\mu_X \cdot X / \mathbf{E}X)(A) = \mathbf{E} \left[\frac{X}{\mathbf{E}X} \mathbf{1}_{[X \in A]} \right].$$

This is another probability distribution on \mathbb{R} , since $\hat{\mu}_X(\mathbb{R}) = \mathbf{E}[X / \mathbf{E}X] = 1$.

Suppose P is Poisson(λ). Then $\mathbf{E}P = \lambda$ so for Borel $B \subset \mathbb{R}$,

$$\begin{aligned} \hat{\mu}_P(B) &= \left(\mu_P \frac{P}{\lambda} \right) (B) = \mathbf{E} \left[\frac{P}{\lambda} \mathbf{1}_{\{P \in B\}} \right] = \sum_{k \geq 1, k \in B} \frac{k}{\lambda} \mathbf{P} \{P = k\} \\ &= \sum_{k \geq 1, k \in B} \frac{k}{\lambda} \frac{\lambda^k e^{-\lambda}}{k!} = \sum_{k \geq 1, k \in B} \frac{\lambda^{k-1} e^{-\lambda}}{(k-1)!} = \mathbf{P} \{P + 1 \in B\} = \mu_{P+1}(B). \end{aligned}$$

In other words, the size-biased distribution of P is just the distribution of $P + 1$. Of course, nothing like this need hold in general.

Next suppose N is a standard Gaussian; so N has density $\Phi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$ with respect to Lebesgue measure. The distribution of $|N|$ is called the *folded normal* distribution; it has density $\psi(x) = 2\Phi(x)\mathbf{1}_{[x \geq 0]} = \sqrt{2/\pi} e^{-x^2/2} \mathbf{1}_{[x \geq 0]}$ with respect to Lebesgue measure.

The size-biasing of the distribution of $|N|$ is $\hat{\mu}_{|N|} = (\mu_{|N|} \cdot |N|/\mathbf{E}|N|)$. To find an explicit formula for this, we first use the change of variables formula to compute

$$\mathbf{E}|N| = \int_{\mathbb{R}} x d\mu_{|N|} = \int_{[0, \infty)} x \cdot \sqrt{\frac{2}{\pi}} e^{-x^2/2} dx = \left[-\sqrt{\frac{2}{\pi}} e^{-x^2/2} \right]_0^\infty = \sqrt{\frac{2}{\pi}}.$$

It follows that for $B \subset [0, \infty)$ Borel,

$$\hat{\mu}_{|N|}(B) = \int \mathbf{1}_{[B]} \cdot \frac{|N|}{\mathbf{E}|N|} d\mu = \int_B \frac{x}{\sqrt{2/\pi}} \sqrt{\frac{2}{\pi}} e^{-x^2/2} dx = \int_B x e^{-x^2/2} dx.$$

Thus, $\hat{\mu}_{|N|}$ has density $x e^{-x^2/2} \mathbf{1}_{[x \geq 0]}$ with respect to Lebesgue measure. This is called the *Rayleigh distribution*.

Exercise 8.4. If X, Y are independent standard Gaussians then $\sqrt{X^2 + Y^2}$ is Rayleigh distributed.

8.1. Product measure and Fubini's theorem. Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space. If $X, Y : \Omega \rightarrow \mathbb{R}$ are independent random variables and $f, g : \mathbb{R} \rightarrow \mathbb{R}$ are bounded Borel functions then $\mathbf{E}f(X)g(Y) = \mathbf{E}f(X)\mathbf{E}g(Y)$. By Exercise 5.2 (e), the pair (X, Y) is $\Omega/\mathcal{B}(\mathbb{R}^2)$ -measurable; in other words, it is an \mathbb{R}^2 -valued random variable. What is its distribution $\mu_{(X,Y)}$? Note that by the factorization formula, if $A, B \in \mathcal{B}(\mathbb{R})$ then

$$\mu_{(X,Y)}(A \times B) = \mathbf{P} \{(X, Y) \in A \times B\} = \mathbf{P} \{X \in A\} \mathbf{P} \{Y \in B\} = \mu_X(A)\mu_Y(B).$$

The collection $\mathcal{B}(\mathbb{R}) \times \mathcal{B}(\mathbb{R}) = \{A \times B : A, B \in \mathcal{B}(\mathbb{R})\}$ generates $\mathcal{B}(\mathbb{R}^2)$, so the preceding formula uniquely identifies $\mu_{(X,Y)}$ as the *product measure* of measures μ_X and μ_Y .

This concrete example is generalized as follows. Fix measurable spaces (M, \mathcal{M}) and (N, \mathcal{N}) . Sets $A \times B \in \mathcal{M} \times \mathcal{N}$ are called *rectangles*. Let $\mathcal{M} \boxtimes \mathcal{N}$ be the *field* generated by $\mathcal{M} \times \mathcal{N}$.

Exercise 8.5. It holds that

$$\begin{aligned} \mathcal{M} \boxtimes \mathcal{N} &= \left\{ \bigcup_{i=1}^n A_i \times B_i : n \geq 1; \forall i \in [n], A_i \times B_i \in \mathcal{M} \times \mathcal{N} \right\} \\ &= \left\{ \bigcup_{i=1}^n A_i \times B_i : n \geq 1; \forall i \in [n], A_i \times B_i \in \mathcal{M} \times \mathcal{N}; A_1 \times B_1, \dots, A_n \times B_n \text{ disjoint} \right\} \end{aligned}$$

The product measurable space $(M \times N, \mathcal{M} \otimes \mathcal{N})$ is defined by setting

$$\mathcal{M} \otimes \mathcal{N} := \sigma(\mathcal{M} \times \mathcal{N}) = \sigma(A \times B : A \in \mathcal{M}, B \in \mathcal{N}) = \sigma(\mathcal{M} \boxtimes \mathcal{N}).$$

If μ and ν are σ -finite measures on (M, \mathcal{M}) and (N, \mathcal{N}) respectively, define a function $\mu \boxtimes \nu$ on $\mathcal{M} \boxtimes \mathcal{N}$ by setting

$$\mu \boxtimes \nu \left(\bigcup_{i=1}^n A_i \times B_i \right) := \sum_{i=1}^n \mu(A_i) \nu(B_i),$$

Rayleigh distribution

Rectangles
 $\mathcal{M} \boxtimes \mathcal{N}$

Product measurable space
 $\mathcal{M} \otimes \mathcal{N}$

for disjoint rectangles $A_1 \times B_1, \dots, A_n \times B_n \in \mathcal{M} \times \mathcal{N}$.

Exercise 8.6. *The function $\mu \boxtimes \nu$ is well-defined, in that if $P \in \mathcal{M} \otimes \mathcal{N}$ may be represented as a disjoint union of rectangles in multiple ways,*

$$P = \bigcup_{i=1}^n A_i \times B_i = \bigcup_{i=1}^m C_i \times D_i$$

then $\sum_{i=1}^n \mu(A_i)\nu(B_i) = \sum_{i=1}^m \mu(C_i)\nu(D_i)$.

Proposition 8.4. *If (M, \mathcal{M}, μ) and (N, \mathcal{N}, ν) are σ -finite measure spaces then $\mu \boxtimes \nu$ is a pre-measure on $\mathcal{M} \boxtimes \mathcal{N}$.*

Assuming the proposition holds, it follows by the Carathéodory Extension Theorem and Dynkin's Uniqueness theorem that $\mu \boxtimes \nu$ extends uniquely to a measure $\mu \otimes \nu$ on $\mathcal{M} \otimes \mathcal{N}$. This extension is called the *product measure* of μ and ν .

Product measure $\mu \otimes \nu$.

Exercise 8.7 (Product measure is commutative). *Suppose (M, \mathcal{M}, μ) and (N, \mathcal{N}, ν) are σ -finite measure spaces. Let $\mu \otimes \nu$ be product measure on $\mathcal{M} \otimes \mathcal{N}$ and let $\nu \otimes \mu$ be product measure on $\mathcal{N} \otimes \mathcal{M}$. Prove that for all $B \in \mathcal{M} \otimes \mathcal{N}$, $B^* := \{(b, a) : (a, b) \in B\} \in \mathcal{N} \otimes \mathcal{M}$, and $(\nu \otimes \mu)(B^*) = (\mu \otimes \nu)(B)$.*

We extract two steps of the proof of Proposition 8.4 as separate lemmas. This proof is based on the one given by Norris in his lecture notes on probability and measure.

Lemma 8.5. *Under the hypotheses of Proposition 8.4, if $f : M \otimes N \rightarrow \mathbb{R}$ is $(\mathcal{M} \otimes \mathcal{N})/\mathcal{B}(\mathbb{R})$ -measurable then for all $a \in \mathcal{M}$, the function $f_a : \mathcal{N} \rightarrow \mathbb{R}$ given by $f_a(b) := f(a, b)$ is $(\mathcal{N}/\mathcal{B}(\mathbb{R}))$ -measurable.*

Proof. Write

$$\mathcal{S} := \{f : M \times N \rightarrow \mathbb{R} : \forall a \in \mathcal{M}, f_a \text{ is } (\mathcal{N}/\mathcal{B}(\mathbb{R}))\text{-measurable}\}.$$

We aim to show \mathcal{S} contains all $(\mathcal{M} \otimes \mathcal{N})/\mathcal{B}(\mathbb{R})$ -measurable functions.

First, if $f = \mathbf{1}_{[A \times B]}$ for $A \times B \in \mathcal{M} \times \mathcal{N}$ then for $a \in A$, $f_a \equiv \mathbf{1}_{[B]}$, and for $a \notin A$, $f_a \equiv 0$. In both cases f_a is measurable so $f \in \mathcal{S}$; thus \mathcal{S} contains indicators of rectangles.

Next, if $f, g \in \mathcal{S}$ and $c \in \mathbb{R}$ then for all $a \in \mathcal{M}$,

$$(cf + g)_a(b) = (cf + g)(a, b) = cf(a, b) + g(a, b) = cf_a(b) + g_a(b) = (cf_a + g_a)(b),$$

so $(cf + g)_a$ is a linear combination of $(\mathcal{N}/\mathcal{B}(\mathbb{R}))$ -measurable functions and so is $(\mathcal{N}/\mathcal{B}(\mathbb{R}))$ -measurable. Therefore \mathcal{S} is closed under linear combinations.

Moreover, if $(f^{(n)}, n \geq 1)$ is a sequence of elements of \mathcal{S} and $0 \leq f^{(n)} \uparrow f$ for some bounded function f , then for all $a \in \mathcal{M}$, $f_a^{(n)} \uparrow f_a$. As a monotone limit of measurable functions, f_a is $(\mathcal{N}/\mathcal{B}(\mathbb{R}))$ -measurable; thus $f \in \mathcal{S}$.

Since rectangles form a π -system generating $\mathcal{M} \otimes \mathcal{N}$, by the monotone class theorem it follows that \mathcal{S} contains all bounded $(\mathcal{M} \otimes \mathcal{N})/\mathcal{B}(\mathbb{R})$ -measurable functions.

Finally, if $f : M \times N \rightarrow \mathbb{R}$ is any $(\mathcal{M} \otimes \mathcal{N})/\mathcal{B}(\mathbb{R})$ -measurable we may write f as a limit of bounded measurable functions $f = \lim_{n \rightarrow \infty} f^{(n)}$ by taking $f^{(n)} = f \mathbf{1}_{\{|f| \leq n\}}$. For all $a \in \mathcal{M}$ we then have $f_a = \lim_{n \rightarrow \infty} f_a^{(n)}$ so f_a is a limit of $(\mathcal{N}/\mathcal{B}(\mathbb{R}))$ -measurable functions and so is $(\mathcal{N}/\mathcal{B}(\mathbb{R}))$ -measurable. Thus $f \in \mathcal{S}$. \square

Lemma 8.6. *Under the hypotheses of Proposition 8.4, Let $f : M \times N \rightarrow \mathbb{R}$ be $(\mathcal{M} \otimes \mathcal{N})/\mathcal{B}(\mathbb{R})$ -measurable and either bounded or non-negative. Define $f_M : M \rightarrow \mathbb{R} \cup \{+\infty\}$ by*

$$f_M(a) := \int_N f(a, b) \nu(db).$$

If $\nu(N) < \infty$ and f is bounded then f_M is bounded and $(\mathcal{M})/\mathcal{B}(\mathbb{R})$ -measurable. Also, if f is non-negative then $f_M : M \rightarrow [0, \infty]$ is $(\mathcal{M})/\mathcal{B}(\mathbb{R}^)$ -measurable.*

Proof. Note that by the definition of f_a we have

$$f_M(a) := \int_N f_a(b)\nu(db),$$

so the integral in the lemma statement at least makes sense by Lemma 8.5.

Suppose $\nu(N) < \infty$. We wish to show that for any bounded $\mathcal{M} \otimes \mathcal{N} / \mathcal{B}(\mathbb{R})$ -measurable function f , the function f_M is $\mathcal{M} / \mathcal{B}(\mathbb{R})$ -measurable.

First, if $f = \mathbf{1}_{[A \times B]}$ for $A \times B \in \mathcal{M} \times \mathcal{N}$, then for $a \in A$,

$$f_M(a) = \int_N \mathbf{1}_{[B]}(b)\nu(db) = \nu(B) < \infty,$$

and for $a \notin A$, $f_M(a) = \int_N 0\nu(db) = 0$. Thus $f_M \equiv \nu(B)\mathbf{1}_{[A]}$ is bounded and $\mathcal{M} / \mathcal{B}(\mathbb{R})$ -measurable. Next, if f, g are bounded functions such that f_M and g_M are $\mathcal{M} / \mathcal{B}(\mathbb{R})$ -measurable and $c \in \mathbb{R}$ then $(cf + g)_M = cf_M + g_M$ by linearity of integration, so $(cf + g)_M$ is bounded and $\mathcal{M} / \mathcal{B}(\mathbb{R})$ -measurable. Finally, if $0 \leq f^{(n)} \uparrow f$ then by the monotone convergence theorem,

$$f_M(a) = \int_N f_a(b)\nu(db) = \lim_{n \rightarrow \infty} \int_N f_a^{(n)}(b)\nu(db) = \lim_{n \rightarrow \infty} f_M^{(n)}(a),$$

so f_M is an increasing limit of measurable functions and thus measurable. The first assertion of the lemma follows by the monotone class theorem.

The second assertion of the lemma follows by a similar argument. \square

In the course of the preceding proof, we derived that if $f = \mathbf{1}_{[A \times B]}$ for a rectangle $A \times B$, then $f_M = \nu(B)\mathbf{1}_{[A]}$, which implies that

$$\int_M \int_N f(a, b)\nu(db)\mu(da) = \int_M f_M(a)\mu(da) = \int_M \nu(B)\mathbf{1}_{[A]}(a)\mu(da) = \nu(B)\mu(A); \quad (8.2)$$

we will use this in the next proof.

Proof of Proposition 8.4. The function $\mu \boxtimes \nu$ is obviously non-negative, and it is additive by definition. To prove $\mu \boxtimes \nu$ is a pre-measure, it suffices to show that it is countably additive.

So suppose that $\bigcup_{i=1}^k A_i \times B_i \in \mathcal{M} \boxtimes \mathcal{N}$ is a finite disjoint union of rectangles which may also be represented as an infinite disjoint union of rectangles,

$$\bigcup_{i=1}^k A_i \times B_i = \bigcup_{i \geq 1} C_i \times D_i.$$

Using (8.2), we have

$$\mu \otimes \nu \left(\bigcup_{i=1}^k A_i \times B_i \right) = \sum_{i=1}^k \mu(A_i)\nu(B_i) = \sum_{i=1}^k \int_M \int_N \mathbf{1}_{[A_i \times B_i]}(a, b)\mu(da)\nu(db).$$

We may use linearity of integration twice to bring the sum inside the two integrals in the final term. Since $\sum_{i=1}^k \mathbf{1}_{[A_i \times B_i]} = \mathbf{1}_{[\bigcup_{i=1}^k A_i \times B_i]}$, it follows that

$$\mu \otimes \nu \left(\bigcup_{i=1}^k A_i \times B_i \right) = \int_M \int_N \mathbf{1}_{[\bigcup_{i=1}^k A_i \times B_i]} d\mu d\nu = \int_M \int_N f d\mu d\nu,$$

where we have taken $f := \mathbf{1}_{[\bigcup_{i=1}^k A_i \times B_i]}$.

Now write $f^{(n)} = \mathbf{1}_{[\bigcup_{i=1}^n C_i \times D_i]}$. Repeating the above logic gives

$$\mu \otimes \nu \left(\bigcup_{i=1}^n C_i \times D_i \right) = \int_M \int_N \mathbf{1}_{[\bigcup_{i=1}^n C_i \times D_i]} d\mu d\nu = \int_M \int_N f^{(n)} d\mu d\nu,$$

Also, $f^{(n)} \uparrow f$ since $\bigcup_{i=1}^{\infty} C_i \times D_i = \bigcup_{i=1}^k A_i \times B_i$, so for all $a \in M$, $f_a^{(n)} \uparrow f_a$, so by the monotone convergence theorem,

$$f_M^{(n)}(a) = \int_N f_a^{(n)}(b) \nu(db) \nearrow \int_N f_a(b) \nu(db) = f_M(a).$$

Since this convergence is monotone, another application of the monotone convergence theorem gives that

$$\int_M \int_N f^{(n)} d\nu d\mu = \int_M f_M^{(n)} d\mu \rightarrow \int_M f d\mu = \int_M \int_N f d\nu d\mu = \mu \otimes \nu \left(\bigcup_{i \geq 1} C_i \times D_i \right).$$

But also

$$\int_M \int_N f^{(n)} d\nu d\mu = \mu \otimes \nu \left(\bigcup_{i=1}^n C_i \times D_i \right) = \sum_{i=1}^n \mu(C_i) \nu(D_i) \rightarrow \sum_{i=1}^{\infty} \mu(C_i) \nu(D_i).$$

Comparing the two preceding displays, we see that

$$\mu \otimes \nu \left(\bigcup_{i \geq 1} C_i \times D_i \right) = \sum_{i=1}^{\infty} \mu(C_i) \nu(D_i) = \sum_{i \geq 1} \mu \otimes \nu(C_i \times D_i);$$

thus $\mu \otimes \nu$ is indeed a pre-measure. □

Theorem 8.7 (Fubini's theorem). *Let (M, \mathcal{M}, μ) and (N, \mathcal{N}, ν) be σ -finite measure spaces, and let $f : M \times N \rightarrow \mathbb{R}$ be $(\mathcal{M} \otimes \mathcal{N})/\mathcal{B}(\mathbb{R})$ -measurable.*

(a) *If $f \geq 0$ then*

$$\int f d(\mu \otimes \nu) = \int_M \int_N f d\nu d\mu. \quad (8.3)$$

(b) *If $f \in L_1(\mu \otimes \nu)$ then with $F := \{a \in M : \int_N |f(a, b)| \nu(db) < \infty\}$, it holds that $\mu(M \setminus F) = 0$. Moreover, setting*

$$\hat{f}_M(a) = \begin{cases} \int_N f(a, b) \nu(db) & \text{if } a \in F \\ 0 & \text{if } a \notin F \end{cases}$$

then $\hat{f}_M \in L_1(\mu)$, and

$$\int \hat{f}_M d\mu = \int f d(\mu \otimes \nu).$$

Part (b) of the theorem implies the following. Set $Z = F \times N$. Then $(\mu \otimes \nu)(Z^c) = \mu(M \setminus F) \nu(N) = 0 \cdot \nu(N) = 0$, so $f \mathbf{1}_{[Z]} : M \times N \rightarrow \mathbb{R}$ is $(\mu \otimes \nu)$ -a.e. equal to f , and

$$\int f d(\mu \otimes \nu) = \int f \mathbf{1}_{[Z]} d(\mu \otimes \nu) = \int_M \int_N f(a, b) \mathbf{1}_{[Z]}(a, b) \nu(db) \mu(da). \quad (8.4)$$

The only thing preventing us from removing the indicator from the double integral is that otherwise there can exist points $a \in M$ where the inner integral is not defined.

Proof. We first assume both measure spaces are finite. First, if $f = \mathbf{1}_{[A \times B]}$ for $A \times B \in \mathcal{M} \times \mathcal{N}$ then the identity holds by (8.2). Write

$$\mathcal{S} = \left\{ f : M \times N \rightarrow \mathbb{R} : \int f d(\mu \otimes \nu) = \int_M \int_N f d\nu d\mu \right\}.$$

Using linearity of integration and the monotone convergence theorem, it is not hard to check the conditions to see that \mathcal{S} satisfies the conditions of the monotone class theorem. It then follows that (8.2) holds for all bounded $(\mathcal{M} \otimes \mathcal{N})/\mathcal{B}(\mathbb{R})$ -measurable functions $f : M \times N \rightarrow \mathbb{R}$.

Next, suppose f is non-negative, and for $n \geq 1$ write $f^{(n)} = \min(f, n)$. Then $f^{(n)} \uparrow f$ so by the monotone convergence theorem

$$\int f^{(n)} d\mu \otimes \nu \uparrow \int f d\mu \otimes \nu.$$

Writing $f_a^{(n)}(b) = f^{(n)}(a, b)$, for all $a \in M$, we also have $f_a^{(n)} \uparrow f_a$, so

$$f_M^{(n)}(a) = \int_N f_a^{(n)}(b) \nu(db) \nearrow \int_N f_a(b) \nu(db),$$

so

$$\int_M \int_N f^{(n)} d\nu d\mu \rightarrow \int_M \int_N f^{(n)} d\nu d\mu.$$

Since $f^{(n)}$ is bounded, we have

$$\int f^{(n)} d(\mu \otimes \nu) = \int_M \int_N f^{(n)} d\nu d\mu$$

for all n , so it follows that $\int f d(\mu \otimes \nu) = \int_M \int_N f d\nu d\mu$, proving (a).

Next suppose $f \in L_1(\mu \otimes \nu)$ and let

$$|f|_M(a) := \int_N |f(a, b)| \nu(db), \quad f_M^+(a) := \int_N f^+(a, b) \nu(db), \quad \text{and} \quad f_M^-(a) = \int_N f^-(a, b) \nu(db).$$

Note that all three functions are $(\mathcal{M} \otimes \mathcal{N})/\mathcal{B}(\mathbb{R})$ -measurable by Lemma 8.6; the lemma only guarantees this with $\mathcal{B}(\mathbb{R})$ replaced by $\mathcal{B}(\mathbb{R}^*)$, but the condition that $f \in L_1(\mu \otimes \nu)$ ensures that everything stays finite. Since $|f| \geq 0$, we may apply part (a) of the theorem to deduce that

$$\int_M |f|_M d\mu = \int_M \int_N |f| d\nu d\mu = \int |f| d(\mu \otimes \nu) < \infty.$$

Thus $|f|_M$ is μ -almost everywhere finite; i.e., $\mu(M \setminus F) = 0$.

Finally, $\hat{f}_M = (f_M^+ - f_M^-) \mathbf{1}_{[F]}$, at least if we are willing to accept the convention that $(\infty - \infty) \cdot 0 = 0$, and so

$$\begin{aligned} \int f d(\mu \otimes \nu) &= \int f^+ d(\mu \otimes \nu) - \int f^- d(\mu \otimes \nu) && \text{linearity of integration} \\ &= \int f_M^+ d\mu - \int f_M^- d\mu && \text{by part (a)} \\ &= \int (f_M^+ - f_M^-) d\mu && \text{linearity of integration} \\ &= \int (f_M^+ - f_M^-) \mathbf{1}_{[F]} d\mu && \text{since } \mathbf{1}_{[F]} \stackrel{\mu\text{-a.e.}}{=} 1 \\ &= \int \hat{f}_M d\mu, \end{aligned}$$

proving (b).

The extension to the case that (M, \mathcal{M}, μ) and (N, \mathcal{N}, ν) are σ -finite follows by letting $(M_k, k \geq 1)$ be measurable sets in \mathcal{M} with $\mu(M_k) < \infty$ and $M_k \uparrow M$, and $(N_k, k \geq 1)$ be measurable sets in \mathcal{N} with $\mu(N_k) < \infty$ and $N_k \uparrow N$. The finite measure case of Fubini's theorem can be applied to the restriction $(M_k \times N_k, \mathcal{M}_k \otimes \mathcal{N}_k, \mu_k \otimes \nu_k)$, where $\mathcal{M}_k = \mathcal{M}|_{M_k}$ and $\mu_k = \mu|_{M_k}$, and N_k and ν_k are defined accordingly. The conclusions of Fubini's theorem in the σ -finite case can then be deduced by letting $k \rightarrow \infty$ and applying the monotone convergence theorem and linearity of integration. \square

By an exactly parallel development to the above, we may obtain an analogue of Fubini's theorem for the product measure $\nu \otimes \mu$, where the iterated integral has \int_M as the inner integral. By Exercise 8.7, it follows that (8.3) extends to the identity

$$\int f d(\mu \otimes \nu) = \int_M \int_N f d\nu d\mu = \int_N \int_M f d\mu d\nu, .$$

Proposition 8.8. *Under the conditions of Fubini's theorem, if $f \in L_1(\mu \otimes \nu)$ then there exists $E \in \mathcal{M} \otimes \mathcal{N}$ with $\mu \otimes \nu(Z^c) = 0$ such that*

$$\int f d(\mu \otimes \nu) = \int_M \int_N f(a, b) \mathbf{1}_{[E]}(a, b) \nu(db) \mu(da) = \int_N \int_M f(a, b) \mathbf{1}_{[E]}(a, b) \nu(da) \mu(db) .$$

Proof. Applying Fubini's theorem, we may obtain sets Z_M and Z_N as in (8.4), i.e., so that $\mu \otimes \nu(Z_M^c) = \mu \otimes \nu(Z_N^c) = 0$ and so that

$$\int f d(\mu \otimes \nu) = \int_M \int_N f(a, b) \mathbf{1}_{[Z_M]}(a, b) \nu(db) \mu(da)$$

and

$$\int f d(\mu \otimes \nu) = \int_N \int_M f(a, b) \mathbf{1}_{[Z_N]}(a, b) \nu(db) \mu(da) .$$

Taking $E = Z_M \cap Z_N$, the result follows. \square

Corollaries added Oct 22

Corollary 8.9. *Let $(\Omega, \mathcal{F}, \mathbf{P})$ and (M, \mathcal{M}) , (N, \mathcal{N}) be measurable spaces, and let $X : \Omega \rightarrow M$ and $Y : \Omega \rightarrow N$ be independent random variables (M -valued and N -valued, respectively), with distributions μ and ν . If $h : M \times N \rightarrow \mathbb{R}$ is $(\mathcal{M} \otimes \mathcal{N} / \mathcal{B}(\mathbb{R}))$ -measurable and either $h \geq 0$ or $\mathbf{E}|h(X, Y)| < \infty$, then*

$$\mathbf{E}h(X, Y) = \int_M \int_N h(x, y) \nu(dy) \mu(dx) .$$

Proof. For all $A \in \mathcal{M}$ and $B \in \mathcal{N}$, by independence,

$$\mathbf{P}\{(X, Y) \in A \times B\} = \mathbf{P}\{X \in A\} \mathbf{P}\{Y \in B\} = \mu(A) \nu(B) .$$

Since $\mathcal{M} \times \mathcal{N}$ is a π -system generating $\mathcal{M} \otimes \mathcal{N}$, it follows that the distribution of (X, Y) is $\mu \otimes \nu$. If either $h \geq 0$ or $\mathbf{E}|h(X, Y)| < \infty$, it then follows by the change of variables formula and Fubini's theorem that

$$\mathbf{E}h(X, Y) = \int_{M \times N} h(x, y) d(\mu \otimes \nu) = \int_M \int_N h(x, y) \nu(dy) \mu(dx) . \quad \square$$

Corollary 8.10. *In the setting of Corollary 8.9, for any $E \in \mathcal{M} \otimes \mathcal{N}$,*

$$\mathbf{P}\{(X, Y) \in E\} = \int_M \mathbf{P}\{(x, Y) \in E\} \mu(dx) .$$

Proof. Apply Corollary 8.9 to the non-negative function $h(x, y) = \mathbf{1}_{[(x, y) \in E]}$ to get

$$\mathbf{P}\{(X, Y) \in E\} = \mathbf{E}h(X, Y) = \int_M \int_N \mathbf{1}_{[(x, y) \in E]} \nu(dy) \mu(dx) = \int_M \mathbf{P}\{(x, Y) \in E\} \mu(dx) ,$$

the last equality holding by change of variables. \square

Exercise 8.8. *If X and Y are independent real random variables, and X and Y have respective densities f and g with respect to Lebesgue measure on \mathbb{R} , then (X, Y) has density $h(x, y) = f(x)g(y)$ with respect to Lebesgue measure on \mathbb{R}^2 .*

The final exercise of the section describes an important instance of the “independence means multiply” heuristic, and provides a natural segue to the following section, which is about sums of independent random variables. Given a random variable X with distribution $\mu_X = \mu$, the *moment generating function* of X is

$$G_X(s) := \mathbf{E}[e^{-sX}] = \int_{\mathbb{R}} e^{-sx} \mu(dx) \in (0, \infty] .$$

Exercise 8.9. *If X, Y are independent random variables then $G_{X+Y} = G_X G_Y$.*

Moment generating function. Notation disagrees with that later.

9. Sums of independent random variables

9.1. **Convolutions.** If μ, ν are Borel measures on \mathbb{R} then the *convolution* $\mu * \nu$ is the Borel measure on \mathbb{R} given by

$$\mu * \nu(B) = \int_{\mathbb{R}} \nu(B - x) \mu(dx),$$

for Borel B , where $B - x := \{b - x : b \in B\}$. (Exercise: to check that this definition makes sense, verify that $x \mapsto \nu(B - x)$ is a Borel function.)

Proposition 9.1. *If X, Y are independent random variables with respective laws μ and ν then $X + Y$ has law $\mu * \nu$.*

Proof. For any Borel $A \subset \mathbb{R}$, by Fubini's theorem,

$$\mathbf{P}\{X + Y \in A\} = \int_{\mathbb{R}} \int_{\mathbb{R}} \mathbf{1}_{[x+y \in A]} \nu(dy) \mu(dx) = \int_{\mathbb{R}} \nu(A - x) \mu(dx). \quad \square$$

If $f, g : \mathbb{R} \rightarrow [0, \infty)$ are Borel functions then we likewise define the convolution of f and g as

$$f * g(x) = \int_{\mathbb{R}} f(x - y)g(y)dy.$$

The next exercise states that the connection between convolution and sums of independent random variables also holds at the level of densities.

Exercise 9.1. *If X and Y are independent real random variables, and X and Y have respective densities f and g with respect to Lebesgue measure on \mathbb{R} , then $X + Y$ has density $f * g$ with respect to Lebesgue measure on \mathbb{R} .*

Exercise 9.2. *Let μ, ν be Borel measures on \mathbb{R} and let $f, g : \mathbb{R} \rightarrow [0, \infty)$ be Borel functions. Prove that $\mu * \nu = \nu * \mu$ and that $f * g = g * f$.*

It's worth seeing an example. For $\alpha, \gamma > 0$, the *Gamma*(α, λ) density is

$$\gamma(x) = \gamma_{\alpha, \lambda}(x) := \frac{\lambda^\alpha x^{\alpha-1} e^{-\lambda x}}{\Gamma(\alpha)} \mathbf{1}_{[x \geq 0]}.$$

Here $\Gamma(\alpha) := \int_{[0, \infty]} x^{\alpha-1} e^{-x} dx$ is the Gamma function. A real random variable X is *Gamma*(α, λ)-distributed if it has density $\gamma_{\alpha, \lambda}$ with respect to Lebesgue measure. The next exercise describes a scaling property of Gamma random variables in the second coordinate.

Exercise 9.3. *If X is Gamma(α, λ)-distributed then λX is Gamma($\alpha, 1$)-distributed.*

Suppose X and Y are independent, X is Gamma(α, λ)-distributed and Y is Gamma(β, λ)-distributed. We claim that $Z = X + Y$ is Gamma($\alpha + \beta, \lambda$)-distributed.

To see this, first note that by Exercise 9.3 we may assume $\lambda = 1$. (I.e. it suffices to show that $\lambda X + \lambda Y$ is Gamma($\alpha + \beta, 1$)-distributed.) We restrict to this case, and then note that by the above exercise, the density of Z with respect to Lebesgue measure is

$$\begin{aligned} f_Z(x) &= \int_{[0, x]} \gamma_{\alpha, 1}(y) \gamma_{\beta, 1}(x - y) dy \\ &= \int_{[0, x]} \frac{y^{\alpha-1} e^{-y}}{\Gamma(\alpha)} \frac{(x - y)^{\beta-1} e^{-(x-y)}}{\Gamma(\beta)} dy \\ &= \frac{e^{-x}}{\Gamma(\alpha)\Gamma(\beta)} \int_0^x y^{\alpha-1} (x - y)^{\beta-1} dy. \end{aligned}$$

Making the change of variables $u = y/x$, this becomes

$$f_Z(x) = \frac{x^{\alpha+\beta-1} e^{-x}}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 u^{\alpha-1} (1 - u)^{\beta-1} du.$$

Everything in this section works for random variables taking values in a separable Banach space, but we restrict to \mathbb{R} for concreteness.

Since $\int_{[0,\infty]} f_Z(x) dx = \mathbf{P}\{Z \geq 0\} = 1$ and, by definition, $\int_{[0,\infty]} x^{\alpha+\beta-1} e^{-x} = \Gamma(\alpha + \beta)$, it follows that

$$1 = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 u^{\alpha-1} (1-u)^{\beta-1} du,$$

which combined with the preceding display gives that

$$f_Z(x) = \frac{x^{\alpha+\beta-1} e^{-x}}{\Gamma(\alpha + \beta)},$$

so Z is indeed Gamma($\alpha + \beta, 1$)-distributed.

Another important example is introduced in the next exercise. For $\alpha \in \mathbb{R}$ and $\sigma > 0$, the $N(\alpha, \sigma^2)$ density is given by

$$\varphi_{\alpha, \sigma^2}(x) := \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\alpha)^2/(2\sigma^2)}.$$

- Exercise 9.4.** (a) Use change of variables and Fubini's theorem to prove that $(\int_{\mathbb{R}} e^{-x^2} dx)^2 = \pi$. (You've perhaps seen this before and know how the proof goes. If not: look for an integral over \mathbb{R}^2 , and consider a switch to polar coordinates.)
- (b) Show that if X and Y are independent normals with densities $\varphi_{\alpha, \sigma^2}(x)$ and φ_{β, τ^2} respectively, then $X + Y$ has density $\varphi_{\alpha+\beta, \sigma^2+\tau^2}$; in particular $X + Y$ is again a normal random variable.

10. Laws of large numbers

In the previous section we saw that summing independent random variables corresponds to convolution of their distributions. What happens if there are large number of summands? If X_1, \dots, X_n are independent random variables with a common distribution μ , then by Proposition 9.1, their sum $S_n := X_1 + \dots + X_n$ has distribution μ^{*n} , the n -fold convolution of μ with itself. *Laws of large numbers* describe the first-order asymptotic behaviour of S_n (or equivalently of μ^{*n}) when $n \rightarrow \infty$.

Rather than jumping straight to the most general results, we start with a result that has an easy proof, and has the advantage of introducing one of the most important techniques for controlling the behaviour of random variables, namely *moment methods*. These are essentially all variants of the following simple inequality

Proposition 10.1 (Markov's inequality). *If X is a non-negative random variable then $\mathbf{P}\{X \geq t\} \leq \mathbf{E}X/t$ for all $t > 0$.*

Proof. Since $X \geq X \mathbf{1}_{[X \geq t]}$, by monotonicity and by linearity of expectation,

$$\mathbf{E}X \geq \mathbf{E}[X \mathbf{1}_{[X \geq t]}] \geq \mathbf{E}[t \mathbf{1}_{[X \geq t]}] = t \mathbf{P}\{X \geq t\}. \quad \square$$

Here are some important special cases. For a random variable X with $\mathbf{E}|X| < \infty$, we write $\mathbf{Var}(X) := \mathbf{E}[(X - \mathbf{E}X)^2] \in [0, \infty]$; the quantity $\mathbf{Var}(X)$ is called the *variance* of X .

Variance of a random variable

Corollary 10.2 (Chebyshev's inequality). *For any random variable X with $\mathbf{E}|X| < \infty$, for all $t > 0$,*

$$\mathbf{P}\{|X - \mathbf{E}X| \geq t\} \leq \frac{\mathbf{Var}(X)}{t^2}.$$

Proof. Note that $|X - \mathbf{E}X| \geq t$ if and only if $(X - \mathbf{E}X)^2 \geq t^2$; then apply Markov's inequality. \square

Corollary 10.3 (Chernoff bound). *For any random variable X , for all $t \in \mathbb{R}$,*

$$\mathbf{P}\{X \geq t\} \leq \inf_{a>0} \frac{\mathbf{E}[e^{aX}]}{e^{at}}.$$

Proof. Fix $c > 0$. Then by Markov's inequality,

$$\mathbf{P} \{X \geq t\} = \mathbf{P} \{e^{aX} \geq e^{at}\} \leq \frac{\mathbf{E} [e^{aX}]}{e^{at}}.$$

Since this bound holds for each $a > 0$, the result follows. \square

In general, if X is a random variable taking values in a (possibly unbounded) interval $I \subseteq \mathbb{R}$ and $\phi : I \rightarrow [0, \infty)$ is strictly increasing, then for any we may use Markov's inequality to obtain that for any $t \in I$,

$$\mathbf{P} \{X \geq t\} = \mathbf{P} \{\phi(X) \geq \phi(t)\} \leq \frac{\mathbf{E}\phi(X)}{\phi(t)};$$

both Chebyshev's inequality and the Chernoff bound are special cases of this general bound.

We next use Chebyshev's inequality and the Chernoff bound to control the deviations of sums of independent random variables from their expected values. Before giving the details, we make a few simple observations.

Let X be a random variable with and let $0 \leq q \leq p$. Then

$$\mathbf{E} [|X|^q] \leq \mathbf{E} [\max(1, |X|^q)] \leq \mathbf{E} [\max(1, |X|^p)] \leq \mathbf{E} [1 + |X|^p], \tag{10.1}$$

so if $\mathbf{E} [|X|^p] < \infty$ then $\mathbf{E} [|X|^q] < \infty$. In particular, if $\mathbf{E} [X^2] < \infty$ then $X \in L_1(\mathbf{P})$ and so by linearity of expectation,

$$\mathbf{Var} (X) = \mathbf{E} [(X - \mathbf{E}X)^2] = \mathbf{E} [X^2 - 2X\mathbf{E}X + (\mathbf{E}X)^2] = \mathbf{E} [X^2] - (\mathbf{E}X)^2 \leq \mathbf{E} [X^2]. \tag{10.2}$$

Also, if a random variable X almost surely satisfies $a \leq X \leq b$ then we always have $|X - \mathbf{E}X| \leq b - a$, and so

$$\mathbf{Var} (X) = \mathbf{E} [(X - \mathbf{E}X)^2] \leq |b - a|^2.$$

Exercise 10.1. Strengthen the above bound to $|b - a|^2/4$.

Example: Gaussian tails for sums of bounded random variables. Fix $C > 1$ and let $(X_i, i \geq 1)$ be independent random variables with $|X_i| \leq C$ for all i . As before, write $x_i = \mathbf{E}X_i$, let $S_n := X_1 + \dots + X_n$ and let $s_n = \mathbf{E}S_n = \sum_{i=1}^n x_i$. Then by the Chernoff bound,

$$\begin{aligned} \mathbf{P} \{|S_n - s_n| \geq t\} &\leq \inf_{a>0} e^{-at} \mathbf{E} [e^{a(S_n - s_n)}] \\ &= \inf_{a>0} e^{-at} \mathbf{E} \left[\prod_{i=1}^n e^{a(X_i - x_i)} \right] \\ &= \inf_{a>0} e^{-at} \prod_{i=1}^n \mathbf{E} [e^{a(X_i - x_i)}], \end{aligned}$$

where we have used the factorization formula in the last step. We now use that if $|y| \leq 1$ then $|e^y - 1 - y| \leq y^2$. Since $|X_i| \leq C$, necessarily $|X_i - x_i| \leq 2C$, so if $a \leq 1/(2C)$ then

$$e^{a(X_i - x_i)} \leq 1 + a(X_i - x_i) + a^2(X_i - x_i)^2.$$

Taking $t = xn^{1/2}$ and $a = x/(2C^2n^{1/2})$, we then obtain

$$\mathbf{P} \{|S_n - s_n| \geq xn^{1/2}\} \leq e^{-x^2/2C^2} \prod_{i=1}^n (\mathbf{E} [1 + a(X_i - x_i) + a^2(X_i - x_i)^2])$$

For each $i \in [n]$, by linearity of expectation and since $\mathbf{E} [X_i - x_i] = 0$ and $\mathbf{Var} (X_i) \leq C^2$,

$$\mathbf{E} [1 + a(X_i - x_i) + a^2(X_i - x_i)^2] = 1 + a^2\mathbf{E} [(X_i - x_i)^2] \leq 1 + a^2C^2 = 1 + \frac{x^2}{4C^2n}.$$

Issue with absolute value in this ex.

Combining this with the preceding bound gives

$$\begin{aligned} \mathbf{P} \left\{ |S_n - s_n| \geq xn^{1/2} \right\} &\leq e^{-x^2/2C^2} \left(1 + \frac{x^2}{4C^2n} \right)^n \\ &\leq e^{-x^2/2C^2} e^{x^2/4C^2} \\ &= e^{-x^2/4C^2}, \end{aligned}$$

where in the second inequality we used that $1 + x \leq e^x$.

The next example introduces the notation of *covariance* of random variables. If X, Y are random variables with $\mathbf{E}|X|, \mathbf{E}|Y|, \mathbf{E}|XY| < \infty$, the covariance of X and Y is defined to be $\text{Cov}(X, Y) = \mathbf{E}[(X - \mathbf{E}X)(Y - \mathbf{E}Y)]$. If $\text{Cov}(X, Y)$ is defined and equals zero, then X and Y are said to be *uncorrelated*. Chebyshev's inequality gives clean bounds for sums of uncorrelated random variables, that are useful frequently enough to be stated as a separate corollary.

Covariance, uncorrelated random variables

Corollary 10.4 (Chebyshev's inequality for sums). *Let $(X_i, i \geq 1)$ be uncorrelated random variables with $\mathbf{E}|X_i| < \infty$ for all $i \geq 1$. Let $S_n := X_1 + \dots + X_n$ and let $s_n = \mathbf{E}S_n$. Then for all $t > 0$,*

$$\mathbf{P} \{ |S_n - s_n| \geq t \} \leq \frac{\sum_{i=1}^n \mathbf{Var}(X_i)}{t^2}.$$

Proof. Write $x_i = \mathbf{E}X_i$, so $s_n = x_1 + \dots + x_n$. Then

$$\begin{aligned} \mathbf{Var}(S_n) &= \mathbf{E}[(S_n - s_n)^2] \\ &= \mathbf{E} \left[\left((X_1 - x_1) + \dots + (X_n - x_n) \right)^2 \right] \\ &= \sum_{i=1}^n \mathbf{E}[(X_i - x_i)^2] + \sum_{1 \leq i \neq j \leq n} \mathbf{E}[(X_i - x_i)(X_j - x_j)]. \end{aligned}$$

Since the random variables are uncorrelated, for $i \neq j$ we have $\mathbf{E}[(X_i - x_i)(X_j - x_j)] = \mathbf{E}[X_i - x_i] \mathbf{E}[X_j - x_j] = 0$, so the second sum vanishes. The first sum is simply $\sum_{i=1}^n \mathbf{Var}(X_i)$, so the result follows by Chebyshev's inequality. \square

Example: weak law of large numbers for uncorrelated random variables. Using Chebyshev's inequality for sums, we can easily prove a first law of large numbers.

Theorem 10.5 (Weak law of large numbers for sums of uncorrelated random variables with bounded variance). *Let $(X_i, i \geq 1)$ be independent random variables with $\sup_{i \geq 1} \mathbf{E}[X_i^2] = C < \infty$. Write $x_i = \mathbf{E}X_i$, let $S_n := X_1 + \dots + X_n$ and let $s_n = \mathbf{E}S_n$. Then*

$$\frac{S_n - s_n}{n} \rightarrow 0 \tag{10.3}$$

in probability.

Proof. By Chebyshev's inequality for sums we have

$$\mathbf{P} \{ |S_n - s_n| > t \} \leq \frac{1}{t^2} \sum_{i=1}^n \mathbf{Var}(X_i) \leq \frac{Cn}{t^2}.$$

In the last line we have used that $\mathbf{Var}(X_i) \leq \mathbf{E}[X_i^2] \leq C$ for each $1 \leq i \leq n$. For any $\epsilon > 0$, taking $t = \epsilon n$ above gives

$$\mathbf{P} \left\{ \left| \frac{S_n - s_n}{n} \right| > \epsilon \right\} = \mathbf{P} \{ |S_n - s_n| > \epsilon n \} \leq \frac{\mathbf{Var}(S_n)}{\epsilon^2 n^2} \leq \frac{C}{\epsilon^2 n} \rightarrow 0$$

as $n \rightarrow \infty$; thus $(S_n - s_n)/n \rightarrow 0$ in probability as claimed. \square

Remarks.

- We have just proved a weak law of large numbers for independent random variables with bounded second moments. We'll next see how to combine this with *truncation* and Markov's inequality to prove that $S_n/s_n \rightarrow 1$ under only a first-moment assumption, but additionally assuming that the random variables are identically distributed.
- If the random variables $(X_i, i \geq 1)$ are also identically distributed, then $s_n = n\mathbf{E}X_1$, in which case (10.3) asserts that $S_n/n \rightarrow \mathbf{E}X_1$ in probability; this is a more classical way to state a law of large numbers.

Exercise 10.2. *Modify the above proof to show that, under the same assumptions, if $f : \mathbb{N} \rightarrow [0, \infty)$ and $f(n) \rightarrow \infty$ as $n \rightarrow \infty$ then*

$$\frac{S_n - s_n}{f(n)n^{1/2}} \rightarrow 0$$

in probability, as $n \rightarrow \infty$.

We now use the same idea for random variables with possibly infinite variance (but additionally assuming the random variables are identically distributed). We obviously can't directly use the same proof in this case; we will instead argue by *truncation*.

Theorem 10.6. *Let $(X_n, n \geq 1)$ be independent identically distributed random variables with $\mathbf{E}|X_n| < \infty$, and write $S_n = X_1 + \dots + X_n$. Then for all $\epsilon > 0$,*

$$\mathbf{P} \left\{ \left| \frac{S_n}{n} - \mathbf{E}X_1 \right| \geq \epsilon \right\} \rightarrow 0,$$

as $n \rightarrow \infty$.

Proof. For fixed $N > 0$, we define $X_k^{\leq N}$ and $X_k^{> N}$ as follows: $X_k^{\leq N} = X_k \mathbf{1}_{|X_k| \leq N}$ and $X_k^{> N} = X_k - X_k^{\leq N}$.

We then have that $|X_1^{\leq N}|$ increases to $|X_1|$ as $N \rightarrow \infty$, so by monotone convergence

$$\mathbf{E}|X_1^{\leq N}| \rightarrow \mathbf{E}|X_1|,$$

again as $N \rightarrow \infty$. Since $|X_1| = |X_1^{\leq N}| + |X_1^{> N}|$ (check if it isn't obvious to you), it follows that as $N \rightarrow \infty$ we also have

$$\mathbf{E}|X_1^{> N}| = \mathbf{E}|X_1| - \mathbf{E}|X_1^{\leq N}| \rightarrow 0.$$

Now fix $\epsilon > 0$, and let N be large enough that $\mathbf{E}|X_1^{> N}| < \epsilon^2/8$. By Chebyshev's inequality for sums, we then have

$$\mathbf{P} \{ |\bar{S}_n^{\leq N} - \mathbf{E}\bar{S}_n^{\leq N}| > \epsilon/2 \} \leq \frac{1}{(\epsilon/2)^2 n} \mathbf{Var}(X_1^{\leq N}) \leq \frac{4N^2}{\epsilon^2 n},$$

the last inequality since $-N \leq X_1^{\leq N} \leq N$ so $\mathbf{Var}\{X_1^{\leq N}\} \leq (2N)^2/4 = N^2$. The last expression is less than $\epsilon/2$ for $n > 8N^2/\epsilon^3$. We then have

$$\begin{aligned} \mathbf{P} \{ |\bar{S}_n^{> N} - \mathbf{E}\bar{S}_n^{> N}| > \epsilon/2 \} &\leq \frac{\mathbf{E} [|\bar{S}_n^{> N} - \mathbf{E}\bar{S}_n^{> N}|]}{(\epsilon/2)} && \text{(Markov's inequality)} \\ &\leq \frac{\mathbf{E} [|\bar{S}_n^{> N}|] + |\mathbf{E}\bar{S}_n^{> N}|}{(\epsilon/2)} && \text{(Triangle inequality)} \\ &\leq \frac{4\mathbf{E} [|\bar{S}_n^{> N}|]}{\epsilon} && \text{(Move absolute value inside expectation)} \\ &\leq \frac{\epsilon}{2} && \text{(Since } \mathbf{E}|\bar{S}_n^{> N}| \leq \mathbf{E}|X_1^{> N}| < \epsilon^2/8 \text{).} \end{aligned}$$

It follows that for $n > 8N^2/\epsilon^3$,

$$\mathbf{P} \{ |\bar{S}_n - \mathbf{E}[X_1]| > \epsilon \} \leq \mathbf{P} \{ |\bar{S}_n^{\leq N} - \mathbf{E}\bar{S}_n^{\leq N}| > \epsilon/2 \} + \mathbf{P} \{ |\bar{S}_n^{> N} - \mathbf{E}\bar{S}_n^{> N}| > \epsilon/2 \} < 2\epsilon.$$

Since $\epsilon > 0$ was arbitrary this shows convergence in probability. \square

This argument was straightforward enough that it's worth seeing if we can squeeze a little more out of it. Our goal is to end up proving a *strong* law of large numbers; we want to prove that

$$\frac{S_n}{n} \xrightarrow{\text{a.s.}} \mathbf{E}X_1,$$

strengthening the convergence in probability shown above. How might we naturally proceed?

Well, we did see one way to deduce almost sure convergence from convergence in probability: Proposition 5.8, which states that if $(Z_n, 1 \leq n \leq \infty)$ are a sequence of random variables with $Z_n \xrightarrow{P} Z_\infty$, then there exists a subsequence $(n_k, k \geq 1)$ such that $Z_{n_k} \xrightarrow{\text{a.s.}} Z_\infty$ as $k \rightarrow \infty$. It is reasonable to ask how “dense” a subsequence we can pick, without working too hard, and obtain a subsequential strong law of large numbers using nothing more than the bounds we proved in the course of proving the weak law.

We say a sequence $(n_k, k \geq 1)$ is *lacunary* if it is increasing and there exists $c > 1$ such that $n_{k+1} > cn_k$ for all k sufficiently large. We will prove the following theorem.

Theorem 10.7 (Lacunary Strong Law of Large Numbers). *Let $(X_n, n \geq 1)$ be independent identically distributed random variables with $\mathbf{E}|X_n| < \infty$, and write $S_n = X_1 + \dots + X_n$. Then for any lacunary sequence of positive integers $(n_k, k \geq 1)$,*

$$\mathbb{P} \left(\lim_{k \rightarrow \infty} \frac{S_{n_k}}{n_k} = \mathbf{E}X_1 \right) = 1.$$

Proof. Let $(n_k, k \geq 1)$ be a lacunary sequence, and let $c > 1$ be such that $n_{k+1} \geq cn_k$ for $k \geq k_0$. As before, for any $\epsilon > 0$, $N > 0$ and $n \geq 1$ we have

$$\mathbf{P} \{ |\bar{S}_n^{\leq N} - \mathbf{E}\bar{S}_n^{\leq N}| > \epsilon/2 \} \leq \frac{1}{(\epsilon/2)^2 n} \mathbf{E} \left[(X_1^{\leq N})^2 \right] \leq \frac{4N^2}{\epsilon^2 n}.$$

We could have used $\mathbf{Var} \left(X_1^{\leq N} \right)$ rather than $\mathbf{E} \left[(X_1^{\leq N})^2 \right]$ above, and obtained a better upper bound; but using $\mathbf{E} \left[(X_1^{\leq N})^2 \right]$ will make things a little cleaner later. From the preceding bound it follows that

$$\begin{aligned} \sum_{k \geq 1} \mathbf{P} \{ |\bar{S}_{n_k}^{\leq N} - \mathbf{E}\bar{S}_{n_k}^{\leq N}| > \epsilon/2 \} &\leq k_0 + \sum_{k > k_0} \mathbf{P} \{ |\bar{S}_{n_k}^{\leq N} - \mathbf{E}\bar{S}_{n_k}^{\leq N}| > \epsilon/2 \} \\ &\leq k_0 + \sum_{k > k_0} \frac{\mathbf{E} \left[(X_1^{\leq N})^2 \right]}{(\epsilon/2)^2 n_k} \\ &\leq k_0 + \sum_{k > k_0} \frac{4N^2}{\epsilon^2 c^{k-k_0} n_{k_0}} \\ &< \infty, \end{aligned} \tag{10.4}$$

the last bound holding since the summands of the final sum are geometrically decreasing. It follows by the first Borel-Cantelli lemma that with probability 1, for all k sufficiently large, $|\bar{S}_n^{\leq N} - \mathbf{E}\bar{S}_n^{\leq N}| \leq \epsilon/2$.

That worked out well. But the situation isn't so good when we turn to the unbounded summands. There we have the bound

$$\mathbf{P} \{ |\bar{S}_n^{> N} - \mathbf{E}\bar{S}_n^{> N}| > \epsilon/2 \} \leq \frac{4\mathbf{E}|\bar{S}_n^{> N}|}{\epsilon}.$$

We can make $\mathbf{E} \left[|\bar{S}_n^{> N}| \right]$ as small as we like by choosing N large, but there are infinitely many summands, so a fixed positive bound on $\mathbf{E} \left[|\bar{S}_n^{> N}| \right]$ isn't good enough to let us apply the Borel-Cantelli lemma. Can we find a more explicit/more useful bound? Well, the triangle inequality is a reasonable attack:

$$|\bar{S}_n^{> N}| = \frac{1}{n} |X_1^{> N} + \dots + X_n^{> N}| \leq \frac{1}{n} |X_1^{> N}| + \dots + |X_n^{> N}|,$$

Say something about the fact that IIDness is key here?

and the summands on the right are IID, so this gives

$$\mathbf{P} \{ |\bar{S}_n^{>N} - \mathbf{E}\bar{S}_n^{>N}| > \epsilon/2 \} \leq \frac{4\mathbf{E}|X_1^{>N}|}{\epsilon}. \quad (10.5)$$

We can make $\mathbf{E}|X_1^{>N}|$ small by taking N large. But we can't make it zero, so this doesn't give a finite bound on summation.

Are we stuck? Well, if a fixed positive bound isn't good enough, maybe we can let N vary. Let's look back at the control we achieved for the bounded summands. If instead of picking a single value N we choose a sequence of different values $(N_k, k \geq 1)$, then we obtain

$$\sum_{k \geq 1} \mathbf{P} \{ |\bar{S}_{n_k}^{\leq N_k} - \mathbf{E}\bar{S}_{n_k}^{\leq N_k}| > \epsilon/2 \} \leq k_0 + \sum_{k > k_0} \frac{\mathbf{E} \left[(X_1^{\leq N_k})^2 \right]}{(\epsilon/2)^2 n_k}.$$

If we can show that the last sum is finite, then by the first Borel-Cantelli lemma we again obtain that almost surely $|\bar{S}_{n_k}^{\leq N_k} - \mathbf{E}\bar{S}_{n_k}^{\leq N_k}| \leq \epsilon/2$ except for finitely many values of k . We'd like to make this argument work for a sequence $(N_k, k \geq 1)$ growing as quickly as possible, since the larger the values N_k the easier our work will be when we turn to controlling the unbounded summands.

At this point the first natural thing to do is to use the bound $\mathbf{E} \left[(X_1^{\leq N_k})^2 \right] \leq N_k^2$. If we do that then we will end up with a finite bound provided we choose N_k such that (N_k^2/n_k) is summable. For example, taking $N_k = n_k^{1/4}$ would yield the bound

$$k_0 + \sum_{k > k_0} \frac{1}{(\epsilon^2/2)^2 n_k^{1/2}} \leq k_0 + \sum_{k > k_0} \frac{4}{\epsilon^2 c^{(k-k_0)/2} n_{k_0}^{1/2}} < \infty.$$

This already looks promising. But we can squeeze out a slightly stronger result, and simultaneously simplify our notation, by explicitly considering which values of k have $X_1^{\leq N_k} \neq 0$. That is, let $J = J(\omega) = \min\{k : N_k \geq |X_1(\omega)|\}$. Then $X_1^{\leq N_k} = 0$ for $k < J$, so

$$\sum_{k=k_0}^{\infty} \frac{\mathbf{E} \left[(X_1^{\leq N_k})^2 \right]}{n_k} = \mathbf{E} \left[\sum_{k=k_0}^{\infty} \frac{X_1^2 \mathbf{1}_{|X_1| \leq N_k}}{n_k} \right] = \mathbf{E} \left[\sum_{k=\max(k_0, J)}^{\infty} \frac{X_1^2}{n_k} \right] \leq \frac{c}{c-1} \mathbf{E} \left[\frac{X_1^2}{n_{\max(k_0, J)}} \right].$$

In the last bound we used that $\frac{n_{k+1}}{n_k} \geq c$ for $k > k_0$, so $\sum_{k=\max(k_0, J)}^{\infty} n_k^{-1} \leq n_{\max(k_0, J)}^{-1} \sum_{i \geq 0} c^{-i}$.

How can we make this bound finite? Well, if $N_k \leq n_k$ then by the definition of J we have $n_{\max(k_0, J)} \geq N_{\max(k_0, J)} \geq |X_1|$, so $\mathbf{E} \left[X_1^2/n_{\max(k_0, J)} \right] \leq \mathbf{E}|X_1| < \infty$. We want to choose N_k to be as large as possible, since this should make our lives easier when it comes to controlling the unbounded summands; so let's take $N_k = n_k$ henceforth.¹² Summarizing the story to date, we now have that

$$\sum_{k \geq 1} \mathbf{P} \{ |\bar{S}_{n_k}^{\leq n_k} - \mathbf{E}\bar{S}_{n_k}^{\leq n_k}| > \epsilon/2 \} \leq k_0 + \sum_{k > k_0} \frac{\mathbf{E} \left[(X_1^{\leq n_k})^2 \right]}{(\epsilon/2)^2 n_k} \leq k_0 + \frac{4c}{\epsilon^2(c-1)} \mathbf{E}|X_1| < \infty,$$

so by the first Borel-Cantelli lemma,

$$\mathbf{P} \{ |\bar{S}_{n_k}^{\leq n_k} - \mathbf{E}\bar{S}_{n_k}^{\leq n_k}| > \epsilon/2 \text{ i.o.} \} = 0. \quad (10.6)$$

We are in good shape for the sums of the bounded parts. For the unbounded summands, from (10.5) we have

$$\mathbf{P} \{ |\bar{S}_{n_k}^{>n_k} - \mathbf{E}\bar{S}_{n_k}^{>n_k}| > \epsilon/2 \} \leq \frac{4\mathbf{E} [|X_1^{>n_k}|]}{\epsilon}.$$

What happens if we sum the right-hand side over $k \geq k_0$?

¹²To make N_k even bigger we could take $N_k = An_k$ for some constant $A > 1$, but given that our proof has to work for all lacunary sequences, it's not hard to see that such a change would not make any difference to the success or failure of our argument.

Well, bad news, friends: the sum may be infinite. For example, it could be that $n_k = 2^k$ (in which case $k_0 = 1$). By linearity of expectation, we would then have

$$\begin{aligned} \sum_{k \geq k_0} \mathbf{E} [|X_1^{>n_k}|] &= \mathbf{E} \left[\sum_{k \geq 1} |X_1| \mathbf{1}_{\{|X_1| > 2^k\}} \right] = \mathbf{E} [|X_1| \lceil \log_2 \max(1, |X_1|) \rceil] \\ &\leq \mathbf{E} [|X_1| \log_2(X_1 + 1)] . \end{aligned}$$

We assumed $\mathbf{E}|X_1| < \infty$, but $\mathbf{E} [|X_1| \log(|X_1| + 1)]$ need not be, so this bound may be useless. On the other hand, it's worth recording now that if n_k is any lacunary sequence then similar logic would yield the bound

$$\sum_{k \geq k_0} \mathbf{E} [|X_1^{>n_k}|] = O(\mathbf{E} [|X_1| \log(|X_1| + 1)]),$$

so if $\mathbf{E} [|X_1| \log(|X_1| + 1)] < \infty$ then we can actually finish the proof along these lines.

At this point the situation may seem bleak. We are stuck trying to bound

$$\mathbf{P} \{ |\bar{S}_{n_k}^{>n_k} - \mathbf{E}\bar{S}_{n_k}^{>n_k}| > \epsilon/2 \} ,$$

and our tricks have all failed. But our sleeves are not yet empty. We will go back to the very basics and try to exploit subadditivity of probabilities. By this I mean the following. Since

$$\bar{S}_{n_k}^{>n_k} = \frac{1}{n_k} S_{n_k}^{>n_k} = \frac{1}{n} \sum_{i=1}^{n_k} X_{n_k}^{>n_k} ,$$

by the triangle inequality,

$$|\bar{S}_{n_k}^{>n_k} - \mathbf{E}\bar{S}_{n_k}^{>n_k}| \leq \frac{1}{n_k} \sum_{i=1}^{n_k} |X_{n_k}^{>n_k} - \mathbf{E}X_{n_k}^{>n_k}| ,$$

so if $|\bar{S}_{n_k}^{>n_k} - \mathbf{E}\bar{S}_{n_k}^{>n_k}| > \epsilon/2$ then there must be $1 \leq i \leq n_k$ such that $|X_{n_k}^{>n_k} - \mathbf{E}X_{n_k}^{>n_k}| > \epsilon/2$. We thus have

$$\begin{aligned} \mathbf{P} \{ |\bar{S}_{n_k}^{>n_k} - \mathbf{E}\bar{S}_{n_k}^{>n_k}| > \epsilon/2 \} &\leq \mathbf{P} \{ \exists i \in [n_k] : |X_{n_k}^{>n_k} - \mathbf{E}X_{n_k}^{>n_k}| > \epsilon/2 \} \\ &\leq n_k \mathbf{P} \{ |X_1^{>n_k} - \mathbf{E}X_1^{>n_k}| > \epsilon/2 \} . \end{aligned}$$

But $n_k \rightarrow \infty$ as $k \rightarrow \infty$, so $\mathbf{E}X_1^{>n_k} \rightarrow 0$, and so there is $k_1 \geq k_0$ such that $|\mathbf{E}X_1^{>n_k}| < \epsilon/2$ for $k \geq k_1$. For such k , if $|X_1^{>n_k} - \mathbf{E}X_1^{>n_k}| > \epsilon/2$ then in particular $X_1^{>n_k} \neq 0$, in which case necessarily $|X_1| > n_k$. Using this observation to bound the final probability above, we obtain

$$\begin{aligned} \mathbf{P} \{ |\bar{S}_{n_k}^{>n_k} - \mathbf{E}\bar{S}_{n_k}^{>n_k}| > \epsilon/2 \} &\leq \mathbf{P} \{ \exists i \in [n_k] : |X_{n_k}^{>n_k} - \mathbf{E}X_{n_k}^{>n_k}| > \epsilon/2 \} \\ &\leq n_k \mathbf{P} \{ |X_1| > n_k \} . \end{aligned}$$

This is excellent news. Because $(n_k, k \geq 1)$ is a lacunary sequence, this expectation is actually finite! To see this, recall that $J = \min\{k : n_k \geq |X_1|\}$; then $|X_1| > n_k$ only for $k < J$, so

$$\begin{aligned} \sum_{k=k_1}^{\infty} n_k \mathbf{P}\{|X_1| > n_k\} &= \mathbf{E} \left[\sum_{k=k_1}^{J-1} n_k \mathbf{1}_{|X_1| > n_k} \right] \\ &= \mathbf{E} \left[\sum_{k=k_1}^{J-1} n_k \right] && (\mathbf{1}_{|X_1| > n_k} = 0 \text{ for } i \geq J) \\ &\leq \mathbf{E} \left[\sum_{k=k_1}^{J-1} c^{-(J-1-k)} X_1 \right] && (\text{lacunarity}) \\ &\leq \mathbf{E} \left[\sum_{k=0}^{\infty} c^{-k} X_1 \right] \\ &= \frac{c}{c-1} \mathbf{E}[X_1] < \infty. \end{aligned}$$

Thus $\sum_{k \geq 1} \mathbf{P}\{|\bar{S}_{n_k}^{>n_k} - \mathbf{E}\bar{S}_{n_k}^{>n_k}| > \epsilon/2\} < \infty$, so again by the first Borel-Cantelli lemma,

$$\mathbf{P}\{|\bar{S}_{n_k}^{>n_k} - \mathbf{E}\bar{S}_{n_k}^{>n_k}| > \epsilon/2 \text{ i.o.}\} = 0.$$

But if $|\bar{S}_{n_k} - \mathbf{E}\bar{S}_{n_k}| > \epsilon$ infinitely often then either the bounded or unbounded partial sums must differ from their mean by at least $\epsilon/2$ infinitely often; so the theorem follows from the preceding equality and (10.6) \square

Let's summarize the situation. We proved the weak law of large numbers, stating conditions which guarantee that $S_n/n \rightarrow \mathbf{E}[X_1]$ in probability. Under these conditions, Proposition 5.8 then guarantees the existence of subsequences $(n_k, k \geq 1)$ along which $S_n/n \xrightarrow{\text{a.s.}} \mathbf{E}[X_1]$. The lacunary law of large numbers gave a quantitative strengthening of Proposition 5.8 in this setting, by showing that $(n_k, k \geq 1)$ can be taken to be any lacunary sequence.

This quantitative bound was not trivial to prove, but it was worth the effort, as the general strong law of large numbers ends up being a quite straightforward consequence. Its proof will proceed by first reducing to the case that the summands are non-negative, then using the monotonicity of the partial sums to relate convergence along lacunary subsequences to convergence of the whole sequence. For the second step, the key analytic fact is described in the following exercise.

Exercise 10.3. Let $(s_n, n \geq 0)$ be a non-decreasing sequence with $s_0 = 0$. Fix $\mu > 0, \epsilon \in (0, 1/3)$, and define a sequence by $n_k = \lceil (1 + \epsilon)^k \rceil$.

- (a) Show that for all n sufficiently large (i.e. $n \geq n_0(\epsilon)$), if $s_n \geq \mu n(1 + 3\epsilon)$ then letting k be such that $n_{k-1} < n \leq n_k$, we have $s_{n_k} \geq \mu n_k(1 + \epsilon)$.
- (b) Show that for all n sufficiently large, if $s_n \leq \mu n(1 - 3\epsilon)$ then letting k be such that $n_{k-1} < n \leq n_k$, we have $s_{n_{k-1}} \leq \mu n_{k-1}(1 - \epsilon)$.
- (c) Conclude that if $\limsup_n |s_n - \mu n|/n > 3\epsilon\mu$ then $\limsup_k |s_{n_k} - \mu n_k|/n_k > \epsilon\mu$.

Theorem 10.8. Let $(X_n, n \geq 1)$ be independent identically distributed random variables with $\mathbf{E}|X_1| < \infty$. Write $S_n = X_1 + \dots + X_n$ for $n \geq 1$. Then

$$\frac{S_n}{n} \xrightarrow{\text{a.s.}} \mathbf{E}X_1.$$

Proof. Write $X_n = X_n^+ - X_n^-$ and $S_n^+ = X_1^+ + \dots + X_n^+$ and $S_n^- = X_1^- + \dots + X_n^-$. If $\omega \in \Omega$ is such that $S_n^+(\omega)/n \rightarrow \mathbf{E}X_1^+$ and $S_n^-(\omega)/n \rightarrow \mathbf{E}X_1^-$ then

$$\frac{S_n(\omega)}{n} = \frac{S_n^+(\omega) + S_n^-(\omega)}{n} \rightarrow \mathbf{E}X_1^+ + \mathbf{E}X_1^- = \mathbf{E}X_1.$$

So we see that to prove $S_n/n \rightarrow \mathbf{E}(X_1)$ almost surely, it suffices to prove that

$$\frac{S_n^+}{n} \xrightarrow{\text{a.s.}} \mathbf{E}(X_1^+),$$

and that $S_n^-/n \xrightarrow{\text{a.s.}} \mathbf{E}[X_1^-]$. The point of this reduction is that summands $(X_n^+, n \geq 1)$ are all non-negative, and likewise for $(X_n^-, n \geq 1)$.

So we may now assume (by replacing $(X_i, i \geq 1)$ by either $(X_i^+, i \geq 1)$ or $(X_i^-, i \geq 1)$) that $\mathbf{P}\{X_1 \geq 0\} = 1$; in this case $(S_n, n \geq 1)$ is almost surely non-decreasing. Fix $\epsilon \in (0, 1/3)$ and for $k \geq 1$ let $n_k = n_k(\epsilon) := \lceil (1 + \epsilon)^k \rceil$. Then by Exercise 10.3,

$$\left\{ \omega : \limsup_{n \rightarrow \infty} \left| \frac{S_n(\omega)}{n} - \mathbf{E}X_1 \right| > 3\epsilon \right\} \subseteq \left\{ \omega : \limsup_{k \rightarrow \infty} \left| \frac{S_{n_k}(\omega)}{n_k} - \mathbf{E}X_1 \right| > \epsilon \right\}.$$

By Theorem 10.7, we have $\mathbf{P}\{\limsup_{k \rightarrow \infty} |S_{n_k}/n_k - \mathbf{E}X_1| > \epsilon\} = 0$; it follows that

$$\limsup_{n \rightarrow \infty} \mathbf{P} \left\{ \left| \frac{S_n}{n} - \mathbf{E}X_1 \right| > 3\epsilon \right\} = 0.$$

Since this holds for all $\epsilon > 0$, it follows that

$$\frac{S_n}{n} \xrightarrow{\text{a.s.}} \mathbf{E}X_1. \quad \square.$$

11. Convexity, inequalities, and L_p spaces

We begin with convexity. A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is convex if $f(px + (1-p)y) \leq pf(x) + (1-p)f(y)$ for all $x, y \in \mathbb{R}$ and $p \in [0, 1]$.

Exercise 11.1. *If $f : \mathbb{R} \rightarrow \mathbb{R}$ is convex then it is continuous, so Borel measurable.*

Theorem 11.1 (Jensen's inequality). *If X is a real random variable with $\mathbf{E}|X| < \infty$, and $f : \mathbb{R} \rightarrow \mathbb{R}$ is convex and $\mathbf{E}f(X)$ is defined, then $f(\mathbf{E}X) \leq \mathbf{E}f(X)$.*

Proof. Fix $h > 0$ and $0 < p < 1$. For any $x \in \mathbb{R}$ we have $x + hp = (1-p)x + p(x+h)$ so

$$f(x + hp) \leq (1-p)f(x) + pf(x+h),$$

which after rearrangement gives

$$\frac{f(x + hp) - f(x)}{hp} \leq \frac{f(x+h) - f(x)}{h}.$$

In other words, $(f(x+h) - f(x))/h$ is increasing in h for all $x \in \mathbb{R}$; we define

$$f'_+(x) := \lim_{h \downarrow 0} \frac{f(x+h) - f(x)}{h}.$$

Likewise, $(f(x) - f(x-h))/h$ is decreasing in h , so the limit

$$f'_-(x) := \lim_{h \downarrow 0} \frac{f(x) - f(x-h)}{h}.$$

Convexity also gives

$$f(x) - f(x-h) \leq f(x+h) - f(x),$$

from which it follows that $f'_-(x) \leq f'_+(x)$.

Now let $c := \mathbf{E}f(X)$, and fix $a \in \mathbb{R}$ with $f'_-(c) \leq a \leq f'_+(c)$. Then the line ℓ given by $\ell(x) = f(c) + a(x-c)$ has $\ell \leq f$ and $\ell(c) = f(c)$. By linearity of expectation and monotonicity,

it follows that

$$\begin{aligned} f(\mathbf{E}(X)) &= f(c) \\ &= \ell(c) \\ &= f(c) + a(\mathbf{E}X - c) \\ &= \mathbf{E}[f(c) + a(X - c)] \\ &= \mathbf{E}\ell(X) \\ &\leq \mathbf{E}f(X), \end{aligned}$$

as required. □

For X a random variable and $p \geq 0$ we write $\|X\|_p := (\mathbf{E}[|X|^p])^{1/p}$, and call $\|X\|_p$ the L_p -norm of X . Jensen's inequality immediately yields monotonicity of the L_p norms: if $0 \leq p \leq q$ then using the convexity of the function $x \mapsto x^{q/p}$,

$$\begin{aligned} \|X\|_q^q &= \mathbf{E}[|X|^q] = \mathbf{E}\left[(|X|^p)^{q/p} \right] = \lim_{n \rightarrow \infty} \mathbf{E}\left[(|X^{\leq n}|^p)^{q/p} \right] \\ &\geq \lim_{n \rightarrow \infty} (\mathbf{E}[|X^{\leq n}|^p])^{q/p} = \mathbf{E}[|X|^p]^{q/p} = \|X\|_p^q, \end{aligned}$$

which in a sense strengthens (10.1). (We had to use the monotone convergence theorem because it's possible that $\mathbf{E}[|X|^q] = \infty$, in which case Jensen's inequality doesn't apply directly.)

Given random variables $(X_n, 1 \leq n \leq \infty)$ defined over a common probability space $(\Omega, \mathcal{F}, \mathbf{P})$, we say that $X_n \rightarrow X_\infty$ in $L_p(\mathbf{P})$, and write $X_n \xrightarrow{L_p} X_\infty$, if $X_n \in L_p(\mathbf{P})$ for all $1 \leq n \leq \infty$ and $\|X_n - X_\infty\|_p \rightarrow 0$ as $n \rightarrow \infty$.

Exercise 11.2. For any $p > 0$ and any random variables $(X_n, n \geq 1), X, Y \in L_p(\mathbf{P})$, if $X_n \rightarrow X$ in $L_p(\mathbf{P})$ and $X_n \rightarrow Y$ in $L_p(\mathbf{P})$ then $X \stackrel{\text{a.s.}}{=} Y$.

The monotonicity of the L_p norms immediately implies that for $0 < q \leq p$, if $X_n \xrightarrow{L_p} X_\infty$ then $X_n \xrightarrow{L_q} X_\infty$. The next proposition states that convergence in L_p is at least as strong as convergence in probability.

Proposition 11.2. Let $(X_n, 1 \leq n \leq \infty)$ be real random variables defined on a common space. For any $p > 0$, if $X_n \xrightarrow{L_p} X_\infty$ then $X_n \xrightarrow{\mathbf{P}} X_\infty$.

Proof. If $X_n \xrightarrow{L_p} X_\infty$ then for any $\epsilon > 0$,

$$\mathbf{P}\{|X_n - X_\infty| \geq \epsilon\} = \mathbf{P}\{|X_n - X_\infty|^p \geq \epsilon^p\} \leq \frac{\mathbf{E}[|X_n - X_\infty|^p]}{\epsilon^p} \rightarrow 0,$$

as $n \rightarrow \infty$. □

The next exercise asks you to analyze examples which show that convergence in probability does not imply convergence in $L_p(\mathbf{P})$, which in turn does not imply almost sure convergence.

Exercise 11.3. Let $(B_n, n \geq 1)$ be independent random variables with B_n Bernoulli $(1/n)$ -distributed. Fix $p > 0$ and for $n \geq 1$ let $X_n = n^{1/p} B_n$. Show that $B_n \xrightarrow{L_p} 0$ but that B_n does not converge to 0 almost surely. Show further that $X_n \xrightarrow{\mathbf{P}} 0$ but that X_n does not converge to 0 in L_p .

Jensen's inequality also allows us to prove Hölder's inequality, which provides a tool for showing that a product of random variables is integrable.

Theorem 11.3 (Hölder's inequality). Fix $p, q \geq 1$ with $1 \leq p, q \leq \infty$. If $1/p + 1/q = 1$ then for any random variables X, Y defined on a common probability space,

$$\|XY\|_1 \leq \|X\|_p \|Y\|_q.$$

Proof. We may assume that $\|X\|_1 > 0$ and that $\|Y\|_1 > 0$ or else the left-hand side is zero. Similarly, we may assume that $\|X\|_p < \infty$ and that $\|Y\|_q < \infty$ or else the right-hand side is infinite. Finally, we may assume that $X \geq 0$ and $Y \geq 0$ since the values of both the left- and right-hand sides are unchanged if we replace X by $|X|$ and Y by $|Y|$.

We now write

$$\begin{aligned} \mathbf{E}|XY| &= \mathbf{E} \left[e^{\log(XY)} \right] \\ &= \mathbf{E} \left[e^{\log X + \log Y} \right] \\ &= \mathbf{E} \left[e^{\frac{1}{p} \log(X^p) + \frac{1}{q} \log(Y^q)} \right] \end{aligned}$$

Since $u \mapsto \log u$ is concave, $\frac{1}{p} \log(X^p) + \frac{1}{q} \log(Y^q) \leq \log(\frac{1}{p} X^p + \frac{1}{q} Y^q)$, so \square

The Cauchy-Schwarz inequality is the case $p = q = 2$ of Hölder's inequality.

Corollary 11.4 (Cauchy-Schwarz inequality for random variables). *For any random variables X, Y defined on a common probability space, $\|XY\|_1 \leq \|X\|_2 \|Y\|_2$.*

The next exercise asks you to verify the “ $p = 1, q = \infty$ ” case of Hölder's inequality. Given a random variable X we write $\text{ess sup } X := \sup\{c \in \mathbb{R} : \mathbf{P}\{X > c\} > 0\}$ and call $\text{ess sup } X$ the *essential supremum* of X . We write $\|X\|_\infty := \text{ess sup } |X|$, and let

$$L_\infty(\Omega, \mathcal{F}, \mathbf{P}) := \{X : \Omega \rightarrow \mathbb{R} : X \text{ is } (\mathcal{F}/\mathcal{B}(\mathbb{R}))\text{-measurable and } \|X\|_\infty < \infty\}.$$

$L_\infty(\Omega, \mathcal{F}, \mathbf{P})$.

Exercise 11.4. *Let X, Y be random variables defined on a common space. Show that $\text{ess sup } |X| = \lim_{p \rightarrow \infty} \|X\|_p$. Then show that*

$$\|XY\|_1 \leq \|X\|_\infty \|Y\|_1.$$

A clever application of Hölder's inequality yields *Minkowski's inequality*, which is the triangle inequality for L_p spaces.

Theorem 11.5 (Minkowski's inequality). *Let X, Y be random variables in $L_1(\mathbf{P})$. Then for all $p \geq 1$, $\|X + Y\|_p \leq \|X\|_p + \|Y\|_p$.*

Proof. When $p = 1$ this follows from monotonicity of expectation, since $|X + Y| \leq |X| + |Y|$ by the triangle inequality. For $p > 1$ we may assume that $X, Y \in L_p(\mathbf{P})$ since otherwise the right-hand side is infinite. We now use the triangle inequality, as follows:

$$\|X + Y\|_p^p = \mathbf{E}[|X + Y|^p] = \mathbf{E}[|X + Y||X + Y|^{p-1}] \leq \mathbf{E}[|X||X + Y|^{p-1}] + \mathbf{E}[|Y||X + Y|^{p-1}].$$

Applying Hölder's inequality to each of the above expectations gives

$$\mathbf{E}[|X||X + Y|^{p-1}] \leq (\mathbf{E}[|X|^p])^{1/p} (\mathbf{E}[|X + Y|^p])^{(p-1)/p} = \|X\|_p \|X + Y\|_p^{(p-1)/p}$$

and

$$\mathbf{E}[|Y||X + Y|^{p-1}] \leq (\mathbf{E}[|Y|^p])^{1/p} (\mathbf{E}[|X + Y|^p])^{(p-1)/p} = \|Y\|_p \|X + Y\|_p^{(p-1)/p},$$

so

$$\|X + Y\|_p^p \leq (\|X\|_p + \|Y\|_p) \|X + Y\|_p^{(p-1)/p}.$$

Dividing by $\|X + Y\|_p^{(p-1)/p}$ completes the proof. \square

We would like to think of $L_p(\Omega, \mathcal{F}, \mathbf{P})$ as a metric space, but this isn't quite right because a metric space (M, d) is supposed to satisfy that $d(x, y) = 0$ if and only if $x = y$. But $\|X - Y\|_p = 0$ provided that $X \stackrel{\text{a.s.}}{=} Y$, and almost sure equality is not the same as equality.

There are two ways to deal with this. The first approach, which is the most standard in probability, is to simply accept that instead of a metric space we only have a pseudo-metric space. (A pseudo-metric space is a pair (M, d) where $d : M \times M \rightarrow [0, \infty)$ is a symmetric function satisfying

the triangle inequality. In a pseudo-metric space it is possible to have $d(x, y) = 0$ for distinct points x, y .) The other is to quotient by almost sure equality. In other words, for $X \in L_p(\Omega, \mathcal{F}, \mathbf{P})$ we may write $[X] = \{Y \in L_p(\Omega, \mathcal{F}, \mathbf{P}) : X \stackrel{\text{a.s.}}{=} Y\}$ and $[L_p(\Omega, \mathcal{F}, \mathbf{P})] = \{[X] : X \in L_p(\Omega, \mathcal{F}, \mathbf{P})\}$.

Exercise 11.5. • Check that almost sure equality is an equivalence relation.
 • Check that if $X, Y \in L_p(\Omega, \mathcal{F}, \mathbf{P})$ and $X' \in [X]$ and $Y' \in [Y]$, then $\|X' - Y'\|_p = \|X - Y\|_p$. That is, L_p distance is a class function for the “almost sure equality” equivalence relation.

The next theorem implies that $[L_p(\Omega, \mathcal{F}, \mathbf{P})]$ is a complete metric space for any probability space $(\Omega, \mathcal{F}, \mathbf{P})$.

Theorem 11.6. Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and $p \geq 1$, and let $(X_n, n \geq 1)$ be a Cauchy sequence in $L_p(\mathbf{P})$. Then there is $X \in L_p(\mathbf{P})$ such that $X_n \xrightarrow{L_p} X$. Moreover, for any other random variable $X' : \Omega \rightarrow \mathbb{R}$, if $\|X_n - X'\|_p \rightarrow 0$ then $X' \stackrel{\text{a.s.}}{=} X$.

Proof. Since $(X_n, n \geq 1)$ is Cauchy, we can find an increasing sequence of integers $(n(k), k \geq 1)$ such that for all $m, n \in \mathbb{N}$, if $m, n \geq n(k)$ then $\|X_m - X_n\|_p \leq 2^{-k}$.

Then write $Y_k = X_{n(k)}$. By our choice of the sequence $(n(k), k \geq 1)$ we have $\|Y_{k+1} - Y_k\|_p \leq 2^{-k}$, so

$$\mathbf{E} \left[\sum_{k \geq 1} |Y_{k+1} - Y_k| \right] = \sum_{k \geq 1} \mathbf{E} |Y_{k+1} - Y_k| = \sum_{k \geq 1} \|Y_{k+1} - Y_k\|_1 \leq \sum_{k \geq 1} \|Y_{k+1} - Y_k\|_p \leq 1.$$

It follows that $\mathbf{P} \left\{ \sum_{k \geq 1} |Y_{k+1} - Y_k| < \infty \right\} = 1$, or in other words that $(Y_{k+1} - Y_k, k \geq 1)$ is almost surely absolutely convergent. Letting $Y := \limsup_{k \rightarrow \infty} Y_k$, it follows that $\mathbf{P} \{ \lim_{k \rightarrow \infty} Y_k = Y \} = 1$.

Now, for $n \geq n(k)$, note that

$$X_{n(k)} + \sum_{\ell \geq k} (Y_{\ell+1} - Y_\ell) = Y_k + \sum_{\ell \geq k} (Y_{\ell+1} - Y_\ell) \stackrel{\text{a.s.}}{=} Y,$$

so for $n \geq n(k)$,

$$\begin{aligned} \|X_n - Y\|_p &= \|X_n - X_{n(k)} - \sum_{\ell \geq k} (Y_{\ell+1} - Y_\ell)\|_p \\ &\leq \|X_n - X_{n(k)}\|_p + \sum_{\ell \geq k} \|Y_{\ell+1} - Y_\ell\|_p \\ &\leq \frac{1}{2^k} + \sum_{\ell \geq k} \frac{1}{2^\ell} = \frac{1}{2^{k-2}}. \end{aligned}$$

Thus $\|X_n - Y\|_p \rightarrow 0$ as $p \rightarrow \infty$. Since $\|Y\|_p = \|X_n - Y - X_n\|_p \leq \|X_n - Y\|_p + \|X_n\|_p$, it follows that $\|Y\|_p < \infty$, so $X_n \xrightarrow{L_p} Y$. Finally, the almost sure uniqueness of the limit is Exercise 11.2. \square

11.1. The geometric structure of L_2 . The space $L_2(\mathbf{P})$ is special because it can be endowed with a natural inner product structure, which allows us to harness the power of geometry. For $X, Y \in L^2(\mathbf{P})$, let $\langle X, Y \rangle := \mathbf{E} [XY]$; that the right-hand side is defined follows from the Cauchy-Schwarz inequality. You should check that $\langle \cdot, \cdot \rangle : L_2(\mathbf{P}) \times L_2(\mathbf{P})$ satisfies the axioms of an inner product (up to almost sure equivalence): it is symmetric and bilinear, and $\langle X, X \rangle = 0$ if and only if $X \stackrel{\text{a.s.}}{=} 0$. (The “true” inner product space is $[L_2(\Omega, \mathcal{F}, \mathbf{P})]$, but we will continue working with random variables, at the cost of occasionally having to use the phrase “almost sure”.)

If $X > 0$ and $Y > 0$ are random variables in $L_2(\mathbf{P})$ then we may use the inner product to define an angle $\theta_{XY} \in [0, \pi]$ by the formula

$$\cos \theta_{XY} = \frac{\langle X, Y \rangle}{\|X\|_2 \|Y\|_2}.$$

Note that $\cos \theta_{XX} = \mathbf{E}[X^2] / \|X\|_2^2 = 1$, so $\theta_{XX} = 0$. This geometric structure is closely related to the covariance of the random variables: recall that for $X, Y \in L_2(\mathbf{P})$,

$$\text{Cov}(X, Y) = \mathbf{E}[(X - \mathbf{E}X)(Y - \mathbf{E}Y)] = \langle X, Y \rangle - \mathbf{E}X\mathbf{E}Y.$$

In the case that $\mathbf{E}X = 0$ or $\mathbf{E}Y = 0$, it follows that X and Y are uncorrelated if and only if $\langle X, Y \rangle = 0$ or equivalently if and only if $\theta_{X,Y} = \pi/2$. We also have that

$$\|X + Y\|_2^2 = \mathbf{E}[(X + Y)^2] = \mathbf{E}[X^2] + 2\mathbf{E}[XY] + \mathbf{E}[Y^2] = \|X\|_2^2 + \langle X, Y \rangle + \|Y\|_2^2,$$

so $\|X + Y\|_2^2 = \|X\|_2^2 + \|Y\|_2^2$ if and only if $\langle X, Y \rangle = 0$.

Exercise 11.6 (Parallelogram Law). For $U, V \in L_2(\mathbf{P})$ we have

$$\|U + V\|_2^2 + \|U - V\|_2^2 = 2\|U\|_2^2 + 2\|V\|_2^2.$$

Covariance has a very direct relation to the geometric structure of $L_2(\Omega, \mathcal{F}, \mathbf{P})$. Another feature of the geometry which we will exploit is the ability to perform *projections* onto subspaces. Consider a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and a sub- σ -field \mathcal{G} of \mathcal{F} . If $X : \Omega \rightarrow \mathbb{R}$ is $(\mathcal{G}/\mathcal{B}(\mathbb{R}))$ -measurable then it is $(\mathcal{F}/\mathcal{B}(\mathbb{R}))$ -measurable, so for any $p \geq 1$, if $Z \in L_p(\Omega, \mathcal{G}, \mathbf{P})$ then $Z \in L_p(\Omega, \mathcal{F}, \mathbf{P})$. In other words, “up to almost sure equality” the space $L_p(\Omega, \mathcal{G}, \mathbf{P})$ is a complete subspace of $L_p(\Omega, \mathcal{F}, \mathbf{P})$. In the case $p = 2$, the existence of a notion of orthogonality then allows us to consider projections onto $L_2(\Omega, \mathcal{G}, \mathbf{P})$.

picture?

We are abusing notation by writing $(\Omega, \mathcal{G}, \mathbf{P})$ rather than $(\Omega, \mathcal{G}, \mathbf{P}|_{\mathcal{G}})$, but this shouldn't cause confusion.

Theorem 11.7. Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and a sub- σ -field \mathcal{G} of \mathcal{F} . Fix $X \in L_2(\Omega, \mathcal{F}, \mathbf{P})$ and let $\Delta = \inf\{\|X - Y\|_2 : Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})\}$. Then there is an almost surely unique $Z \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ such that $\|X - Z\|_2 = \Delta$.

Proof. Let $(Y_n, n \geq 1)$ be elements of $L_2(\Omega, \mathcal{G}, \mathbf{P})$ with $\|X - Y_n\|_2 \leq \Delta + 1/n$. For $m, n \geq 1$, we apply the parallelogram law with $U + V = X - Y_n$ and $U - V = X - Y_m$. This means $2U = 2X - (Y_n + Y_m)$ and $2V = Y_m - Y_n$, so we obtain

$$\|X - Y_n\|_2^2 + \|X - Y_m\|_2^2 = 2\|X - (Y_n + Y_m)/2\|_2^2 + \frac{1}{2}\|Y_m - Y_n\|_2^2.$$

The left-hand side is at most $2\Delta^2 + 1/m + 1/n$ by our choice of Y_m and Y_n . Also, $(Y_n + Y_m)/2 \in L_2(\Omega, \mathcal{G}, \mathbf{P})$, so by the definition of Δ we have $\|X - (Y_n + Y_m)/2\|_2^2 \geq \Delta^2$. From the above equality combined with these two bounds we obtain

$$\frac{1}{2}\|Y_m - Y_n\|_2^2 \leq \frac{1}{m} + \frac{1}{n},$$

so $(Y_n, n \geq 1)$ is a Cauchy sequence. By Theorem 11.6, it follows that there is $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ such that $\|Y_n - Y\|_2 \rightarrow 0$. For any $n \geq 1$, by the triangle inequality, we then have

$$\|X - Y\|_2 \leq \|X - Y_n\|_2 + \|Y_n - Y\|_2 \leq \Delta + \frac{1}{n} + \|Y_n - Y\|_2,$$

and taking $n \rightarrow \infty$ shows that $\|X - Y\|_2 \leq \Delta$; by the definition of Δ we must then have $\|X - Y\|_2 = \Delta$.

Finally, suppose Z is another random variable with $\|X - Z\|_2 = \Delta$. Then apply the parallelogram law with $U + V = X - Y$ and $U - V = X - Z$. We then obtain

$$2\Delta^2 = \|X - Y\|_2^2 + \|X - Z\|_2^2 = 2\|X - (Y + Z)/2\|_2^2 + \frac{1}{2}\|Y - Z\|_2^2 \geq 2\Delta^2 + \frac{1}{2}\|Y - Z\|_2^2,$$

so it must be that $\|Y - Z\|_2 = 0$. \square

Let's call the (a.s. unique) minimizer Z in the above theorem the *closest \mathcal{G} -measurable random variable to X* . (This is cumbersome - we'll introduce a shorter name shortly.)

Corollary 11.8. *With the setup of Theorem 11.7, for $Z \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ we have $\|X - Z\|_2 = \Delta$ if and only if $\langle Y, X - Z \rangle = 0$ for all $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$.*

Proof. First, suppose that $Z \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ is such that $\langle Y, X - Z \rangle = 0$ for all $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$. Then for any $Z' \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ we have

$$\begin{aligned} \mathbf{E}[(X - Z')^2] &= \mathbf{E}[(X - Z - (Z - Z'))^2] \\ &= \mathbf{E}[(X - Z)^2] - 2\langle X - Z, Z - Z' \rangle + \mathbf{E}[(Z - Z')^2] \\ &= \mathbf{E}[(X - Z)^2] + \mathbf{E}[(Z - Z')^2] \\ &\geq \mathbf{E}[(X - Z)^2]. \end{aligned}$$

Conversely, suppose that $\|X - Z\|_2 = \Delta$, and fix any $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$. Then for any $t \in \mathbb{R}$ we have $Z + tY \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ so

$$\begin{aligned} \Delta^2 &\leq \mathbf{E}[(X - Z - tY)^2] \\ &= \mathbf{E}[(X - Z)^2] - 2t\mathbf{E}[Y(X - Z)] + t^2\mathbf{E}[Y^2] \\ &= \Delta^2 - 2t\langle Y, X - Z \rangle + t^2\mathbf{E}[Y^2]. \end{aligned}$$

The only way this can hold for small t is if $\langle Y, X - Z \rangle = 0$. □

The condition that $\langle Y, X - Z \rangle = 0$ for all $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ may be rewritten as saying that

$$\mathbf{E}[Y(X - Z)] = 0$$

for all $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$. This is easy enough to verify that we can start doing examples.

Examples. The first **several** examples will relate to a sequence $(X_i, i \geq 1)$ of independent random variables in $L_2(\Omega, \mathcal{F}, \mathbf{P})$.

1. Suppose that $\mathbf{E}[X_i] = 0$ for all i . Fix $n \geq 1$, let $X = \sum_{i=1}^n X_i$, and let $\mathcal{G} = \sigma(X_1, \dots, X_{n-1})$. We claim that $Z = X_1 + \dots + X_{n-1}$ is the closest \mathcal{G} -measurable random variable to X . To see this, we fix any $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ and compute

$$\mathbf{E}[Y(X - Z)] = \mathbf{E}[YX_n] = \mathbf{E}[Y] \mathbf{E}[X_n] = 0.$$

The second equality holds since $\mathcal{G} = \sigma(X_1, \dots, X_{n-1})$ and $\sigma(X_n)$ are independent, so Y and X_n are independent.

2. Again take $X = \sum_{i=1}^n X_i$, but this time don't assume the random variables have zero mean. Write $\mathbf{E}[X_i] = c_i$, fix some set $S \subset [n]$ and let $\mathcal{G} = \sigma(X_i, i \in S)$. If $Z_0 = \sum_{i \in S} X_i$ then for $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ we have

$$\mathbf{E}[Y(X - Z_0)] = \mathbf{E}\left[Y\left(\sum_{i \notin S} X_i\right)\right] = \mathbf{E}[Y] \mathbf{E}\left[\sum_{i \notin S} X_i\right] = \mathbf{E}[Y] \cdot \sum_{i \notin S} c_i.$$

This need not be zero - we should *recenter* Z_0 to take account of what direction the remaining summands are heading in. Taking $Z = Z_0 + \sum_{i \notin S} c_i$, we then get

$$\mathbf{E}[Y(X - Z)] = \mathbf{E}[Y(X - Z_0)] - \mathbf{E}\left[Y \cdot \sum_{i \notin S} c_i\right] = 0,$$

so the closest \mathcal{G} -measurable random variable to X is $\sum_{i \in S} X_i + \sum_{i \notin S} c_i$.

3. Let $X = \prod_{i=1}^n X_i$ and take $\mathcal{G} = \sigma(X_1, \dots, X_{n-1})$. Then with $c = \mathbf{E}X_n$, the closest \mathcal{G} -measurable random variable to X is $Z = c \cdot \prod_{i=1}^{n-1} X_i$. To see this, choose $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$. Since the random variables X_1, \dots, X_n are independent, both X and Z are in $L_2(\Omega, \mathcal{F}, \mathbf{P})$ (**exercise!**).

It follows by Cauchy-Schwarz that $YZ \in L_1(\Omega, \mathcal{F}, \mathbf{P})$; since X_n is independent of YZ and $YX = YZX_n/c$, we thus have

$$\mathbf{E}[YX] = \mathbf{E}[YZX_n/c] = \mathbf{E}[YZ] \mathbf{E}[X_n]/c = \mathbf{E}[YZ].$$

4. Fix $c \in \mathbb{R}$ and suppose that $\mathbf{E}X_i = c$ and (to avoid integrability issues) that $X_i \geq 0$ for all $i \geq 1$. Let N be a positive integer random variable independent of $(X_i, i \geq 1)$, and take $X = \sum_{i=1}^N X_i$ and $\mathcal{G} = \sigma(N)$. We claim that $Z = cN$. To see this, we transform the random sum into a deterministic sum by writing

$$X = \sum_{i=1}^N X_i = \sum_{n=1}^{\infty} (\mathbf{1}_{[N=n]} \cdot \sum_{i=1}^n X_i).$$

For $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ we then have

$$\begin{aligned} \mathbf{E}[YX] &= \mathbf{E} \left[Y \cdot \sum_{n=1}^{\infty} (\mathbf{1}_{[N=n]} \cdot \sum_{i=1}^n X_i) \right] \\ &= \sum_{n=1}^{\infty} \mathbf{E} \left[Y \mathbf{1}_{[N=n]} \cdot \sum_{i=1}^n X_i \right]. \end{aligned}$$

Now, $\mathcal{G} = \sigma(N)$ and $\sigma(X_i, i \geq 1)$ are independent, so the random variables $Y \mathbf{1}_{[N=n]}$ and $\sum_{i=1}^n X_i$ are independent. Applying the factorization formula to the right-hand side above then gives

$$\begin{aligned} \mathbf{E}[YX] &= \sum_{n=1}^{\infty} \mathbf{E} [Y \mathbf{1}_{[N=n]}] \cdot \mathbf{E} \left[\sum_{i=1}^n X_i \right] \\ &= \sum_{n=1}^{\infty} \mathbf{E} [Y \mathbf{1}_{[N=n]}] \cdot cn \\ &= \mathbf{E} \left[\sum_{n=1}^{\infty} Y \mathbf{1}_{[N=n]} \cdot cn \right] \\ &= \mathbf{E} [Y \cdot cN]. \end{aligned}$$

5. This example is chattier. The idea behind it is a bit different from the others, and is quite important. Let Ω be the set of all individuals who filed an income tax return in Canada in 2018, and let \mathbf{P} be the uniform measure on $(\Omega, 2^\Omega)$. Define a random variable $X : \Omega \rightarrow \mathbb{R}$ by taking $X(\omega)$ to be the amount of income tax paid by individual ω .

Define another random variable $R : \Omega \rightarrow \{1, \dots, 13\}$ by taking $R(\omega)$ to be the province or territory of residence of individual ω in 2018¹³, and let $\mathcal{G} = \sigma(R)$. This means that (for example) $\Omega_1 := R^{-1}(1)$ is the set of taxpayers in Alberta, and $\Omega_{13} := R^{-1}(13)$ is the set of taxpayers in the Yukon.

Note that \mathcal{G} is generated by the sets $\Omega_1, \dots, \Omega_{13}$. That means a random variable $U : \Omega \rightarrow \mathbb{R}$ is $(\mathcal{G}/\mathcal{B}(\mathbf{R}))$ -measurable if and only if for any $B \in \mathcal{B}(\mathbf{R})$, the set $U^{-1}(B)$ is a union of some or all of the sets $\Omega_1, \dots, \Omega_{13}$. In other words, whether $U(\omega) \in B$ must only depend on the value of $R(\omega)$, so $U(\omega)$ must be the same for every taxpayer in a given province.

What is the closest \mathcal{G} -measurable random variable to X ? We seek a random variable Z which assigns the same value to every taxpayer in a province, and satisfies

$$\mathbf{E}[XY] = \mathbf{E}[ZY]$$

¹³Order the provinces and territories in some way..

for any other random variable Y which also assigns the same value to every taxpayer in a province. Suppose Y has that property, so we may represent Y as $Y(i) = \sum_{i=1}^{13} y_i \mathbf{1}_{[\omega \in \Omega_i]}$. Then

$$\mathbf{E}[XY] = \frac{1}{|\Omega|} \sum_{\omega \in \Omega} X(\omega)Y(\omega) = \frac{1}{|\Omega|} \sum_{i=1}^{13} y_i \cdot \sum_{\omega \in \Omega_i} X(\omega).$$

If $Z = \sum_{i=1}^{13} z_i \mathbf{1}_{[\omega \in \Omega_i]}$ then

$$\mathbf{E}[ZY] = \frac{1}{|\Omega|} \sum_{\omega \in \Omega} Z(\omega)Y(\omega) = \frac{1}{|\Omega|} \sum_{i=1}^{13} |\Omega_i| y_i z_i.$$

To make the last two expressions equal, we see that we should take $z_i = |\Omega_i|^{-1} \sum_{\omega \in \Omega_i} X_i(\omega)$. This last value is just the average tax paid by taxpayers in province/territory i ! Calling that value μ_i , we then have

$$Z(\omega) = \sum_{i=1}^{13} \mu_i \mathbf{1}_{[\Omega_i]}(\omega).$$

It's worth comparing this example to the previous ones. In examples 1, 2 and 3, the closest \mathcal{G} -measurable random variable to X ended up being obtained by essentially “replacing the part not lying in the subspace by its expected value”. In example 4, the “expectation of the independent part” also came into the picture, but in a more involved way, as X did not separate as cleanly as in the first three cases. In example 5, the random variable Z is a sort of “coarsening” of X , obtained by taking expectations over subsets whenever \mathcal{G} gives us no information about the behaviour of X within those subsets. If you think of X as a lookup table where the first column lists taxpayers and the second lists the amount they paid, then Z is a table which only lists the average (expected) amount paid per province or territory.

This motivates the name we will use for such random variables for the rest of the class; rather than calling Z the closest \mathcal{G} -measurable random variable to X , we call it the *conditional expectation of X given \mathcal{G}* , and denote it $\mathbf{E}[X | \mathcal{G}]$.

This notation takes some time to get used to. The conditional expectation $\mathbf{E}[X | \mathcal{G}]$ is not a *number*: it is a random variable, which “tries to be like X ” but is forced to be simpler than X by the constraint that it must be $(\mathcal{G}/\mathcal{B}(\mathbb{R}))$ -measurable. The next section is devoted to conditional expectation and its properties.

12. Conditional expectation

Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space and let $\mathcal{G} \subset \mathcal{F}$ be a sub- σ -algebra. For a random variable $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$, we say that $Z : \Omega \rightarrow \mathbb{R}$ is a version of $\mathbf{E}[X | \mathcal{G}]$ if

- (a) $Z \in L_1(\Omega, \mathcal{G}, \mathbf{P})$, and
- (b) For all $E \in \mathcal{G}$, $\mathbf{E}[X \mathbf{1}_{[E]}] = \mathbf{E}[Z \mathbf{1}_{[E]}]$.

We will momentarily show existence and (almost sure) uniqueness of conditional expectations of L_1 random variables. First, however, we establish a monotonicity property of conditional expectations. We use the result of the following easy exercise.

Exercise 12.1. Suppose that random variables U and V have $\mathbf{P}\{U \geq V\} = 1$ and $\mathbf{E}U \leq \mathbf{E}V$. Then $U \stackrel{\text{a.s.}}{=} V$.

Proposition 12.1 (Monotonicity of conditional expectation). Suppose that $X, X' \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ satisfy $\mathbf{P}\{X \leq X'\} = 1$, and that Z, Z' are versions of $\mathbf{E}[X | \mathcal{G}]$ and $\mathbf{E}[X' | \mathcal{G}]$, respectively. Then $\mathbf{P}\{Z \leq Z'\} = 1$.

Proof. Since $Z, Z' \in L_1(\Omega, \mathcal{G}, \mathbf{P})$ we have $Z - Z' \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ so $\{Z \geq Z'\} = \{Z - Z' \geq 0\} \in \mathcal{G}$. Thus, by the defining property (b) of conditional expectation and monotonicity of expectation,

$$\mathbf{E}[Z \mathbf{1}_{[Z \geq Z']}] = \mathbf{E}[X \mathbf{1}_{[Z \geq Z']}] \leq \mathbf{E}[X' \mathbf{1}_{[Z \geq Z']}] = \mathbf{E}[Z' \mathbf{1}_{[Z \geq Z']}] .$$

Did I already state this exercise or a close variant earlier?

But $Z\mathbf{1}_{[Z \geq Z']} \geq Z'\mathbf{1}_{[Z \geq Z']}$, so it follows by the above exercise that $\mathbf{P}\{Z\mathbf{1}_{[Z \geq Z']} = Z'\mathbf{1}_{[Z \geq Z]}\} = 1$, which is equivalent to the assertion that $\mathbf{P}\{Z \leq Z'\} = 1$. \square

Theorem 12.2 (Existence of conditional expectation). *For any random variable $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ and any sub- σ -algebra $\mathcal{G} \subset \mathcal{F}$, there exists a version of $\mathbf{E}[X | \mathcal{G}]$. Moreover, if Y and Y' are two versions of $\mathbf{E}[X | \mathcal{G}]$ then $Y \stackrel{\text{a.s.}}{=} Y'$.*

Proof. We first prove the uniqueness claim. Suppose that Z, Z' are two versions of $\mathbf{E}[X | \mathcal{G}]$. Applying Proposition 12.1 with $X' = X$ we have $\mathbf{P}\{Z \leq Z'\} = 1$; by symmetry we then have $\mathbf{P}\{Z = Z'\} = 1$, establishing the uniqueness claimed in the theorem statement.

To prove existence, first suppose $X \in L_2(\Omega, \mathcal{F}, \mathbf{P})$. Then by Corollary 11.8, there is $Z \in L_2(\Omega, \mathcal{G}, \mathbf{P})$ such that

$$\mathbf{E}[XY] = \mathbf{E}[ZY]$$

for all $Y \in L_2(\Omega, \mathcal{G}, \mathbf{P})$. In particular this holds when $Y = \mathbf{1}_{[E]}$ for $E \in \mathcal{G}$, so Z is a version of $\mathbf{E}[X | \mathcal{G}]$.

Now suppose $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ is non-negative. Then for each $n \geq 1$, since $X^{\leq n}$ is bounded it is in $L_2(\Omega, \mathcal{F}, \mathbf{P})$, so we may find a version Z_n of $\mathbf{E}[X^{\leq n} | \mathcal{G}]$. By monotonicity of conditional expectation, the random variables $(Z_n, n \geq 1)$ are almost surely increasing. Set $Z = \limsup_{n \rightarrow \infty} Z_n$, so that Z_n almost surely increases to Z . Then for any event $E \in \mathcal{G}$, by two applications of the monotone convergence theorem we then have

$$\mathbf{E}[X\mathbf{1}_{[E]}] = \lim_{n \rightarrow \infty} \mathbf{E}[X^{\leq n}\mathbf{1}_{[E]}] = \lim_{n \rightarrow \infty} \mathbf{E}[Z_n\mathbf{1}_{[E]}] = \mathbf{E}[Z\mathbf{1}_{[E]}],$$

so Z is a version of $\mathbf{E}[X | \mathcal{G}]$.

Finally, for arbitrary $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ we may write $X = X^+ - X^-$ and let Z_+ and Z_- be versions of $\mathbf{E}[X^+ | \mathcal{G}]$ and $\mathbf{E}[X^- | \mathcal{G}]$, respectively. Then for $E \in \mathcal{G}$, using linearity of expectation we have

$$\mathbf{E}[X\mathbf{1}_{[E]}] = \mathbf{E}[X^+\mathbf{1}_{[E]}] - \mathbf{E}[X^-\mathbf{1}_{[E]}] = \mathbf{E}[Z_+\mathbf{1}_{[E]}] - \mathbf{E}[Z_-\mathbf{1}_{[E]}] = \mathbf{E}[(Z_+ - Z_-)\mathbf{1}_{[E]}],$$

so $Z_+ - Z_-$ is a version of $\mathbf{E}[X | \mathcal{G}]$. \square

It is immediate from the definition that in the five examples with which we concluded the preceding section, the “closest \mathcal{G} -measurable random variables to X ” were in fact versions of $\mathbf{E}[X | \mathcal{G}]$.

Exercise 12.2. *Use the monotone class theorem to show that if Z is a version of $\mathbf{E}[X | \mathcal{G}]$ then for any $Y \in L_\infty(\Omega, \mathcal{G}, \mathbf{P})$, $\mathbf{E}[XY] = \mathbf{E}[ZY]$.*

We’ll sometimes start writing $\mathbf{E}[X | \mathcal{G}]$ rather than referring to versions of $\mathbf{E}[X | \mathcal{G}]$. Also, if $\mathcal{G} = \sigma(V)$ for some random variable V , it’s standard to write $\mathbf{E}[X | V]$ rather than $\mathbf{E}[X | \mathcal{G}]$ or $\mathbf{E}[X | \sigma(V)]$.

More examples.

1. Our first example generalizes the last example from the previous section. Suppose $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$. Let $(\Omega_n, n \geq 1)$ be a partition of Ω with all parts in \mathcal{F} , and let $\mathcal{G} = \sigma(\{\Omega_n, n \geq 1\})$. Write $z_n = \mathbf{E}[X\mathbf{1}_{[\Omega_n]}] / \mathbf{P}\{\Omega_n\}$. Then the random variable $Z = \sum_{n \geq 1} z_n \mathbf{1}_{[\Omega_n]}$ is a version of $\mathbf{E}[X | \mathcal{G}]$. To see this is easy: any event E in \mathcal{G} may be written as

$$E = \sum_{i \in S} \Omega_i$$

for some $S \subset \mathbb{N}$, and then

$$\mathbf{E}[X\mathbf{1}_{[E]}] = \sum_{n \in S} \mathbf{E}[X\mathbf{1}_{[\Omega_n]}] = \sum_{i \in S} z_n \mathbf{P}\{\Omega_n\} = \sum_{n \in S} z_n \mathbf{E}[\mathbf{1}_{[\Omega_n]}] = \mathbf{E}\left[\sum_{n \in S} Z\mathbf{1}_{[\Omega_n]}\right] = \mathbf{E}[Z\mathbf{1}_{[E]}]$$

2. Say that $X, Y \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ have *joint density* f if $f : \mathbb{R}^2 \rightarrow [0, \infty)$ is a Borel function which is a density for the \mathbb{R}^2 -valued random variable (X, Y) ; that is, for any $B \in \mathcal{B}(\mathbb{R}^2)$,

$$\mathbf{P} \{(X, Y) \in B\} = \int_B f(x, y) dx \otimes dy,$$

where we use $dx \otimes dy$ to denote Lebesgue measure on \mathbb{R}^2 . Suppose X and Y have joint density f .

The “natural formula” for $\mathbf{E}[X | Y = y]$ would be

$$\mathbf{E}[X | Y = y] = \frac{\int_{\mathbb{R}} x f(x, y) dx}{\int_{\mathbb{R}} f(x, y) dx}.$$

The top is just the integral along the slice, and the bottom is a normalization factor. If we

Now, cast your mind back to the development of product measure and Fubini’s theorem. Lemma 8.5 tells us that $f(x, y)$ is a Borel function of x ; Lemma 8.6 tells us that

$$\int_{\mathbb{R}} f(x, y) dx$$

is a Borel function of y , and Fubini’s theorem tells us that

$$\int_{\mathbb{R}} |f(x, y)| dx < \infty$$

almost everywhere. We can thus define

$$\phi(y) = \begin{cases} \int_{\mathbb{R}} f(x, y) dx & \text{if } \int_{\mathbb{R}} |f(x, y)| dx < \infty \\ 0 & \text{otherwise,} \end{cases}$$

and by Fubini’s theorem, for $A \in \mathcal{B}(\mathbb{R})$ we have

$$\mathbf{P} \{Y \in A\} = \mathbf{P} \{(X, Y) \in \mathbb{R} \times A\} = \int_{\mathbb{R} \times A} f(x, y) dx \otimes dy = \int_A \phi(y) dy.$$

In other words, ϕ is a density for Y . Now let

$$f(x|y) = \begin{cases} \frac{f(x, y)}{\phi(y)} & \text{if } \phi(y) > 0 \\ 0 & \text{otherwise.} \end{cases}$$

This is a Borel function from \mathbb{R}^2 to \mathbb{R} (**exercise!**), and so is a Borel function in either coordinate (when the other is held fixed).

We now want to define $\mathbf{E}[X | Y] = \int_{\mathbb{R}} x f(x|Y) dx$. We really need to work on the set that the integral is defined, so let $Z = \int_{\mathbb{R}} x f(x|Y) dx \cdot \mathbf{1}_{[Y \in F]}$. Then Z is a composition of Y with a Borel map, so is $\mathcal{G}/\mathcal{B}(\mathbb{R})$ measurable, and using the change of variables formula and monotonicity of probability (or Jensen’s inequality),

$$\begin{aligned} \mathbf{E}|Z| &= \mathbf{E} \left[\left| \int_{\mathbb{R}} x f(x|Y) dx \right| \mathbf{1}_{[Y \in F]} \right] \\ &= \int_{\mathbb{R}} \left| \int_{\mathbb{R}} x f(x|y) dx \right| \phi(y) dy \\ &= \int_{\mathbb{R}} \left| \int_{\mathbb{R}} x f(x, y) dx \right| dy \\ &\leq \int_{\mathbb{R}} \int_{\mathbb{R}} x f(x, y) dx dy \\ &= \mathbf{E}|X| < \infty. \end{aligned}$$

Thus $Z \in L_1(\Omega, \mathcal{G}, \mathbf{P})$.

Finally, for any $E \in \sigma(Y)$ we can write $E = \{Y \in A\}$ for some $A \in \mathcal{B}(\mathbb{R})$, so by two applications of the change of variables formula,

$$\begin{aligned} \mathbf{E}[X\mathbf{1}_{[Y \in A]}] &= \int_{\mathbb{R} \times A} xf(x, y) dx \otimes dy \\ &= \int_A \int_{\mathbb{R}} xf(x, y) dx dy \\ &= \int_A \int_{\mathbb{R}} xf(x|y)\phi(y) dx dy \\ &= \int_A \int_{\mathbb{R}} xf(x|y) dx \phi(y) dy \\ &= \mathbf{E}\left[\int_{\mathbb{R}} xf(x|Y)\mathbf{1}_{[Y \in A]}\right]. \end{aligned}$$

Thus $\int_{\mathbb{R}} xf(x|Y) dx$ is indeed a version of $\mathbf{E}[X | Y]$.

3. Suppose that X and Y are independent and that $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a Borel function with $\mathbf{E}|\phi(X, Y)| < \infty$. Let $F = \{y \in \mathbb{R} : \mathbf{E}|\phi(X, y)| < \infty\}$ and set $g(y) := (\mathbf{E}\phi(X, y))\mathbf{1}_{[y \in F]}$. We claim that $g(Y) \stackrel{\text{a.s.}}{=} \mathbf{E}\phi(X, Y) | Y$.

To see this, first note that by the change of variables formula and monotonicity of integration,

$$\begin{aligned} \mathbf{E}|g(Y)| &= \int_F |\mathbf{E}\phi(X, y)| d\mu_Y \\ &= \int_F \left| \int_{\mathbb{R}} \phi(x, y) d\mu_X \right| d\mu_Y \\ &\leq \int_F \int_{\mathbb{R}} |\phi(x, y)| d\mu_X d\mu_Y. \end{aligned}$$

Since X and Y are independent, the pair (X, Y) has joint law $d_X \otimes d_Y$, so by Fubini's theorem and another application of the change of variables formula,

$$\begin{aligned} \int_F \int_{\mathbb{R}} |\phi(x, y)| d\mu_X d\mu_Y &= \int_{\mathbb{R}^2} |\phi(x, y)| d(\mu_X \otimes \mu_Y) \\ &= \mathbf{E}|\phi(X, Y)| < \infty, \end{aligned}$$

so $g(Y) \in L_1(\Omega, \sigma(Y), \mathbf{P})$. Next, for any $A \in \mathcal{B}(\mathbb{R})$, again using change of variables and Fubini we have

$$\begin{aligned} \mathbf{E}[g(Y)\mathbf{1}_{[Y \in A]}] &= \int_F \mathbf{E}\phi(X, y)\mathbf{1}_{[A]}(y) d\mu_Y \\ &= \int_F \int_{\mathbb{R}} \phi(x, y)\mathbf{1}_{[A]}(y) d\mu_X d\mu_Y \\ &= \int_{A \times \mathbb{R}} \phi(x, y) d(\mu_X \otimes \mu_Y) \\ &= \mathbf{E}[\phi(X, Y)\mathbf{1}_{[Y \in A]}], \end{aligned}$$

as required.

4. This example is a straightforward generalization of the previous one to more than two random variables, and a detailed justification is omitted (only the construction is given). Suppose (X_1, \dots, X_n) are independent random variables on a probability space $(\Omega, \mathcal{F}, \mathbf{P})$. Fix any Borel function $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ such that $\phi(X_1, \dots, X_n) \in L_1(\Omega, \mathcal{F}, \mathbf{P})$. Fix $1 \leq i \leq n$ and let $\mathcal{G} = \sigma(X_1, \dots, X_i)$.

Let

$$F = \{(x_1, \dots, x_i) \in \mathbb{R}^i : \mathbf{E}|\phi(x_1, \dots, x_i, X_{i+1}, \dots, X_n)| < \infty\}.$$

Define $g : \mathbb{R}^i \rightarrow \mathbb{R}$ by

$$g(x_1, \dots, x_i) = \mathbf{E} [\phi(x_1, \dots, x_i, X_{i+1}, \dots, X_n)] \mathbf{1}_{[(x_1, \dots, x_i) \in F]}.$$

Then $g(X_1, \dots, X_i)$ is a version of $\mathbf{E} [\phi(X_1, \dots, X_n) \mid \mathcal{G}]$. This construction subsumes¹⁴ examples 1 through 4 from Section 11.1.

Exercise 12.3. Revisit the examples of Section 11.1, considering the projections from the perspective of conditional expectations. Check that you see why the projections satisfy the defining properties of conditional expectation for their respective random variables.

12.1. Properties of conditional expectation. In this section we record a litany¹⁵ of basic properties satisfied by conditional expectation. We always assume $(\Omega, \mathcal{F}, \mathbf{P})$ is a probability space, that $X, Y \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ and that \mathcal{G} is a sub- σ -field of \mathcal{F} .

- (i) $\mathbf{E} [\mathbf{E} \{X \mid \mathcal{G}\}] = \mathbf{E} X$
- (ii) If $\sigma(X) \subset \mathcal{G}$ then $\mathbf{E} \{X \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} X$.

Proving the first two properties is an exercise in understanding and applying the definition of conditional expectation.

- (iii) If $\sigma(X)$ and \mathcal{G} are independent then $\mathbf{E} \{X \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} \mathbf{E} X$.

Proof: For all $A \in \mathcal{G}$, by the independence assumption,

$$\mathbf{E} [X \mathbf{1}_{[A]}] = \mathbf{E} X \cdot \mathbf{E} [\mathbf{1}_{[A]}] = \mathbf{E} [(\mathbf{E} X) \cdot \mathbf{1}_{[A]}]. \quad \square$$

- (iv) **Linearity of conditional expectation.** For all $a \in \mathbb{R}$, $\mathbf{E} \{aX + Y \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} a\mathbf{E} \{X \mid \mathcal{G}\} + \mathbf{E} \{Y \mid \mathcal{G}\}$.

- (v) **Monotonicity.** If $X \leq Y$ almost surely then $\mathbf{E} \{X \mid \mathcal{G}\} \leq \mathbf{E} \{Y \mid \mathcal{G}\}$ almost surely.

The last fact is just a restatement of Proposition 12.1.

For the next three properties, we additionally require a sequence $(X_n, n \geq 1)$ of random variables over $(\Omega, \mathcal{F}, \mathbf{P})$. The first is left as an **exercise**.

- (v) **Conditional Monotone Convergence Theorem.** If $0 \leq X_n \uparrow X$ almost surely then $\mathbf{E} \{X_n \mid \mathcal{G}\} \uparrow \mathbf{E} \{X \mid \mathcal{G}\}$ almost surely.
- (vi) **Conditional Fatou's Lemma.** If $X_n \geq 0$ for all n then $\mathbf{E} \{\liminf_{n \rightarrow \infty} X_n \mid \mathcal{G}\} \leq \liminf_{n \rightarrow \infty} \mathbf{E} \{X_n \mid \mathcal{G}\}$.

Proof: For any $n \geq 1$, for all $n' \geq n$ we have $X_{n'} \geq \inf_{m \geq n} X_m$, and it follows by monotonicity of conditional expectation that

$$\inf_{m \geq n} \mathbf{E} \{X_m \mid \mathcal{G}\} \stackrel{\text{a.s.}}{\geq} \mathbf{E} \left\{ \inf_{m \geq n} X_m \mid \mathcal{G} \right\}.$$

Taking $n \rightarrow \infty$ on both sides gives

$$\liminf_{n \rightarrow \infty} \mathbf{E} \{X_n \mid \mathcal{G}\} \stackrel{\text{a.s.}}{\geq} \lim_{n \rightarrow \infty} \mathbf{E} \left\{ \inf_{m \geq n} X_m \mid \mathcal{G} \right\} \stackrel{\text{a.s.}}{=} \mathbf{E} \left\{ \liminf_{n \rightarrow \infty} X_n \mid \mathcal{G} \right\},$$

where the almost sure equality follows from the conditional monotone convergence theorem. □

- (vii) **Conditional Dominated Convergence Theorem.** If $X_n \xrightarrow{\text{a.s.}} X$ almost surely and $|X_n| \leq Y$ almost surely for all n , then $\mathbf{E} \{X_n \mid \mathcal{G}\} \xrightarrow{\text{a.s.}} \mathbf{E} \{X \mid \mathcal{G}\}$.

¹⁴Subsume, v.: 6. transitive. a. To take up or absorb (a concept, thing, person, etc.) into another, esp. one which is larger or higher; to include in. b. To bring (an idea, principle, etc.) under another; to instance or include (a case, term, etc.) under a rule, category, etc. –Oxford English Dictionary

¹⁵Litany, n.: 2. *transferred*. A form of supplication (e.g. in non-Christian worship) resembling a litany; also, a continuous repetition or long enumeration resembling those of litanies. –Oxford English Dictionary

The conditional dominated convergence theorem follows from the conditional Fatou's lemma in essentially the same way as the dominated convergence theorem follows from Fatou's lemma.

We next turn to inequalities related to convexity.

- (viii) **Conditional Jensen's inequality.** If $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is convex and $\varphi(X) \in L_1(\Omega, \mathcal{F}, \mathbf{P})$, then $\varphi(\mathbf{E}\{X \mid \mathcal{G}\}) \stackrel{\text{a.s.}}{\leq} \mathbf{E}\{\varphi(X) \mid \mathcal{G}\}$.

Proof: We may fix a sequence of linear functions $\ell_n(x) = a_n x + b_n$ such that for all $x \in \mathbb{R}$, $\varphi(x) = \sup_{n \geq 1} (a_n x + b_n)$. We then have $\varphi(X) \geq a_n X + b_n$ for all n , so by monotonicity and linearity of conditional expectation,

$$\mathbf{E}\{\varphi(X) \mid \mathcal{G}\} \stackrel{\text{a.s.}}{\geq} \mathbf{E}\{a_n X + b_n \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} a_n \mathbf{E}\{X \mid \mathcal{G}\} + b_n.$$

Taking a supremum over $n \geq 1$ gives

$$\mathbf{E}\{\varphi(X) \mid \mathcal{G}\} \stackrel{\text{a.s.}}{\geq} \sup_{n \geq 1} (a_n \mathbf{E}\{X \mid \mathcal{G}\} + b_n) = \varphi(\mathbf{E}\{X \mid \mathcal{G}\}). \quad \square$$

- (ix) For all $p \geq 1$, $\|\mathbf{E}\{X \mid \mathcal{G}\}\|_p \leq \|X\|_p$.

Proof: This is obvious if $\|X\|_p = \infty$. Otherwise, by the conditional Jensen's inequality applied to the function $\phi(x) = |x|^p$ with the random variable $X^p \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ we have $|\mathbf{E}\{X \mid \mathcal{G}\}|^p \leq \mathbf{E}\{|X|^p \mid \mathcal{G}\}$. It follows by monotonicity and the definition of conditional expectation that

$$\begin{aligned} \|\mathbf{E}\{X \mid \mathcal{G}\}\|_p^p &= \mathbf{E}[|\mathbf{E}\{X \mid \mathcal{G}\}|^p] \\ &\leq \mathbf{E}[\mathbf{E}\{|X|^p \mid \mathcal{G}\}] \\ &= \mathbf{E}[\mathbf{E}\{|X|^p \mid \mathcal{G}\} \mathbf{1}_{[\Omega]}] = \mathbf{E}[|X|^p \mathbf{1}_{[\Omega]}] = \|X\|_p^p. \end{aligned} \quad \square$$

- (x) **Conditional Hölder's inequality.** For $p, q \geq 1$ with $1/p + 1/q = 1$. If $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ and $Y \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ then $XY \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ and

$$\mathbf{E}\{|XY| \mid \mathcal{G}\} \leq (\mathbf{E}\{|X|^p \mid \mathcal{G}\})^{1/p} (\mathbf{E}\{|Y|^q \mid \mathcal{G}\})^{1/q}.$$

We briefly delay the proof of Hölder's inequality as it uses a property of conditional expectation we have not yet seen.

The next three properties are perhaps less “intuitive”, as they are not simply conditional versions of facts you have already seen. The first is related to the fact that the projection operation is idempotent. The second says that if $\sigma(Y) \subset \mathcal{G}$ then Y “acts like a constant” with respect to conditional expectations given \mathcal{G} . The third says (informally) that conditioning a conditional expectation on another independent σ -field doesn't change anything.

- (xi) **The tower property.** If $\mathcal{H} \subset \mathcal{G}$ is another σ -field, then

$$\mathbf{E}\{\mathbf{E}\{X \mid \mathcal{G}\} \mid \mathcal{H}\} \stackrel{\text{a.s.}}{=} \mathbf{E}\{X \mid \mathcal{H}\} \stackrel{\text{a.s.}}{=} \mathbf{E}\{\mathbf{E}\{X \mid \mathcal{H}\} \mid \mathcal{G}\}.$$

Proof: First, $\mathbf{E}\{X \mid \mathcal{H}\}$ is $\mathcal{H}/\mathcal{B}(\mathbb{R})$ -measurable, and $\mathcal{H} \subset \mathcal{G}$, so by property (ii),

$$\mathbf{E}\{\mathbf{E}\{X \mid \mathcal{H}\} \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} \mathbf{E}\{X \mid \mathcal{H}\}.$$

Next, let Z be a version of $\mathbf{E}\{\mathbf{E}\{X \mid \mathcal{G}\} \mid \mathcal{H}\}$. Then by definition, $Z \in L_1(\Omega, \mathcal{H}, \mathbf{P})$ and for all $A \in \mathcal{H}$,

$$\mathbf{E}[Z \mathbf{1}_{[A]}] = \mathbf{E}[\mathbf{E}\{X \mid \mathcal{G}\} \mathbf{1}_{[A]}].$$

By the definition of $\mathbf{E}\{X \mid \mathcal{G}\}$, we also have $\mathbf{E}[\mathbf{E}\{X \mid \mathcal{G}\} \mathbf{1}_{[A]}] = \mathbf{E}[X \mathbf{1}_{[A]}]$. It follows that $\mathbf{E}[Z \mathbf{1}_{[A]}] = \mathbf{E}[X \mathbf{1}_{[A]}]$, so Z is a version of $\mathbf{E}\{X \mid \mathcal{H}\}$. \square

- (xii) **Moving variables out of conditional expectations.** For any random variable $Z \in L_\infty(\Omega, \mathcal{G}, \mathbf{P})$ it holds that $\mathbf{E}\{XZ \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} \mathbf{E}\{X \mid \mathcal{G}\} Z$.

Proof: Let

$$\mathcal{S} = \left\{ Z \in L_\infty(\Omega, \mathcal{G}, \mathbf{P}) : \mathbf{E}\{XZ \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} \mathbf{E}\{X \mid \mathcal{G}\} Z \right\}.$$

We aim to show that $\mathcal{S} = L_\infty(\Omega, \mathcal{G}, \mathbf{P})$. First suppose $Z = \mathbf{1}_{[B]}$ for some $B \in \mathcal{G}$. Then for all $A \in \mathcal{G}$,

$$\mathbf{E} [\mathbf{E} \{X \mid \mathcal{G}\} Z \mathbf{1}_{[A]}] = \mathbf{E} [\mathbf{E} \{X \mid \mathcal{G}\} \mathbf{1}_{[A \cap B]}] = \mathbf{E} [X \mathbf{1}_{[A \cap B]}] = \mathbf{E} [X Z \mathbf{1}_{[A]}],$$

so by the definition of conditional expectation, $\mathbf{E} \{X \mid \mathcal{G}\} Z$ is a version of $\mathbf{E} \{XZ \mid \mathcal{G}\}$ and therefore $Z \in \mathcal{S}$.

Next, if $Z, Z' \in \mathcal{S}$ and $a \in \mathbb{R}$ then by linearity of conditional expectation,

$$\begin{aligned} \mathbf{E} \{X \mid \mathcal{G}\} (aZ + Z') &= a\mathbf{E} \{X \mid \mathcal{G}\} Z + \mathbf{E} \{X \mid \mathcal{G}\} Z' \\ &\stackrel{\text{a.s.}}{=} \mathbf{E} \{aXZ + XZ' \mid \mathcal{G}\} = \mathbf{E} \{X(aZ + Z') \mid \mathcal{G}\}, \end{aligned}$$

so $aZ + Z' \in \mathcal{S}$.

Next, if $0 \leq Z_n \in \mathcal{S}$ for $n \geq 1$ and $Z_n \uparrow Z$ as $n \rightarrow \infty$, then $X^+ Z_n \uparrow X^+ Z$ as $n \rightarrow \infty$, so by the conditional monotone convergence theorem,

$$\mathbf{E} \{X^+ \mid \mathcal{G}\} Z = \lim_{n \rightarrow \infty} \mathbf{E} \{X^+ \mid \mathcal{G}\} Z_n \stackrel{\text{a.s.}}{=} \lim_{n \rightarrow \infty} \mathbf{E} \{X^+ Z_n \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} \mathbf{E} \{X^+ Z \mid \mathcal{G}\}.$$

Likewise $\mathbf{E} \{X^- \mid \mathcal{G}\} Z \stackrel{\text{a.s.}}{=} \mathbf{E} \{X^- Z \mid \mathcal{G}\}$, so by linearity of conditional expectation,

$$\mathbf{E} \{X \mid \mathcal{G}\} Z \stackrel{\text{a.s.}}{=} (\mathbf{E} \{X^+ \mid \mathcal{G}\} + \mathbf{E} \{X^- \mid \mathcal{G}\}) Z \stackrel{\text{a.s.}}{=} \mathbf{E} \{X^+ Z + X^- Z \mid \mathcal{G}\} = \mathbf{E} \{XZ \mid \mathcal{G}\}.$$

Thus $Z \in \mathcal{S}$. It follows by the monotone class theorem that $\mathcal{S} = L_\infty(\Omega, \mathcal{G}, \mathbf{P})$. \square

(xiii) **Adding an independent conditioning changes nothing.** If $\mathcal{H} \subset \mathcal{F}$ is another σ -algebra and $\sigma(X, \mathcal{G}) := \sigma(\sigma(X) \cup \mathcal{G})$ is independent of \mathcal{H} then $\mathbf{E} \{X \mid \sigma(\mathcal{G}, \mathcal{H})\} = \mathbf{E} \{X \mid \mathcal{G}\}$.

The last property is left as an **exercise**. We also state the following strengthening of (xiii) as an exercise.

Exercise 12.4. Prove that if $Z : \Omega \rightarrow \mathbb{R}$ is $\mathcal{G}/\mathcal{B}(\mathbb{R})$ -measurable and $X, XZ \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ then $\mathbf{E} \{XZ \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} \mathbf{E} \{X \mid \mathcal{G}\} Z$.

Proof of Holder's inequality. To avoid issues of integrability and dividing by zero, for $\epsilon \geq 0$ write $U_\epsilon = (\mathbf{E} \{|X|^p \mid \mathcal{G}\} + \epsilon)^{1/p}$ and $V_\epsilon = (\mathbf{E} \{|Y|^q \mid \mathcal{G}\} + \epsilon)^{1/q}$. Then let $X_\epsilon = \frac{|X|}{U_\epsilon}$ and $Y_\epsilon = \frac{Y}{V_\epsilon}$.

For $\epsilon > 0$ we then have

$$X_\epsilon Y_\epsilon = \exp \left(\frac{1}{p} \log(X_\epsilon^p) + \frac{1}{q} \log(Y_\epsilon^q) \right) \leq \exp \left(\log \frac{X_\epsilon^p}{p} + \frac{Y_\epsilon^q}{q} \right) = \frac{X_\epsilon^p}{p} + \frac{Y_\epsilon^q}{q},$$

so

$$\begin{aligned} \mathbf{E} \{X_\epsilon Y_\epsilon \mid \mathcal{G}\} &\leq \frac{1}{p} \mathbf{E} \{X_\epsilon^p \mid \mathcal{G}\} + \frac{1}{q} \mathbf{E} \{Y_\epsilon^q \mid \mathcal{G}\} \\ &= \frac{1}{p} \mathbf{E} \{|X|^p U_\epsilon^{-p} \mid \mathcal{G}\} + \frac{1}{q} \mathbf{E} \{|Y|^q V_\epsilon^{-q} \mid \mathcal{G}\} \end{aligned}$$

The terms U_ϵ^{-p} and V_ϵ^{-q} are in $L_\infty(\Omega, \mathcal{G}, \mathbf{R})$, so by (xiii) we may move them outside the conditional expectations. The previous bound then beomes

$$\mathbf{E} \{X_\epsilon Y_\epsilon \mid \mathcal{G}\} \leq \frac{1}{p} \frac{\mathbf{E} \{|X|^p \mid \mathcal{G}\}}{U_\epsilon^p} + \frac{1}{q} \frac{\mathbf{E} \{|Y|^q \mid \mathcal{G}\}}{V_\epsilon^q} = \frac{1}{p} \frac{U_0^p}{U_\epsilon^p} + \frac{1}{q} \frac{V_0^q}{V_\epsilon^q} \leq 1.$$

Again using that U_ϵ^{-p} and V_ϵ^{-q} are in $L_\infty(\Omega, \mathcal{G}, \mathbf{R})$, we also have

$$\mathbf{E} \{X_\epsilon Y_\epsilon \mid \mathcal{G}\} \stackrel{\text{a.s.}}{=} \frac{\mathbf{E} \{|XY| \mid \mathcal{G}\}}{U_\epsilon V_\epsilon},$$

which combined with the previous inequality gives

$$\mathbf{E} \{|XY| \mid \mathcal{G}\} \stackrel{\text{a.s.}}{\leq} U_\epsilon V_\epsilon.$$

Taking $\epsilon \downarrow 0$, the result follows. \square

We conclude the section with a slight extension of the domain of definition of conditional expectations. When defining conditional expectation, in the proof of Theorem 12.2, for a non-negative random variable $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ we defined $\mathbf{E}\{X \mid \mathcal{G}\}$ as the almost sure increasing limit of $\mathbf{E}\{X^{\leq n} \mid \mathcal{G}\}$. More generally, if $X : \Omega \rightarrow [0, \infty)$ is $\mathcal{F}/\mathcal{B}(\mathbb{R})$ -measurable (even if $\mathbf{E}X = \infty$), we let $Z = \liminf_{n \rightarrow \infty} \mathbf{E}\{X^{\leq n} \mid \mathcal{G}\}$. We call Z , or any random variable in its almost sure equivalence class, a version of $\mathbf{E}\{X \mid \mathcal{G}\}$.

It is important to note if we do not insist that $\mathbf{E}[X] < \infty$ then it can occur that that $\mathbf{E}\{X \mid \mathcal{G}\}$ takes the value $+\infty$ with positive probability. For example, suppose that $\mathbf{P}\{X = k\} = \frac{1}{\zeta(2)} \frac{1}{k^2}$ for $k \geq 1$. Let $\mathcal{G} = \sigma(\{X = 1\}) = \{\emptyset, \{X = 1\}, \{X > 1\}, \Omega\}$. Noting that $\mathbf{P}\{X = k \mid X > 1\} = \frac{\zeta(2)-1}{\zeta(2)} \frac{1}{k^2}$ for $k > 1$, it follows that

$$\mathbf{E}\{X^{\leq n} \mid \mathcal{G}\}(\omega) = \begin{cases} 1 & \text{if } X(\omega) = 1 \\ \frac{\zeta(2)-1}{\zeta(2)} \sum_{k=2}^n \frac{1}{k} & \text{if } X(\omega) > 1, \end{cases}$$

so

$$\mathbf{E}\{X \mid \mathcal{G}\}(\omega) = \begin{cases} 1 & \text{if } X(\omega) = 1 \\ \infty & \text{if } X(\omega) > 1. \end{cases}$$

Having accepted the fact that conditional expectations of non-negative random variables can take the value $+\infty$, we can even go a little further. If $X : \Omega \rightarrow [0, \infty]$ is $\mathcal{F}/\mathcal{B}(\mathbb{R})$ -measurable, so X is a non-negative extended real random variable, then we define

$$\begin{aligned} \mathbf{E}\{X \mid \mathcal{G}\} &= \mathbf{E}\{X \mathbf{1}_{[X < \infty]} \mid \mathcal{G}\} + \infty \cdot \mathbf{E}\{\mathbf{1}_{[X = \infty]} \mid \mathcal{G}\} \\ &= \lim_{n \rightarrow \infty} \mathbf{E}\{X \mathbf{1}_{[X < n]} \mid \mathcal{G}\} + \infty \cdot \mathbf{E}\{\mathbf{1}_{[X = \infty]} \mid \mathcal{G}\}. \end{aligned}$$

The preceding equalities should be understood at the level of almost sure equivalence classes. This definition generalizes the definition given for non-negative real random variables and agrees with that definition when $\mathbf{P}\{X = \infty\} = 0$. It will come up later, in particular when exploring the connection between martingales and the Radon-Nikodym theorem.

Proposition 12.3 (Conditional monotone convergence theorem). *If $0 \leq X_n \uparrow X \leq \infty$ then $\lim_{n \rightarrow \infty} \mathbf{E}[X_n \mid \mathcal{G}] = \mathbf{E}[X \mid \mathcal{G}]$ almost surely.*

Proof. By monotonicity of conditional expectations, there exists an a.s. increasing limit $Y = \lim_{n \rightarrow \infty} \mathbf{E}[X_n \mid \mathcal{G}]$. For any event $E \in \mathcal{G}$,

$$\begin{aligned} \mathbf{E}[Y \mathbf{1}_{[E]}] &= \mathbf{E}\left[\lim_{n \rightarrow \infty} \mathbf{E}[X_n \mid \mathcal{G}] \mathbf{1}_{[E]}\right] \\ &= \lim_{n \rightarrow \infty} \mathbf{E}[\mathbf{E}[X_n \mid \mathcal{G}] \mathbf{1}_{[E]}] && \text{(Monotone convergence)} \\ &= \lim_{n \rightarrow \infty} \mathbf{E}[\mathbf{E}[X_n \mathbf{1}_{[E]} \mid \mathcal{G}]] && \text{(Since } E \in \mathcal{G}\text{)} \\ &= \lim_{n \rightarrow \infty} \mathbf{E}[X_n \mathbf{1}_{[E]}] && \text{(Def. of cond. expectation)} \\ &= \mathbf{E}\left[\lim_{n \rightarrow \infty} X_n \mathbf{1}_{[E]}\right] && \text{(Monotone convergence)} \\ &= \mathbf{E}[X \mathbf{1}_{[E]}]. \end{aligned}$$

Thus Y is a version of $\mathbf{E}[X \mid \mathcal{G}]$. \square

Exercise 12.5 (Conditional extended Fatou's lemma). *If $(X_n, n \geq 0)$ is any sequence of non-negative extended real-valued random variables, then $\mathbf{E}[\liminf_{n \rightarrow \infty} X_n \mid \mathcal{G}] \leq \liminf_{n \rightarrow \infty} \mathbf{E}[X_n \mid \mathcal{G}]$.*

12.2. Conditional expectations, tightness and uniform integrability. Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and a collection $X = (X_i, i \in I)$ of random variables in $L_1(\Omega, \mathcal{F}, \mathbf{P})$.

tight
tight

Write μ_i for the distribution of X_i . The family $(\mu_i, i \in I)$ of probability measures is *tight* if for all $\epsilon > 0$ there is a compact set $K \subset \mathbb{R}$

$$\sup_{i \geq 1} \mu_i(\mathbb{R} \setminus K) < \epsilon.$$

The collection X is *uniformly integrable* with respect to \mathbf{P} if for all $\epsilon > 0$ there is a compact set $K \subset \mathbb{R}$ such that

$$\sup_{i \in I} \mathbf{E} [|X_i| \mathbf{1}_{[|X_i| \notin K]}] < \epsilon.$$

The two conditions are syntactically similar. They are connected by using the *size-biasing* operation introduced earlier. Write $\hat{\mu}_i$ for the size-biasing of μ_i , so

$$\hat{\mu}_i(B) = \left(\mu_i \cdot \frac{|X_i|}{\mathbf{E}|X_i|} \right) (B) = \mathbf{E} [|X_i| \mathbf{1}_{[X_i \in B]}].$$

Exercise 12.6. A collection $X = (X_i, i \in I)$ of random variables in $L_1(\Omega, \mathcal{F}, \mathbf{P})$ is uniformly integrable if and only if $(\hat{\mu}_i, i \in I)$ is tight.

Exercise 12.7. Let $(\mu_n, n \geq 1)$ be a tight family of probability measures. Then there exists a subsequence $(n_k, k \geq 1)$ such that μ_{n_k} converges in distribution. (I.e. such that if X_{n_k} has distribution μ_{n_k} then X_{n_k} converges in distribution.)

Exercise 12.8. Let $(X_n, 1 \leq n \leq \infty)$ be random variables in $L_1(\Omega, \mathcal{F}, \mathbf{P})$ such that $X_n \xrightarrow{\mathbf{P}} X_\infty$. Prove that the following are equivalent: (a) $X_n \xrightarrow{L_1} X_\infty$; (b) $(X_n, 1 \leq n \leq \infty)$ is uniformly integrable; (c) $\mathbf{E}|X_n| \rightarrow \mathbf{E}|X_\infty|$.

The next proposition connects uniform integrability and conditional expectations, and is the first step toward martingales and martingale convergence theorems.

Proposition 12.4. Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and a random variable $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$. Then $\{\mathbf{E}\{X \mid \mathcal{G}\} : \mathcal{G} \subset \mathcal{F} \text{ a sub-}\sigma\text{-field}\}$ is a uniformly integrable collection of random variables.

Lemma 12.5. If $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ then for all $\epsilon > 0$ there is $\delta > 0$ such that for all $A \in \mathcal{F}$, if $\mathbf{P}\{A\} \leq \delta$ then $\mathbf{E}[|X| \mathbf{1}_{[A]}] < \epsilon$.

Proof. Suppose that the assertion of the lemma is false. Then we may find $\epsilon > 0$ and events $(A_n, n \geq 1)$ in \mathcal{F} with $\mathbf{P}\{A_n\} \leq 2^{-n}$ such that $\mathbf{E}[|X| \mathbf{1}_{[A_n]}] \geq \epsilon$ for all n .

We now show this implies that $\mathbf{E}[|X| \mathbf{1}_{[A_n \text{ i.o.}]}] \geq \epsilon$. By definition, $\{A_n \text{ i.o.}\} = \bigcap_{n \geq 1} \bigcup_{m \geq n} A_m$, so $\mathbf{1}_{[A_n \text{ i.o.}]} = \mathbf{1}_{[\bigcap_{n \geq 1} \bigcup_{m \geq n} A_m]} = \lim_{n \rightarrow \infty} \mathbf{1}_{[\bigcup_{m \geq n} A_m]}$

For any event $E \in \mathcal{F}$ we have $|X| \mathbf{1}_{[E]} \leq |X|$, so by the dominated convergence theorem,

$$\mathbf{E}[|X| \mathbf{1}_{[A_n \text{ i.o.}]}] = \mathbf{E}\left[\lim_{n \rightarrow \infty} |X| \mathbf{1}_{[\bigcup_{m \geq n} A_m]}\right] = \lim_{n \rightarrow \infty} \mathbf{E}[|X| \mathbf{1}_{[\bigcup_{m \geq n} A_m]}] \geq \epsilon.$$

On the other hand, $\sum_{n \geq 1} \mathbf{P}\{A_n\} = 1 < \infty$, so by the first Borel-Cantelli lemma, $\mathbf{P}\{A_n \text{ i.o.}\} = 0$ and thus $\mathbf{E}[|X| \mathbf{1}_{[A_n \text{ i.o.}]}] = 0$, a contradiction. \square

Proof of Proposition 12.4. Fix $\epsilon > 0$ and let $\delta > 0$ be such that $\mathbf{E}[|X| \mathbf{1}_{[A]}] < \epsilon$ whenever $\mathbf{P}\{A\} \leq \delta$; such δ exists by the lemma. Then for any sub- σ -field $\mathcal{G} \subset \mathcal{F}$, by the conditional Jensen's inequality

$$|\mathbf{E}\{X \mid \mathcal{G}\}| \leq \mathbf{E}\{|X| \mid \mathcal{G}\},$$

so

$$\mathbf{E}[|\mathbf{E}\{X \mid \mathcal{G}\}|] \leq \mathbf{E}[\mathbf{E}\{|X| \mid \mathcal{G}\}] = \mathbf{E}|X|,$$

Taking $K = [-\mathbf{E}|X|/\delta, \mathbf{E}|X|/\delta]$, it follows that

$$\begin{aligned} \mathbf{P} \{ |\mathbf{E}\{X \mid \mathcal{G}\}| \notin K \} &= \mathbf{P} \{ |\mathbf{E}\{X \mid \mathcal{G}\}| > \mathbf{E}|X|/\delta \} \\ &\leq \mathbf{E} [|\mathbf{E}\{X \mid \mathcal{G}\}| \mathbf{E}|X|/\delta] \\ &\leq \delta. \end{aligned}$$

It follows that

$$\begin{aligned} \mathbf{E} [|\mathbf{E}\{X \mid \mathcal{G}\}| \mathbf{1}_{\{|\mathbf{E}\{X \mid \mathcal{G}\}| \notin K\}}] &\leq \mathbf{E} [\mathbf{E}\{|X| \mid \mathcal{G}\} \mathbf{1}_{\{|\mathbf{E}\{X \mid \mathcal{G}\}| \notin K\}}] \\ &= \mathbf{E} [\mathbf{E}\{|X| \mathbf{1}_{\{|\mathbf{E}\{X \mid \mathcal{G}\}| \notin K\}} \mid \mathcal{G}\}] \\ &= \mathbf{E} [|X| \mathbf{1}_{\{|\mathbf{E}\{X \mid \mathcal{G}\}| \notin K\}}] \\ &\leq \epsilon, \end{aligned}$$

the last bound holding since $\mathbf{P} \{ |\mathbf{E}\{X \mid \mathcal{G}\}| \notin K \} \leq \delta$. □

13. Martingales

A stochastic process is simply a family of random variables $(X_i, i \in I)$ defined over a common probability space $(\Omega, \mathcal{F}, \mathbf{P})$. Martingales are stochastic processes which model “fair games”, or random systems which evolve in time without a bias in any particular direction. They are one of the most important general classes of stochastic processes; the next part of these notes is devoted to defining martingales and understanding their properties.

A *filtration* is an increasing sequence of σ -algebras $(\mathcal{F}_n)_{n \geq 0}$ over a common ground set. A *filtered probability space* is a tuple $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$, where $(\mathcal{F}_n)_{n \geq 0}$ is a filtration over Ω and $\mathcal{F}_n \subset \mathcal{F}$ for all $n \geq 0$.

A sequence $X = (X_n)_{n \geq 0}$ of random variables is (\mathcal{F}_n) -*adapted* if X_n is $\mathcal{F}_n/\mathcal{B}(\mathbb{R})$ -measurable for all $n \geq 0$. It is *integrable* if $X_n \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ for all $n \geq 1$. Finally, it is an (\mathcal{F}_n) -*martingale* (or just a martingale for short) if it is integrable and adapted and satisfies the *martingale property*: for all $n > 0$,

$$\mathbf{E}\{X_n \mid \mathcal{F}_{n-1}\} = X_{n-1}.$$

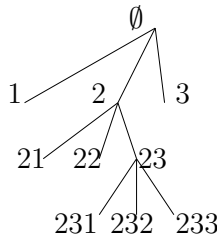
If you think of $(X_n)_{n \geq 1}$ as a stock value (for example), then the martingale property states that the best prediction for the stock’s future value given its past performance is simply its present value.¹⁶

Example: Simple random walk. Let $(Z_i, i \geq 1)$ be iid random variables in $L_1(\Omega, \mathcal{F}, \mathbf{P})$ with $\mathbf{E}Z_1 = 0$ and let $X_n = Z_1 + \dots + Z_n$. Then with

$$\mathcal{F}_n = \sigma(Z_1, \dots, Z_n) = \sigma(X_1, \dots, X_n),$$

the sequence $X = (X_n, n \geq 0)$ is an (\mathcal{F}_n) -martingale.

Example: branching processes. Let μ be a probability measure on \mathbb{R} with $\mu(\mathbb{N}) = 1$.



- ★ Start from the root (call it \emptyset), let B_\emptyset have law μ .
- ★ Give \emptyset children $1, \dots, B_\emptyset$.
- ★ Independently for each $i = 1, \dots, B_\emptyset$, let B_i have law μ .
- ★ Give i children $i1, i2, \dots, iB_i$.
- ★ Repeat forever or until done; call the resulting random tree \mathcal{T}_B .

Let Z_n be the number of individuals in the n 'th generation (the individuals of the n 'th generation are those whose name is n characters long), and write $|\mathcal{T}_B| = \sum_{n=0}^\infty Z_n$ for the total number of individuals. We say the survival occurs if $Z_n > 0$ for all n , and otherwise that say that extinction occurs. Equivalently, survival occurs if $|\mathcal{T}_B| = \infty$, and extinction occurs if $|\mathcal{T}_B| < \infty$.

¹⁶The efficient markets hypothesis in economics is essentially a statement that stocks behave like martingales.

For $n \geq 0$ let $\mathcal{F}_n = \sigma(Z_0, \dots, Z_n)$. Now fix $n \geq 0$ and let \mathcal{S} be the set of nodes in the n 'th generation of \mathcal{T}_B ; for $n \geq 1$ this is a random subset of $\mathbb{N}^n = \{b_1 b_2 \dots b_n : b_i \in \mathbb{N}, 1 \leq i \leq n\}$. Then $Z_{n+1} = \sum_{v \in \mathcal{S}} B_v$, so for any *fixed* subset S of \mathbb{N}^n ,

$$\mathbf{E} \{ Z_{n+1} \mid \mathcal{S} = S \} = \mathbf{E} \left\{ \sum_{v \in S} B_v \mid \mathcal{S} = S \right\} = \sum_{v \in S} \mathbf{E} B = |S| \cdot \mathbf{E} B.$$

The second inequality holds by linearity of expectation and since $\mathbf{E} \{ B_v \mid \mathcal{S} = S \} = \mathbf{E} B$. (We have been a bit informal about the last fact but it should be intuitively clear; we will be more precise about this later in the notes.) Since $Z_n = |\mathcal{S}|$, it follows that $\mathbf{E} \{ Z_{n+1} \mid \mathcal{F}_n \} = Z_n \cdot \mathbf{E} B$. Therefore, if $\mathbf{E} B = 1$ then $(Z_n, n \geq 0)$ is an \mathcal{F}_n -martingale. More generally, setting $M_n = Z_n / (\mathbf{E} B)^n$, then $(M_n, n \geq 0)$ is always an \mathcal{F}_n -martingale. We will analyze this example further in Section 14.

Exercise 13.1. Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space and let $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$. Write $X_n \stackrel{\text{a.s.}}{=} \mathbf{E} \{ X \mid \mathcal{F}_n \}$. Show that $(X_n, n \geq 0)$ is a martingale relative to the filtration $(\mathcal{F}_n, n \geq 0)$.

The main goal of the section is to find conditions which guarantee that a martingale $(X_n)_{n \geq 0}$ converges to some limit X in some sense. However, the convergence theory is not the only point and, in fact, the theory will be easier to understand and will appear better motivated if we first approach the subject from a more applied point of view.

Continuing the analogy with stock prices, suppose that $X = (X_n)_{n \geq 0}$ is an (\mathcal{F}_n) -adapted process, and think of it as tracking a stock price over time. At time n you can choose to invest some amount money C_{n+1} for the next unit of time, based on your observation of the stock's behaviour to date. In the next unit of time your profit/loss will then be $C_{n+1}(X_{n+1} - X_n)$.

An integrable stochastic process $(C_n)_{n \geq 1}$ is (\mathcal{F}_n) -previsible if $C_{n+1} \in L_1(\Omega, \mathcal{F}_n, \mathbf{P})$ for all $n \geq 0$. In the stock market analogy, saying that that C_{n+1} should be chosen based on past observations precisely means that means that $(C_n)_{n \geq 1}$ should be (\mathcal{F}_n) -previsible. If this is the case, then by the properties of conditional expectation (and assuming the random variables C_n are bounded), the profit/loss in the time unit from n to $n + 1$ is

$$\begin{aligned} \mathbf{E} [C_{n+1}(X_{n+1} - X_n)] &= \mathbf{E} [\mathbf{E} \{ C_{n+1}(X_{n+1} - X_n) \mid \mathcal{F}_n \}] \\ &= \mathbf{E} [\mathbf{E} \{ C_{n+1} X_{n+1} \mid \mathcal{F}_n \}] - \mathbf{E} [\mathbf{E} \{ C_{n+1} X_n \mid \mathcal{F}_n \}] \\ &= \mathbf{E} [C_{n+1} \mathbf{E} \{ X_{n+1} \mid \mathcal{F}_n \}] - \mathbf{E} [C_{n+1} X_n]. \end{aligned}$$

We've used that C_{n+1} and X_n are $(\mathcal{F}_n/\mathcal{B}(\mathbb{R}))$ -measurable to extract them from the conditional expectation. If X is an (\mathcal{F}_n) -martingale then by the martingale property, the last line equals zero, which means that gambling on this stock yields no expected profit or loss.

In the above setup, the total profit/loss by time n is

$$\sum_{i=1}^n C_i (X_i - X_{i-1})$$

which is our first glimpse at stochastic integration; it looks like a discrete analogue of an integral $\int_0^n X_i dC_i$. This perspective has been fruitfully developed into an entire academic discipline.

The theory of martingales is, among other things, a computational tool. Basic facts about martingales allow some expected values to be identified by appeal to general theory rather than via ad hoc calculations. For example, imagine that $(R_n)_{n \geq 0}$ tracks the dollar value of your current bankroll¹⁷ in a gambling game. You may choose to stop gambling at the first time T that either $R_n \geq 1000$ or $R_n = 0$. You will then return home with R_T dollars, and may care to know the expected value $\mathbf{E} [R_T]$. The *optional stopping theorem* says that, if you were playing a fair game, then $\mathbf{E} [R_T] = R_0$; you expect to walk out with whatever you brought in. Of course, most casinos don't

¹⁷Bankroll, n. originally and chiefly U.S. A roll of banknotes; (in extended use) the money a person possesses; funds, financial resources; (Gambling) the amount of money a person sets aside for a given session or period of gambling. Frequently with possessive. –Oxford English Dictionary

offer fair games. (If you are inclined to split hairs¹⁸, there are other issues with this as a model for gambling play; what is the meaning of R_n for $n > T$, for example?)

To state the optional stopping theorem we first need to define stopping times, and take the opportunity to state some elementary facts about them. Given a filtered probability space $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ a random variable $T : \Omega \rightarrow \mathbb{N} \cup \{+\infty\}$ is an (\mathcal{F}_n) -stopping time (or just “stopping time”) if for all $n \in \mathbb{N}$, the event $\{T \leq n\} \in \mathcal{F}_n$. The idea to have in mind is, if \mathcal{F}_n is the information available to you at time n , then saying T is a stopping time means that, if you are trying to stop at time T , then enough information is available to you that you will know when to stop (gambling, riding the bus, owning a stock, . . .).

In the gambling example from two paragraphs ago, we could take $T^* = \inf\{n : R_n \in \mathbb{R} \setminus (0, 1000)\}$, which could also be written as $T^* = \min(T_0, T_1)$, where $T_0 = \inf\{n : R_n \leq 0\}$ and $T_1 = \inf\{n : R_n \geq 1000\}$. All three of T_0, T_1 and T^* are stopping times. An example of a non-stopping time would be this: “I’ll play for 100 rounds and stop whenever my bankroll is largest”. This corresponds to the random variable $T_2 = \arg \max(R_n, 0 \leq n \leq 100)$; ¹⁹ but to stop at time T_2 would require foreknowledge of $(R_n, T_2 \leq n \leq 100)$. Laws against insider trading are in a sense legislating that decisions about when to buy and sell stocks must be stopping times.

Exercise 13.2. Show that T_0, T_1 and T^* defined above are all stopping times with respect to the filtration $\mathcal{F}_n = \sigma(R_m, 0 \leq m \leq n)$.

Given an (\mathcal{F}_n) -stopping time T , we define the *stopped σ -field* as follows: let $\mathcal{F}_\infty = \sigma(\bigcup_{n \geq 0} \mathcal{F}_n)$, and let

$$\mathcal{F}_T := \{A \in \mathcal{F}_\infty : \forall n \geq 0, A \cap \{T \leq n\} \in \mathcal{F}_n\}.$$

For example, the event A that (R_n) first exceeds 1000 before first reaching 0 is in \mathcal{F}_{T^*} , since (informally) at time T^* we know which of 0 and 1000 was first reached by (R_n) . The next exercise is a special case of a fact from the subsequent proposition, but is perhaps worth doing separately to make sure you’re comfortable with these basic ideas.

Exercise 13.3. Let $A = \{T_1 \leq T_0\}$. Show that A is in $\mathcal{F}_{T_0}, \mathcal{F}_{T_1}$, and \mathcal{F}_{T^*} .

Proposition 13.1 (Basic facts about stopping times). Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space, let $(X_n)_{n \geq 0}$ be an (\mathcal{F}_n) -adapted process, and let S, T be two (\mathcal{F}_n) -stopping times. Then the following all hold.

- (1) $\min(S, T)$ is a stopping time.
- (2) \mathcal{F}_T is a σ -field.
- (3) If $S \leq T$ then $\mathcal{F}_S \subseteq \mathcal{F}_T$.
- (4) $X_T \mathbf{1}_{\{T < \infty\}}$ is $\mathcal{F}_T / \mathcal{B}(\mathbb{R})$ -measurable.
- (5) The process $(X_{\min(T, n)}, n \geq 0)$ is (\mathcal{F}_n) -adapted.²⁰
- (6) If $(X_n)_{n \geq 0}$ is integrable then $(X_{\min(T, n)})_{n \geq 0}$ is integrable.

The proofs of the facts stated in the proposition are left as **exercises**.

An (\mathcal{F}_n) -adapted integrable process $(X_n)_{n \geq 0}$ is a *supermartingale* if $\mathbf{E}\{X_n \mid \mathcal{F}_m\} \stackrel{\text{a.s.}}{\leq} X_m$ for all $0 \leq m \leq n$. It is a *submartingale* if $\mathbf{E}\{X_n \mid \mathcal{F}_m\} \stackrel{\text{a.s.}}{\geq} X_m$ for all $0 \leq m \leq n$. You might expect the inequalities to go the other way; in its current form it is more in line with the definitions of sub/superharmonic functions, but this is hard to explain rigorously without a large digression. So for the time being you’ll just have to find your own way to remember.

Submartingale,
supermartingale

Theorem 13.2 (Optional stopping theorem). Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space and let $(X_n)_{n \geq 0}$ be an \mathcal{F}_n -supermartingale. Then for any bounded stopping times $0 \leq S \leq T$, it holds that $\mathbf{E}X_T \leq \mathbf{E}X_S$.

¹⁸Couper les cheveux en quatre (fr)/Fendre les cheveux en quatre (qc)/S’enfarger dans les fleurs du tapis (qc)

¹⁹Given a finite collection $(x_i, i \in I)$ of real numbers, $\arg \max(x_i, i \in I)$ returns the value of i for which x_i is largest.

²⁰Hairs unsplit.

Before proving the theorem, we note its immediate corollary for martingales.

Corollary 13.3. *Suppose $(X_n)_{n \geq 0}$ is in fact an \mathcal{F}_n -martingale. Then for any bounded stopping times $0 \leq S \leq T$, it holds that $\mathbf{E}X_T = \mathbf{E}X_S$.*

To prove the corollary, note that if $(X_n)_{n \geq 0}$ is a martingale then both $(X_n)_{n \geq 0}$ and $(-X_n)_{n \geq 0}$ are supermartingales; then apply Theorem 13.2. The optional stopping theorem is a consequence of the following theorem, which lists three necessary and sufficient conditions for an adapted integrable process to be a supermartingale.

Theorem 13.4. *Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space and let $(X_n)_{n \geq 0}$ be an (\mathcal{F}_n) -adapted integrable process. Then the following are equivalent.*

- (a) $(X_n)_{n \geq 0}$ is an (\mathcal{F}_n) -supermartingale.
- (b) For any bounded (\mathcal{F}_n) -stopping time T and any (\mathcal{F}_n) -stopping time S , $\mathbf{E}\{X_T \mid \mathcal{F}_S\} \stackrel{\text{a.s.}}{\leq} X_{\min(S,T)}$.
- (c) For any (\mathcal{F}_n) -stopping time T , the process $(X_{\min(T,n)})_{n \geq 0}$ is an (\mathcal{F}_n) -supermartingale.
- (d) For any bounded (\mathcal{F}_n) -stopping times S and T with $S \stackrel{\text{a.s.}}{\leq} T$, $\mathbf{E}X_T \leq \mathbf{E}X_S$.

Proof. [(c) \Rightarrow (a)]. Let T be the constant function which is identically equal to n . Then $X_{\min(T,n)} = X_n$ and $X_{\min(T,n-1)} = n - 1$, so by the assumption in (c),

$$\mathbf{E}\{X_n \mid \mathcal{F}_{n-1}\} = \mathbf{E}\{X_{\min(T,n)} \mid \mathcal{F}_{n-1}\} \stackrel{\text{a.s.}}{\leq} X_{\min(T,n-1)} = X_{n-1},$$

so X is an (\mathcal{F}_n) -supermartingale.

[(b) \Rightarrow (c)]. Let T be a stopping time, fix $n \geq 1$ and take $S \equiv n - 1$. Then $\mathcal{F}_S = \mathcal{F}_{n-1}$ (exercise), and $\min(T, n)$ is a bounded stopping time, so by (b),

$$\mathbf{E}\{X_{\min(T,n)} \mid \mathcal{F}_{n-1}\} = \mathbf{E}\{X_{\min(T,n)} \mid \mathcal{F}_S\} \stackrel{\text{a.s.}}{\leq} X_{\min(\min(T,n),S)} = X_{\min(T,n-1)}, \stackrel{\text{a.s.}}{\leq}$$

so $(X_{\min(T,n)})_{n \geq 0}$ is an (\mathcal{F}_n) -supermartingale.

[(b) \Rightarrow (d)]. If S and T are both bounded stopping times with $S \stackrel{\text{a.s.}}{\leq} T$ then (b) gives us that

$$\mathbf{E}\{X_T \mid \mathcal{F}_S\} \stackrel{\text{a.s.}}{\leq} X_{\min(S,T)} \stackrel{\text{a.s.}}{=} X_S.$$

Taking expectations on both sides, it follows that $\mathbf{E}X_T \leq \mathbf{E}X_S$.

[(a) \Rightarrow (b)]. Suppose $(X_n)_{n \geq 0}$ is a supermartingale, let S be a stopping time and T be a bounded stopping time, and choose $n \in \mathbb{N}$ such that $\mathbf{P}\{T \leq n\} = 1$. Then we can write

$$X_T = X_{\min(S,T)} + \sum_{k=0}^n (X_{k+1} - X_k) \mathbf{1}_{[S \leq k < T]}.$$

For any event $A \in \mathcal{F}_S$ and $k \in \mathbb{N}$, by definition we have $A \cap \{S \leq k\} \in \mathcal{F}_k$. Also, $\{T \leq k\} \in \mathcal{F}_k$ so $\{T > k\} \in \mathcal{F}_k$. Using this measurability together with the tower law, and then using supermartingale property, it follows that

$$\begin{aligned} \mathbf{E}[(X_{k+1} - X_k) \mathbf{1}_{[S \leq k < T]} \mathbf{1}_{[A]}] &= \mathbf{E}[\mathbf{E}\{(X_{k+1} - X_k) \mathbf{1}_{[T > k]} \mathbf{1}_{[A \cap \{S \leq k\}]} \mid \mathcal{F}_k\}] \\ &= \mathbf{E}[\mathbf{E}\{X_{k+1} - X_k \mid \mathcal{F}_k\} \mathbf{1}_{[T > k]} \mathbf{1}_{[A \cap \{S \leq k\}]}] \\ &\leq \mathbf{E}[0 \cdot \mathbf{1}_{[T > k]} \mathbf{1}_{[A \cap \{S \leq k\}]}] \\ &= 0. \end{aligned}$$

Combined with the previous identity for X_T , it follows that for any event $A \in \mathcal{F}_S$,

$$\mathbf{E}[X_T \mathbf{1}_{[A]}] \leq \mathbf{E}[X_{\min(S,T)} \mathbf{1}_{[A]}].$$

From this and the definition of conditional expectation, it follows that

$$\mathbf{E}[\mathbf{E}\{X_T \mid \mathcal{F}_S\} \mathbf{1}_{[A]}] = \mathbf{E}[X_T \mathbf{1}_{[A]}] \leq \mathbf{E}[X_{\min(S,T)} \mathbf{1}_{[A]}].$$

Since both $\mathbf{E}\{X_T \mid \mathcal{F}_S\}$ and $X_{\min(S,T)}$ are $\mathcal{F}_S/\mathcal{B}(\mathbb{R})$ -measurable, it follows that $\mathbf{E}\{X_T \mid \mathcal{F}_S\} \stackrel{\text{a.s.}}{\leq} X_{\min(S,T)}$, so (b) holds.

[(d) \Rightarrow (a)]. We must show that for all $n \geq 1$

$$\mathbf{E}\{X_n \mid \mathcal{F}_{n-1}\} \stackrel{\text{a.s.}}{\leq} X_{n-1}$$

To establish this inequality, it suffices to show that for any event $A \in \mathcal{F}_{n-1}$,

$$\mathbf{E}[\mathbf{E}\{X_n \mid \mathcal{F}_{n-1}\} \mathbf{1}_{[A]}] \leq \mathbf{E}[X_{n-1} \mathbf{1}_{[A]}].$$

So fix $n \geq 1$ and $A \in \mathcal{F}_{n-1}$. Let $T \equiv n$ and let $S = (n-1)\mathbf{1}_{[A]} + n\mathbf{1}_{[A^c]}$. It is not hard to see that S is a stopping time **exercise**. Since $S \leq T$, it follows from (d) that

$$\mathbf{E}X_n = \mathbf{E}X_T \leq \mathbf{E}X_S.$$

But $X_S = X_{n-1}\mathbf{1}_{[A]} + X_n\mathbf{1}_{[A^c]}$, so this gives

$$\mathbf{E}X_n \leq \mathbf{E}[X_{n-1}\mathbf{1}_{[A]}] + \mathbf{E}[X_n\mathbf{1}_{[A^c]}];$$

rearranging gives $\mathbf{E}[X_n\mathbf{1}_{[A]}] \leq \mathbf{E}[X_{n-1}\mathbf{1}_{[A]}]$. But by definition, $\mathbf{E}[X_n\mathbf{1}_{[A]}] = \mathbf{E}[\mathbf{E}\{X_n \mid \mathcal{F}_{n-1}\} \mathbf{1}_{[A]}]$, so the required inequality follows. \square

13.1. Martingale convergence theorems. Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space, and let $X = (X_n, n \geq 0)$ be an \mathcal{F}_n -martingale. Then for any bounded stopping time T , the optimal stopping theorem tells us that $\mathbf{E}X_T = \mathbf{E}X_0$. In gambling or stock market terminology, this says that we should not expect to increase our initial fortune during any fixed time window (our lifetime, say).

Here is a potential counterargument: what about “buy low, sell high”, the oldest trading principle in the book? In other words, what prevents a trader from fixing two values $a < b$, then buying whenever the stock price dips below a , and selling whenever the price rises above b ? If it isn't given too much thought, this strategy seems pretty good. Let $T_0 = 0$, and for $k \geq 0$ let $S_{k+1} = \inf\{m \geq T_k : X_m < a\}$ and $T_{k+1} = \inf\{m \geq S_{k+1} : X_m > b\}$. For $i \geq 1$ the interval $[S_i, T_i]$ is the i 'th *upcrossing* of the interval $[a, b]$ by the martingale X . Our hypothetical trader earns $b - a$ for each upcrossing they can apply their strategy to.

upcrossing

The next theorem is the reason why this strategy isn't as great as it looks. For $n \geq 0$ write $U_n[a, b] = \max\{i : T_i \leq n\}$ for the number of upcrossings of $[a, b]$ by X by time n , and let $U[a, b] = \lim_{n \rightarrow \infty} U_n[a, b]$ be the total number of upcrossings of $[a, b]$ by X .

Theorem 13.5 (Doob's upcrossing inequality). *Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space, and let $X = (X_n, n \geq 0)$ be an \mathcal{F}_n -supermartingale, and fix $a < b \in \mathbb{R}$. Then*

$$\mathbf{E}[U[a, b]] \leq \frac{1}{b-a} \sup_{k \geq 0} \mathbf{E}[(X_k - a)^-] \leq \frac{1}{b-a} \sup_{k \geq 0} (\mathbf{E}X_k^- + a).$$

Proof. First, note that the number of upcrossings of $[a, b]$ by X is the same as the number of upcrossings of $[0, b-a]$ by $X - a := (X_n - a, n \geq 0)$; we can thus assume $a = 0$, and must then prove that

$$\mathbf{E}[U[a, b]] \leq \frac{1}{b} \sup_{n \geq 0} \mathbf{E}[X_n^-].$$

The idea is that $X_{T_i} - X_{S_i} \geq b$ for all i , so $\sum_{i \geq 1} X_{T_i} - X_{S_i} \geq bU[0, b]$. On the other hand, X is a supermartingale and $S_i \leq T_i$, so the expectations of the terms $X_{T_i} - X_{S_i}$ should be non-positive. This doesn't quite work as written because the stopping times S_i and T_i need not be bounded. To make the idea work, we need to “localize” – work with the bounded stopping times $S_i \wedge n$ and $T_i \wedge n$ – and then taking a limit.

For any $n \in \mathbb{N}$, by the optional stopping theorem we do have $\mathbf{E}[X_{T_i \wedge n} - X_{S_i \wedge n}] \leq 0$. Below write $m = U_n[0, b] = \max i : T_i \leq n$ for the number of upcrossings completed before time n . For $i \leq m$ we have

$$X_{T_i \wedge n} - X_{S_i \wedge n} = X_{T_i} - X_{S_i} \geq b.$$

For $i > m + 1$ we have $S_i > T_{m+1} > n$ so

$$X_{T_i \wedge n} - X_{S_i \wedge n} = X_n - X_n = 0.$$

Similarly, if $i = m + 1$ and $S_i \geq n$ then $X_{T_i \wedge n} - X_{S_i \wedge n} = 0$. Finally, if $i = m + 1$ and $S_i \leq n$ then since $X_{S_i} \leq 0$ we have (draw a picture!)

$$X_{T_i \wedge n} - X_{S_i \wedge n} = X_n - X_{S_i} \geq -X_n^-. \tag{13.1}$$

Combining the above bounds we obtain that

$$\sum_{i \geq 0} (X_{T_i \wedge n} - X_{S_i \wedge n}) \geq bU_n[0, b] - X_n^-;$$

taking expectations, the optimal stopping theorem then gives that $b\mathbf{E}[U_n[0, b]] \leq \mathbf{E}[X_n^-] \leq \sup_{k \geq 0} \mathbf{E}[X_k^-]$. The result now follows by the monotone convergence theorem. \square

Note that at (13.1) there was some flexibility how to bound $X_{T_i \wedge n} - X_{S_i \wedge n}$ from below. A natural approach would have been to simply throw away the term $-X_{S_i}$ and use X_n as a lower bound. You can check that this leads to the bound

$$\mathbf{E}[U[a, b]] \leq \frac{1}{b-a} \sup_{k \geq 0} \mathbf{E}[a - X_k].$$

Here is an easy corollary of Doob's upcrossing inequality.

Corollary 13.6. *Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space, and let $X = (X_n, n \geq 0)$ be an \mathcal{F}_n -martingale with $\sup_{k \geq 0} \mathbf{E}|X_k| < \infty$. Then $\mathbf{E}[U[a, b]] < \infty$ for all $a < b \in \mathbb{R}$.*

Proof. By the upcrossing inequality we have

$$(b-a)\mathbf{E}[U[a, b]] \leq \sup_{k \geq 0} (\mathbf{E}X_k^- + a) \leq \sup_{k \geq 0} (\mathbf{E}|X_k| + a) = a + \sup_{k \geq 0} \mathbf{E}|X_k| < \infty. \quad \square$$

This allows us to quickly prove our first martingale convergence theorem.

Theorem 13.7 (L_1 martingale convergence theorem). *Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space, and let $X = (X_n, n \geq 0)$ be an \mathcal{F}_n -martingale with $\sup_{k \geq 0} \mathbf{E}|X_k| < \infty$. Then there exists a $\mathcal{F}_\infty/\mathcal{B}(\mathbb{R})$ -measurable random variable X_∞ with $\mathbf{E}|X_\infty| \leq \sup_{k \geq 0} \mathbf{E}|X_k|$ such that $X_n \xrightarrow{\text{a.s.}} X_\infty$ as $n \rightarrow \infty$.*

The next exercise contains a straightforward analytic fact we will use in the course of the proof.

Exercise 13.4. *A sequence $(x_n, n \geq 0)$ of real numbers converges (possibly to $\pm\infty$) if and only if for all $a < b \in \mathbb{Q}$, the number of upcrossings of $[a, b]$ by $(x_n, n \geq 0)$ is finite.*

Proof of Theorem 13.7. For $a < b \in \mathbb{Q}$ let $E_{a,b}$ be the event that $U[a, b] < \infty$. By Corollary 13.6 and countable subadditivity we have

$$\mathbf{P} \left\{ \bigcap_{a < b \in \mathbb{Q}} E_{a,b} \right\} = 1,$$

so by Exercise 13.4 (a), the sequence of random variables $(X_n, n \geq 0)$ converges almost surely. Write $X = \liminf_{n \rightarrow \infty} X_n$; then $X_n \xrightarrow{\text{a.s.}} X$. Moreover, since each X_n is $\mathcal{F}_\infty/\mathcal{B}(\mathbb{R})$ -measurable, it follows from Exercise 5.2 (b) that X is $\mathcal{F}_\infty/\mathcal{B}(\mathbb{R}^*)$ -measurable.

Since $X_n \xrightarrow{\text{a.s.}} X$, also $|X_n| \xrightarrow{\text{a.s.}} |X|$, so by Fatou's lemma,

$$\mathbf{E}|X| \leq \liminf_{n \rightarrow \infty} \mathbf{E}|X_n| \leq \sup_{k \geq 0} \mathbf{E}|X_k| < \infty,$$

so $|X|$ is a.s. finite. Taking $X_\infty = X\mathbf{1}_{[|X|<\infty]}$, we then have that $X_n \xrightarrow{\text{a.s.}} X_\infty$ as well, and X_∞ is $\mathcal{F}_\infty/\mathcal{B}(\mathbb{R})$ -measurable. \square

Here is a quite important special case.

Corollary 13.8 (Non-negative martingale convergence theorem). *Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space, and let $\mathbf{X} = (X_n, n \geq 0)$ be a non-negative \mathcal{F}_n -martingale. Then there exists a non-negative $\mathcal{F}_\infty/\mathcal{B}(\mathbb{R})$ -measurable random variable X_∞ with $\mathbf{E}X_\infty \leq \mathbf{E}X_0$ such that $X_n \xrightarrow{\text{a.s.}} X_\infty$ as $n \rightarrow \infty$.*

Proof. Since $X_n \geq 0$ almost surely for all n , necessarily the limit in the martingale convergence theorem is almost surely non-negative. For the expectation bound, simply note that since all the random variables are non-negative, we have $\sup_{n \geq 0} \mathbf{E}|X_n| = \sup_{n \geq 0} \mathbf{E}X_n = \mathbf{E}X_0$, where we have used the martingale property for the last equality. \square

Examples.

1. (Branching processes.) Will be developed at some length, below.

2. (Simple random walk stopped at zero.) Let $(S_n, n \geq 0)$ be a symmetric simple random walk started from $S_0 = 1$. Let $N = \min\{n : S_n \leq 0\}$ and let $M_n = S_{N \wedge n}$. Then $(M_n, n \geq 0)$ is a non-negative martingale so has an almost sure limit M_∞ with $\mathbf{E}M_\infty \leq 1$.

An integer-valued sequence which converges to a finite value is eventually constant, so it must be that $M_n = M_\infty$ for all sufficiently large n . This implies that $M_\infty \stackrel{\text{a.s.}}{=} 0$, since if $M_n \neq 0$ then $|M_{n+1} - M_n| = 1$. So $M_n \xrightarrow{\text{a.s.}} 0$; however, $\mathbf{E}M_n = \mathbf{E}M_0 = 1$ for all n .

3. (Expected boundary value of a Markov chain.) Fix a finite or countable set V . A *Markov chain* with state space V is a sequence of V -valued random variables $(X_n, n \geq 0)$ with the property that for all $n \geq 0$ and any sequence $(v_0, \dots, v_n, v_{n+1})$ of elements of V ,

$$\mathbf{P}\{X_{n+1} = v_{n+1} \mid X_i = v_i, i \leq n\} = \mathbf{P}\{X_{n+1} = v_{n+1} \mid X_n = v_n\}.$$

In other words, at each time n we can specify a V -by- V matrix of *time- n transition probabilities* $(P_n(u, v))_{u, v \in V}$, and whatever the values (v_0, \dots, v_{n+1}) above we will have

$$\mathbf{P}\{X_{n+1} = v_{n+1} \mid X_i = v_i, i \leq n\} = P_n(v_n, v_{n+1}).$$

The matrices P_n must have all row-sums equal to one for this to make sense: $\sum_{w \in V} P_n(v, w) = 1$. By far the most commonly studied Markov chains are *time-homogeneous*: there is a single transition probability matrix P such that $P_n = P$ for all $n \geq 0$.

Now fix a *boundary* $S \subset V$, and boundary values $b : S \rightarrow \mathbb{R}$. Let $\mathbf{X} = (X_n, n \geq 0)$ be a time-homogeneous Markov chain with transition matrix P , and let $\tau = \inf\{n \geq 0 : X_n \in S\}$. If we think of $b(v)$ as a reward (or penalty) associated with the state $v \in S$, then a natural way to extend b to all of V is

$$b(v) := \mathbf{E}\{b(X_\tau)\mathbf{1}_{[\tau < \infty]} \mid X_0 = v\};$$

the is the *expected reward* starting from state v . (As usual, $b(X_\tau)\mathbf{1}_{[\tau < \infty]}$ is to be interpreted as taking the value 0 if $\tau = \infty$.) Note that

$$\begin{aligned} b(v) &= \sum_{w \in V} \mathbf{E}\{b(X_\tau)\mathbf{1}_{[\tau < \infty]} \mid X_0 = v, X_1 = w\} \mathbf{P}\{X_1 = w \mid X_0 = v\} \\ &= \sum_{w \in V} \mathbf{E}\{b(X_\tau)\mathbf{1}_{[\tau < \infty]} \mid X_0 = w\} P(v, w) \\ &= \sum_{w \in V} P(v, w)b(w); \end{aligned}$$

the expected reward starting from v is a weighted average of the expected rewards of other sites, weighted according to the likelihood of moving to those sites.

Set $M_n = b(X_{n \wedge \tau})$. This is adapted to the filtration $(\mathcal{F}_n)_{n \geq 0}$ generated by \mathbf{X} . Moreover,

$$\mathbf{E}\{M_{n+1} \mid \mathcal{F}_n\} = \begin{cases} b(X_{n \wedge \tau}) & \text{if } \tau \leq n \\ \mathbf{E}\{b(X_{n+1}) \mid \mathcal{F}_n\} & \text{if } T > n. \end{cases}$$

But

$$\mathbf{E} \{ b(X_{n+1}) \mid X_n = v \} = \sum_{w \in V} P(v, w) b(w) = b(v),$$

so on the event $\tau > n$ we have $\mathbf{E} \{ b(X_{n+1}) \mid \mathcal{F}_n \} = b(X_n) = b(X_{n \wedge \tau}) = b(M_n)$. When $\tau \leq n$ we have $b(X_{n \wedge \tau}) = b(M_n)$, so the above identities show that $(M_n, n \geq 0)$ is an \mathcal{F}_n -martingale.

4. Product of IID. Let $(X_n, n \geq 1)$ be IID non-negative mean-one random variables in $L_1(\Omega, \mathcal{F}, \mathbf{P})$. Set $M_0 = 1$ and for $n \geq 1$ let $M_n = \prod_{i=1}^n X_i$. Then $(M_n, n \geq 1)$ is a martingale with respect to the filtration generated by $(X_n, n \geq 1)$, and $\mathbf{E}|M_n| = \mathbf{E} \prod_{i=1}^n X_i = 1$ for all $n \geq 1$, so M_n has an almost sure limit M_∞ with $\mathbf{E}M_\infty \leq 1$, by the martingale convergence theorem.

Let $Y_n = \log X_n$. If $Y_n \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ then it follows by the law of large numbers that $n^{-1} \log M_n = n^{-1} \sum_{i=1}^n Y_i \rightarrow \mathbf{E}Y_1$. Jensen's inequality tells us that $\mathbf{E}Y_1 \leq \log \mathbf{E}X_1 = 0$, and the inequality is strict unless X_1 is a.s. constant. Provided that X_1 are not a.s. constant, it follows that $\log M_n \rightarrow -\infty$, so $M_\infty \stackrel{\text{a.s.}}{=} 0$.

Exercise 13.5. Let $(X_n, n \geq 1)$ be IID non-negative mean-one random variables in $L_1(\Omega, \mathcal{F}, \mathbf{P})$. Suppose that $\log X_1$ is not in $L_1(\Omega, \mathcal{F}, \mathbf{P})$. Show that $n^{-1} \log \prod_{i=1}^n X_i \rightarrow -\infty$ almost surely.

Exercise 13.6. Let $(X_n, n \geq 0)$ be a time-homogeneous Markov chain with finite state space V and transition probability matrix $P = (P(u, v))_{u, v \in V}$. Fix a boundary $S \subset V$ and boundary values $b : S \rightarrow (0, \infty)$, and extend b to V as in the above example. Say that P is irreducible if for any $A \subset V$ there exist $u \in A$ and $v \in V \setminus A$ such that $P(u, v) > 0$.

- (a) Suppose P is irreducible. Prove that $b(u) > 0$ for all $u \in V$.
- (b) Write $T = \{v \in V : b(v) = \max\{b(w) : w \in V\}\}$ for the sites with maximum expected reward. Prove that $T \cap S$ is non-empty.
- (c) Suppose that for any $A \subset V \setminus S$ there exist $u \in A$ and $v \in (V \setminus S) \setminus A$ such that $P(u, v) > 0$. Prove that under this condition, either $b_{V \setminus S}$ is a constant function or else $T \subset S$. This is a discrete version of the maximum principle for harmonic functions.

Define a matrix $Q = Q(P, b)$ as follows: for $u, v \in V$ let

$$Q(u, v) = \begin{cases} 1 & \text{if } u = v \text{ and } b(u) = 0 \\ 0 & \text{if } u \neq v \text{ and } b(u) = 0 \\ 1 & \text{if } u = v \text{ and } u \in S \\ 0 & \text{if } u \neq v \text{ and } u \in S \\ \frac{P(u, v)b(v)}{b(u)} & \text{otherwise.} \end{cases}$$

- (d) Show that $Q = Q(P, b)$ is a transition probability matrix.
- (e) More to come. The idea is to iterate this and see what happens (extend b to V using Q instead of P , then repeat).

Uniform integrability. We've now seen a couple of examples of martingales which converge almost surely but which do not converge in expectation. The missing ingredient for convergence in expectation is *uniform integrability*. Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space and let $X = (X_i, i \in I)$ be $(\mathcal{F}/\mathcal{B}(\mathbb{R}))$ -measurable random variables. We say the collection F is *uniformly integrable* if

$$\lim_{M \rightarrow \infty} \sup_{i \in I} \mathbf{E} [|X_i| \mathbf{1}_{|X_i| > M}] = 0.$$

Exercise 13.7.

- (a) Prove that if $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ then $\lim_{M \rightarrow \infty} \mathbf{E} [|X| \mathbf{1}_{|X| > M}] = 0$.
- (b) Prove that if $X \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ then for any $\epsilon > 0$ there is $\delta = \delta(\epsilon, X) > 0$ such that for all $B \in \mathcal{F}$ with $\mathbf{P}(B) < \delta$,

$$\mathbf{E} [X \mathbf{1}_{[B]}] := \int_B X d\mathbf{P} < \epsilon.$$

(c) Prove that $X = (X_i, i \in I)$ is uniformly integrable if and only if (i) $\sup_{i \in I} \mathbf{E}|X_i| < \infty$ and (ii) for all $\epsilon > 0$ there is $\delta = \delta(\epsilon) > 0$ such that for all $B \in \mathcal{F}$ with $\mathbf{P}(B) < \delta$,

$$\mathbf{E}[X\mathbf{1}_{[B]}] < \epsilon.$$

(d) Show by example that neither (i) nor (ii) in part (c) implies the other.

Proposition 13.9. Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space and let $(X_n, n \geq 1)$ be uniformly integrable random variables over $(\Omega, \mathcal{F}, \mathbf{P})$. If $X_n \xrightarrow{d} X_\infty$ then $X_\infty \in L_1(\Omega, \mathcal{F}, \mathbf{P})$ and $\mathbf{E}X_n \rightarrow \mathbf{E}X_\infty$.

In the proof we'll use the notation introduced in the proof of Theorem 10.6: for a random variable Y and for real $C > 0$ we write $Y^{\leq C} := Y\mathbf{1}_{\{|Y| \leq C\}}$, and likewise define $Y^{< C}$, $Y^{\geq C}$ and $Y^{> C}$.

Proof. We first suppose that in fact $X_n \xrightarrow{\text{a.s.}} X$; we'll explain how to remove this assumption at the end.

Let M be large enough that $\sup_{n \geq 1} \mathbf{E}[|X_n|\mathbf{1}_{\{|X_n| > M\}}] < 1$. Then

$$\mathbf{E}|X_\infty| \leq \liminf_{n \geq 1} \mathbf{E}|X_n| \leq \sup_{n \geq 1} \mathbf{E}|X_n| \leq \sup_{n \geq 1} (\mathbf{E}[|X_n|\mathbf{1}_{\{|X_n| \leq M\}}] + 1) \leq M + 1,$$

so $X_\infty \in L_1(\Omega, \mathcal{F}, \mathbf{P})$.

Next, for any $C > 0$ with $\mathbf{P}\{|X_\infty| = C\} = 0$, we have $X_n^{\leq C} \xrightarrow{\text{a.s.}} X^{\leq C}$, so by the bounded convergence theorem $\mathbf{E}[|X_n^{\leq C} - X^{\leq C}|] \rightarrow 0$ as $n \rightarrow \infty$. For any such C , writing

$$|X_n - X_\infty| = |X_n^{\leq C} - X_\infty^{\leq C} + X_n^{> C} - X_\infty^{> C}| \leq |X_n^{\leq C} - X_\infty^{\leq C}| + |X_n^{> C}| + |X_\infty^{> C}|,$$

we obtain that

$$\limsup_{n \rightarrow \infty} \mathbf{E}[|X_n - X_\infty|] \leq \limsup_{n \rightarrow \infty} \mathbf{E}[|X_n^{> C}| + |X_\infty^{> C}|].$$

For any $\epsilon > 0$ we may find C with $\mathbf{P}\{|X_\infty| = C\} = 0$ large enough that $\sup_{1 \leq n \leq \infty} \mathbf{E}|X_n^{> C}| < \epsilon/2$, so that the preceding equation is less than ϵ . It follows that $\mathbf{E}[|X_n - X_\infty|] \rightarrow 0$ as $n \rightarrow \infty$, so $\mathbf{E}X_n \rightarrow \mathbf{E}X_\infty$.

This handles the case that $X_n \xrightarrow{\text{a.s.}} X_\infty$. In general, provided that $X_n \xrightarrow{d} X_\infty$ then by the Skorohod representation theorem, Theorem 5.10, there exists a coupling $(Y_n, 1 \leq n \leq \infty)$ of $(X_n, 1 \leq n \leq \infty)$ so that $Y_n \xrightarrow{\text{a.s.}} Y_\infty$. By the result already proved we then have that $\lim_{n \rightarrow \infty} \mathbf{E}X_n = \lim_{n \rightarrow \infty} \mathbf{E}Y_n = \mathbf{E}Y_\infty = \mathbf{E}X_\infty$. \square

Corollary 13.10 (UI martingale convergence theorem). Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space, and let $X = (X_n, n \geq 0)$ be a uniformly integrable \mathcal{F}_n -martingale. Then there exists $X_\infty \in L_1(\Omega, \mathcal{F}_\infty, \mathbf{P})$ with $\mathbf{E}X_\infty = \mathbf{E}X_0$ such that $X_n \xrightarrow{L_1} X_\infty$ as $n \rightarrow \infty$.

Proof. By the L_1 martingale convergence theorem, there exists $X_\infty \in L_1(\Omega, \mathcal{F}_\infty, \mathbf{P})$ such that $X_n \xrightarrow{\text{a.s.}} X_\infty$. We then also have $|X_n| \xrightarrow{\text{a.s.}} |X_\infty|$ and the family $(|X_n|, 1 \leq n \leq \infty)$ is also uniformly integrable, so by the proposition, $\mathbf{E}|X_n| \rightarrow \mathbf{E}|X_\infty|$. It follows by Exercise 12.8 that $X_n \xrightarrow{L_1} X_\infty$, which also implies that

$$|\mathbf{E}X_0 - \mathbf{E}X_\infty| = \lim_{n \rightarrow \infty} |\mathbf{E}X_n - \mathbf{E}X_\infty| \leq \limsup_{n \rightarrow \infty} \mathbf{E}|X_n - X_\infty| = 0.$$

\square

Corollary 13.11. Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space, let $X = (X_n, n \geq 0)$ be a uniformly integrable \mathcal{F}_n -martingale, and let X_∞ be the a.s. and L_1 limit of X . Then $\mathbf{E}\{X_\infty \mid \mathcal{F}_n\} \xrightarrow{\text{a.s.}} X_\infty$ for all $0 \leq n < \infty$. Moreover, for any $Y \in L_1(\Omega, \mathcal{F}_\infty, \mathbf{P})$, $\mathbf{E}\{Y \mid \mathcal{F}_n\} \rightarrow Y$ almost surely and in L_1 .

Proof. For any $0 \leq n \leq m$ we have $X_n \stackrel{\text{a.s.}}{=} \mathbf{E}\{X_m \mid \mathcal{F}_n\}$, so by linearity of expectation and Jensen's inequality (both conditional),

$$\begin{aligned} \|X_n - \mathbf{E}\{X_\infty \mid \mathcal{F}_n\}\|_1 &= \|\mathbf{E}\{X_m \mid \mathcal{F}_n\} - \mathbf{E}\{X_\infty \mid \mathcal{F}_n\}\|_1 \\ &= \|\mathbf{E}\{X_m - X_\infty \mid \mathcal{F}_n\}\|_1 \\ &= \mathbf{E}|\mathbf{E}\{X_m - X_\infty \mid \mathcal{F}_n\}| \\ &\stackrel{\text{a.s.}}{\leq} \mathbf{E}[\mathbf{E}\{|X_m - X_\infty| \mid \mathcal{F}_n\}] \\ &= \mathbf{E}|X_m - X_\infty|. \end{aligned}$$

Taking $m \rightarrow \infty$, the right-hand side tends to zero, from which it follows that $X_n \stackrel{\text{a.s.}}{=} \mathbf{E}\{X_\infty \mid \mathcal{F}_n\}$.

Next suppose that $Y \in L_1(\Omega, \mathcal{F}_\infty, \mathbf{P})$, and write $Y_n = \mathbf{E}\{Y \mid \mathcal{F}_n\}$ for $n \geq 0$. Then $(Y_n, n \geq 0)$ is a uniformly integrable martingale so there is $Y_\infty \in L_1(\Omega, \mathcal{F}_\infty, \mathbf{P})$ such that $Y_n \rightarrow Y_\infty$ almost surely and in L_1 . Moreover, by the first assertion of the corollary we know that $\mathbf{E}\{Y_\infty \mid \mathcal{F}_n\} \stackrel{\text{a.s.}}{=} Y_n$, so $\mathbf{E}\{Y_\infty \mid \mathcal{F}_n\} \stackrel{\text{a.s.}}{=} \mathbf{E}\{Y \mid \mathcal{F}_n\}$, for all $n \geq 0$.

Now write $\mathcal{P} = \bigcup_{n \geq 0} \mathcal{F}_n$; note that \mathcal{P} is a π -system generating \mathcal{F}_∞ . Fix $A \in \mathcal{P}$ and let $n \in \mathbb{N}$ be large enough that $A \in \mathcal{F}_n$. Using the definition of conditional expectation and the fact that $\mathbf{1}_{[A]}$ is $\mathcal{F}_n/\mathcal{B}(\mathbb{R})$ -measurable, we then have

$$\begin{aligned} \mathbf{E}[Y_\infty \mathbf{1}_{[A]}] &= \mathbf{E}[\mathbf{E}\{Y_\infty \mathbf{1}_{[A]} \mid \mathcal{F}_n\}] \\ &= \mathbf{E}[\mathbf{E}\{Y_\infty \mid \mathcal{F}_n\} \mathbf{1}_{[A]}] \\ &= \mathbf{E}[\mathbf{E}\{Y \mid \mathcal{F}_n\} \mathbf{1}_{[A]}] \\ &= \mathbf{E}[\mathbf{E}\{Y \mathbf{1}_{[A]} \mid \mathcal{F}_n\}] \\ &= \mathbf{E}[Y \mathbf{1}_{[A]}] \end{aligned}$$

Thus $\mathbf{E}[Y_\infty \mathbf{1}_{[A]}] = \mathbf{E}[Y \mathbf{1}_{[A]}]$ for all $A \in \mathcal{P}$. Exercise 6.5 now yields that $Y_\infty \stackrel{\text{a.s.}}{=} Y$. \square

This corollary is important. Proposition 12.4 already told us that for any random variable X and filtration $(\mathcal{F}_n, n \geq 0)$, the martingale $(\mathbf{E}\{X \mid \mathcal{F}_n\}, n \geq 0)$ is uniformly integrable. Corollary 13.11 tells us that every uniformly integrable martingale has this form. In other words, we have characterized the UI martingales: they are precisely those that may be obtained from a single L_1 random variable by taking conditional expectations along a filtration.

Uniform integrability gives us precisely the control we need to extend the optional stopping theorem to unbounded stopping times (i.e. stopping times which may take arbitrarily large values). This is the content of the *optional sampling theorem*.

Theorem 13.12 (Optional sampling theorem). *Let $X = (X_n, n \geq 0)$ be a uniformly integrable martingale relative to a filtration $(\mathcal{F}_n, n \geq 0)$. Then for any almost surely finite \mathcal{F}_n -stopping times $0 \leq S \leq T$, it holds that $\mathbf{E}\{X_T \mid \mathcal{F}_S\} \stackrel{\text{a.s.}}{=} X_S$, and $\mathbf{E}[X_T] = \mathbf{E}[X_0]$.*

Proof. The idea of the proof is to localize and then use Fatou's lemma; the uniform integrability is simply what we need for Fatou's lemma to be effective.

So fix stopping times $S \leq T$ as in the theorem statement. For any $n \geq 0$, by Theorem 13.4 (b) applied to X and to $-X$ (which are both supermartingales), we have that

$$\mathbf{E}\{X_{T \wedge n} \mid \mathcal{F}_S\} \stackrel{\text{a.s.}}{=} X_{(T \wedge n) \wedge S}.$$

Since T is almost surely finite, $X_{(T \wedge n) \wedge S} \stackrel{\text{a.s.}}{\rightarrow} X_{T \wedge S}$, and Fatou's lemma then gives that

$$\begin{aligned} \mathbf{E}|X_{T \wedge S} - \mathbf{E}\{X_T \mid \mathcal{F}_S\}| &= \mathbf{E}\left[\liminf_{n \rightarrow \infty} |X_{(T \wedge n) \wedge S} - \mathbf{E}\{X_T \mid \mathcal{F}_S\}|\right] \\ &\leq \liminf_{n \rightarrow \infty} \mathbf{E}\left[|X_{(T \wedge n) \wedge S} - \mathbf{E}\{X_T \mid \mathcal{F}_S\}|\right] \\ &= \liminf_{n \rightarrow \infty} \mathbf{E}\left[|\mathbf{E}\{X_{T \wedge n} \mid \mathcal{F}_S\} - \mathbf{E}\{X_T \mid \mathcal{F}_S\}|\right]. \end{aligned}$$

To prove the first assertion of the theorem, it thus suffices to establish that $\mathbf{E}\{X_{T \wedge n} \mid \mathcal{F}_S\} \xrightarrow{L_1} \mathbf{E}\{X_T \mid \mathcal{F}_S\}$, and we now turn to this.

By the UI martingale convergence theorem there exists $X_\infty \in L_1(\Omega, \mathcal{F}_\infty, \mathbf{P})$ such that $X_n \rightarrow X_\infty$ almost surely and in L_1 , and $\mathbf{E}\{X_\infty \mid \mathcal{F}_n\} \stackrel{\text{a.s.}}{=} X_n$ for all $n \geq 0$.

Next, note that that $\mathbf{P}\{X_{T \wedge n} = X_T\} \geq \mathbf{P}\{T \wedge n = T\} \rightarrow 1$ as $n \rightarrow \infty$, so $X_{T \wedge n} \xrightarrow{\text{a.s.}} X_T$ as $n \rightarrow \infty$. Moreover, since X is a martingale, by Theorem 13.4 (b) we have that

$$X_{T \wedge n} \stackrel{\text{a.s.}}{=} \mathbf{E}\{X_n \mid \mathcal{F}_{T \wedge n}\} \stackrel{\text{a.s.}}{=} \mathbf{E}\{\mathbf{E}\{X_\infty \mid \mathcal{F}_n\} \mid \mathcal{F}_{T \wedge n}\} \stackrel{\text{a.s.}}{=} \mathbf{E}\{X_\infty \mid \mathcal{F}_{T \wedge n}\},$$

the last equality holding by the tower law since $\mathcal{F}_{T \wedge n} \subset \mathcal{F}_T$. This implies that $(X_{T \wedge n}, n \geq 0)$ is a uniformly integrable martingale, so $X_{T \wedge n} \xrightarrow{L_1} X_T$ as $n \rightarrow \infty$. It follows that

$$\mathbf{E}[|\mathbf{E}\{X_{T \wedge n} \mid \mathcal{F}_S\} - \mathbf{E}\{X_T \mid \mathcal{F}_S\}|] \leq \mathbf{E}|X_{T \wedge n} - X_T| \rightarrow 0,$$

so $\mathbf{E}\{X_{T \wedge n} \mid \mathcal{F}_S\} \xrightarrow{L_1} \mathbf{E}\{X_T \mid \mathcal{F}_S\}$ as required.

Finally, since $T \wedge n$ is a bounded stopping time, by the optional stopping theorem we have $\mathbf{E}X_{T \wedge n} = \mathbf{E}X_0$. Taking $n \rightarrow \infty$ and using that $X_{T \wedge n} \xrightarrow{L_1} X_T$ gives that $\mathbf{E}X_T = \mathbf{E}X_0$. \square

13.2. Maximal inequalities and the L_p martingale convergence theorem. Given a sequence $X = (X_n, n \geq 0)$ of random variables defined on a common probability space, write $X^* = \sup(|X_n|, n \geq 0)$, and for $n \geq 0$ let $X_n^* = \sup(|X_k|, 0 \leq k \leq n)$. Doob's maximal inequalities provide a way to control the tail behaviour of X^* and X_n^* for martingales. The intuition is this: if X is a martingale then $(|X_n|, n \geq 0)$ is a submartingale. So if $|X_k|$ is large for some k , then by the submartingale property $|X_n|$ should will also be large in expectation for $n \geq k$; so the conditional expectation of $|X_n|$, given that X_n^* is large, should also be large.

Theorem 13.13 (Doob's maximal inequality). *Let $X = (X_n, n \geq 0)$ be a martingale or a non-negative submartingale defined on a probability space $(\Omega, \mathcal{F}, \mathbf{P})$. Then for all $\lambda > 0$,*

$$\lambda \cdot \mathbf{P}\{X^* \geq \lambda\} \leq \sup_{n \geq 0} \mathbf{E}|X_n|,$$

and for all $n \geq 0$,

$$\lambda \cdot \mathbf{P}\{X_n^* \geq \lambda\} \leq \sup_{0 \leq k \leq n} \mathbf{E}|X_k|.$$

Proof. First note that $X_n^* \uparrow X^*$ almost surely, so by the monotone convergence theorem $\lambda \mathbf{P}\{X_n^* \geq \lambda\} \rightarrow \lambda \mathbf{P}\{X^* \geq \lambda\}$. Moreover, $\sup_{0 \leq k \leq n} \mathbf{E}|X_k| \rightarrow \sup_{n \geq 0} \mathbf{E}|X_n|$. Together these two convergence facts yield that the first inequality is implied by the second, so it suffices to prove the second bound.

Next, note that if X is a martingale then $|X| = (|X_n|, n \geq 0)$ is a non-negative submartingale, by the conditional version of Jensen's inequality. Replacing X by $|X|$ doesn't change the value of X_n^* , so we may assume that X is itself a non-negative submartingale (this is just so we don't have to carry absolute value signs around).

Now fix $n \geq 0$, and let $T = n \wedge \inf(m \geq 0 : X_m \geq \lambda)$. We can bound $\mathbf{E}X_T$ from below as follows. Write $X_T = X_T(\mathbf{1}_{[X_n^* < \lambda]} + \mathbf{1}_{[X_n^* \geq \lambda]})$. If $X_n^* < \lambda$ then $T = n$ so $X_T \mathbf{1}_{[X_n^* < \lambda]} = X_n \mathbf{1}_{[X_n^* < \lambda]}$. If $X_n^* \geq \lambda$ then $\inf(m \geq 0 : X_m \geq \lambda) \leq n$, so $X_T \mathbf{1}_{[X_n^* \geq \lambda]} \geq \lambda \mathbf{1}_{[X_n^* \geq \lambda]}$. It follows that

$$\begin{aligned} \mathbf{E}X_T &= \mathbf{E}[X_T(\mathbf{1}_{[X_n^* < \lambda]} + \mathbf{1}_{[X_n^* \geq \lambda]})] \\ &\geq \mathbf{E}[X_n \mathbf{1}_{[X_n^* < \lambda]}] + \lambda \mathbf{P}\{X_n^* \geq \lambda\}. \end{aligned}$$

On the other hand, since $T \leq n$, by the optional stopping theorem $\mathbf{E}X_T \leq \mathbf{E}X_n$; combining this with the previous inequality and rearranging gives

$$\lambda \mathbf{P}\{X_n^* \geq \lambda\} \leq \mathbf{E}X_n - \mathbf{E}[X_n \mathbf{1}_{[X_n^* < \lambda]}] = \mathbf{E}[X_n \mathbf{1}_{[X_n^* \geq \lambda]}] \leq \mathbf{E}X_n.$$

Since X is a submartingale, $\mathbf{E}X_n = \sup_{0 \leq k \leq n} \mathbf{E}X_k$, so the result follows. \square

Corollary 13.14 (Doob's L_p inequality). *Under the hypotheses of Theorem 13.13, for all $p > 1$*

$$\|X^*\|_p \leq \frac{p}{p-1} \sup_{n \geq 0} \|X_n\|_p,$$

and for all $n \in \mathbb{N}$,

$$\|X_n^*\|_p \leq \frac{p}{p-1} \|X_n\|_p = \frac{p}{p-1} \sup_{k \leq n} \|X_k\|_p.$$

Proof. It again suffices to consider the case that X is a non-negative submartingale. The first bound again follows from the second by the monotone convergence theorem. Also, equality in the second display clearly holds, and we focus our attention on proving the first.

Fix $n \in \mathbb{N}$. For any $m \in \mathbb{N}$ and $x \geq 0$ we have $(x \wedge m)^p = \int_0^{x \wedge m} p\lambda^{p-1} d\lambda = \int_0^m p\lambda^{p-1} \mathbf{1}_{[x \geq \lambda]} d\lambda$, so by Fubini's theorem,

$$\begin{aligned} \|X_n^* \wedge m\|_p^p &= \mathbf{E} [(X_n^* \wedge m)^p] \\ &= \mathbf{E} \left[\int_0^m p\lambda^{p-1} \mathbf{1}_{[X_n^* \geq \lambda]} d\lambda \right] \\ &= \int_0^m p\lambda^{p-1} \mathbf{P} \{X_n^* \geq \lambda\} d\lambda. \end{aligned}$$

By Doob's maximal inequality the integrand is at most $p\lambda^{p-2} \mathbf{E} [X_n \mathbf{1}_{[X_n^* \geq \lambda]}]$, so by monotonicity and another application of Fubini's theorem we have

$$\begin{aligned} \|X_n^* \wedge m\|_p^p &\leq \mathbf{E} \left[\int_0^m p\lambda^{p-2} X_n \mathbf{1}_{[X_n^* \geq \lambda]} d\lambda \right] \\ &= \mathbf{E} \left[\frac{p}{p-1} X_n (X_n^* \wedge m)^{p-1} \right]. \end{aligned}$$

Finally, by Hölder's inequality,

$$\begin{aligned} \mathbf{E} [X_n (X_n^* \wedge m)^{p-1}] &\leq \|X_n\|_p \| (X_n^* \wedge m)^{p-1} \|_{p/(p-1)} \\ &= \|X_n\|_p \left(\mathbf{E} [((X_n^* \wedge m)^{p-1})^{p/(p-1)}] \right)^{(p-1)/p} \\ &= \|X_n\|_p \| (X_n^* \wedge m) \|_p^{p-1}. \end{aligned}$$

Combining the above bounds gives that $\|X_n^* \wedge m\|_p^p \leq \frac{p}{p-1} \|X_n\|_p$. Since X_n^* is non-negative, $X_n^* \wedge m \uparrow X_n^*$ as $m \rightarrow \infty$, so it follows by the monotone convergence theorem that $\|X_n^*\|_p^p \leq \frac{p}{p-1} \|X_n\|_p$. \square

The preceding inequality allows us to show that L_p -bounded martingales in fact converge in L_p .

Theorem 13.15 (L_p martingale convergence theorem). *Fix $p > 1$. Let $(\Omega, \mathcal{F}, (\mathcal{F}_n)_{n \geq 0}, \mathbf{P})$ be a filtered probability space and let $X = (X_n, n \geq 0)$ be an \mathcal{F}_n -martingale such that $\sup_{n \geq 0} \mathbf{E} [|X_n|^p] < \infty$. Then X is a uniformly integrable martingale, its a.s. limit X_∞ is in $L_p(\Omega, \mathcal{F}_\infty, \mathbf{P})$, and $X_n \xrightarrow{L_p} X_\infty$.*

Proof. First, for all $n \in \mathbb{N}$ and all $\lambda > 0$,

$$\mathbf{E} [|X_n| \mathbf{1}_{[|X_n| \geq \lambda]}] \leq \frac{1}{\lambda^{p-1}} \mathbf{E} [|X_n|^p \mathbf{1}_{[|X_n| \geq \lambda]}] \leq \frac{1}{\lambda^{p-1}} \sup_{n \geq 1} \mathbf{E} [|X_n|^p] < \infty,$$

so for any $\epsilon > 0$, if $\lambda^{p-1} \geq \sup_{n \geq 1} \mathbf{E} [|X_n|^p] / \epsilon$ then $\mathbf{E} [|X_n| \mathbf{1}_{[|X_n| \geq \lambda]}] \leq \epsilon$. It follows that X is uniformly integrable, so there is $X_\infty \in L_1(\Omega, \mathcal{F}_\infty, \mathbf{P})$ such that $X_n \rightarrow X_\infty$ almost surely and in L_1 .

To obtain the L_p convergence claimed in the theorem, write $X^* = \sup_{n \geq 0} |X_n|$ as before. For all $n \geq 0$,

$$|X_n - X_\infty|^p \leq (|X_n| + |X_\infty|)^p \leq (2X^*)^p.$$

Doob's L_p inequality gives

$$\mathbf{E} [(2X^*)^p] = 2^p \|X^*\|_p^p \leq 2^p \left(\frac{p}{p-1} \sup_{n \geq 0} \|X_n\|_p \right)^p < \infty,$$

so by the dominated convergence theorem it follows that

$$\lim_{n \rightarrow \infty} \mathbf{E} [|X_n - X_\infty|^p] = \mathbf{E} \left[\lim_{n \rightarrow \infty} |X_n - X_\infty|^p \right] = 0. \quad \square$$

Exercise 13.8 (UI, conditional expectations and L^p convergence). *Fix a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and a real number $p \geq 1$. Let S be the set of all sub- σ -fields of \mathcal{F} .*

- Prove that for any random variable $X \in L^p(\Omega, \mathcal{F}, \mathbb{P})$, the collection $\{\mathbf{E}[X | \mathcal{G}]^p : \mathcal{G} \in S\}$ is a UI family.*
- Prove that for any sequence $(X_n, n \geq 1)$ of non-negative random variables such that $(X_n^p, n \geq 0)$ is uniformly integrable, if $X_n \rightarrow X$ in probability for some random variable X , then also $X_n \xrightarrow{L^p} X$.*
- Under the hypotheses of the L_p martingale convergence theorem, prove that for any random variable $Y \in L_p(\Omega, \mathcal{F}_\infty, \mathbf{P})$ it holds that $\mathbf{E}\{Y | \mathcal{F}_n\} \xrightarrow{L_p} Y$ as $n \rightarrow \infty$.*

13.3. Filtrations and changes of measure.

Recall that if $\mu \ll \nu$ are two σ -finite measures on (Ω, \mathcal{F}) , the Radon-Nikodým derivative $X = d\mu/d\nu : \Omega \rightarrow [0, \infty)$ is the ν -a.e. unique $\mathcal{F}/\mathcal{B}(\mathbb{R})$ -measurable function X such that for all $E \in \mathcal{F}$,

$$\mu(E) = \int_E d\mu = \int_E X d\nu. \quad (13.2)$$

Some useful notation. Write $\mu = X\nu$ if X satisfies (13.2). Of course, this means X is a version of $d\mu/d\nu$; but proving the existence and uniqueness of that derivative (via martingales) is the point of this section.

Theorem 13.16 (Radon-Nikodým theorem). *Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space. Suppose that \mathbf{Q} is a finite measure on (Ω, \mathcal{F}) such that $\mathbf{Q} \ll \mathbf{P}$, in that for all $F \in \mathcal{F}$, if $\mathbf{P}\{F\} = 0$ then $\mathbf{Q}(F) = 0$. Then there exists a \mathbf{P} -a.s. unique, nonnegative random variable $X \in L^1(\Omega, \mathcal{F}, \mathbf{P})$ such that $\mathbf{Q} = X\mathbf{P}$.*

Note. We often use probabilistic “expectation” notation, writing e.g. $\mathbf{E}_{\mathbf{Q}}\{X\} := \int X d\mathbf{Q}$, even if \mathbf{Q} is not a probability measure.

We only prove Theorem 13.16 in the case that \mathcal{F} is separable, in that there exists a countable collection $\{F_n, n \geq 0\} \subset \mathcal{F}$ such that $\mathcal{F} = \sigma(\{F_n, n \geq 0\})$.

Write $\mathcal{F}_n = \sigma(\{F_1, \dots, F_n\})$ for $n \geq 0$, and $\mathbf{P}_n = \mathbf{P}|_{\mathcal{F}_n}$, $\mathbf{Q}_n = \mathbf{Q}|_{\mathcal{F}_n}$. Say a set E is an *atom* of \mathcal{F}_n if $E = G_1 \cap G_2 \cap \dots \cap G_n$ where each G_i is either F_i or F_i^c . There are at most 2^n atoms (the representations need not be unique). List the atoms as

$$A_{n,1}, \dots, A_{n,r(n)};$$

then each set of \mathcal{F}_n is a union of some collection of atoms.

Define $X_n : \Omega \rightarrow \mathbb{R}$ by setting, for $\omega \in A_{n,k}$

$$X_n(\omega) = \begin{cases} 0 & \text{if } \mathbf{P}\{A_{n,k} = 0\} \\ \frac{\mathbf{Q}(A_{n,k})}{\mathbf{P}(A_{n,k})} & \text{if } \mathbf{P}\{A_{n,k} > 0\}. \end{cases}$$

Lemma 13.17. *We have $\mathbf{Q}_n = X_n \mathbf{P}_n$.*

Proof. Clearly $X_n \in L^1(\Omega, \mathcal{F}, \mathbf{P})$, since X_n only takes finitely many values. Now fix any $F \in \mathcal{F}_n$; there is a unique representation of F as $F = \bigcup_{j=1}^k A_{n,i(j)}$, for some $k \leq r(n)$ and $i(1) <$

$i(2) \dots < i(k) \leq r(n)$. We then have

$$\begin{aligned} \mathbf{E}_{\mathbf{P}_n} \{X_n \mathbf{1}_{[F]}\} &= \mathbf{E}_{\mathbf{P}} \{X_n \mathbf{1}_{[F]}\} \\ &= \sum_{j=1}^k \mathbf{E}_{\mathbf{P}} \{X_n \mathbf{1}_{[A_{n,k}]}\} \\ &= \sum_{j \leq k: \mathbf{P}\{A_{n,k}\} > 0} \frac{\mathbf{Q}(A_{n,k})}{\mathbf{P}(A_{n,k})} \mathbf{E}_{\mathbf{P}} \{\mathbf{1}_{[A_{n,k}]}\} \\ &= \sum_{j \leq k: \mathbf{P}\{A_{n,k}\} > 0} \mathbf{Q}(A_{n,k}) \\ &= \mathbf{Q}(F) \\ &= \mathbf{Q}_n(F). \end{aligned}$$

Since $F \in \mathcal{F}_n$ was arbitrary, the result follows. □

It is straightforward that $(X_n, n \geq 0)$ is an \mathcal{F}_n -martingale for \mathbf{P} ; see Exercise 13.9, below. Moreover, it is non-negative, so defining $X_\infty = \limsup_{n \rightarrow \infty} X_n$, it follows that \mathbf{P} -almost surely $X_n \rightarrow X_\infty$.

Lemma 13.18. *For all $\epsilon > 0$ there exists $\delta > 0$ such that $\mathbf{P}\{F\} < \delta \Rightarrow \mathbf{Q}(F) < \epsilon$.*

Proof. Otherwise, there exists $\epsilon > 0$ and a sequence of sets F_n with $\mathbf{P}\{F_n\} \leq 2^{-n}$ such that $\mathbf{Q}(F_n) \geq \epsilon$ for all n . Let $F = \limsup F_n = \bigcap_{n \geq 1} \bigcup_{m \geq n} F_m$. Then by (reverse) Fatou's lemma, $\mathbf{Q}(F) \geq \limsup_n \mathbf{Q}(F_n) \geq \epsilon$. But $\sum_{n \geq 0} \mathbf{P}\{F_n\} < \infty$ so by the first Borel-Cantelli lemma, $\mathbf{P}\{F\} = 0$, contradicting the assumption that $\mathbf{Q} \ll \mathbf{P}$. □

Proof of Theorem 13.16, separable case. We first claim that $(X_n, n \geq 0)$ is a uniformly integrable martingale for \mathbf{P} . To see this, fix any $\epsilon > 0$ and let $\delta > 0$ be as in Lemma 13.18. Then for $K > \mathbf{Q}(\Omega)/\delta$, for all $n \geq 0$, since $\mathbf{Q}_n = X_n \mathbf{P}_n$ we have

$$\mathbf{P}\{X_n > K\} \leq \frac{\mathbf{E}_{\mathbf{P}}\{X_n\}}{K} = \frac{\mathbf{E}_{\mathbf{Q}}\{1\}}{K} = \frac{\mathbf{Q}(\Omega)}{K} < \delta.$$

Using again that $\mathbf{Q}_n = X_n \mathbf{P}_n$ gives

$$\mathbf{E}_{\mathbf{P}}\{X_n \mathbf{1}_{[X_n > K]}\} = \mathbf{E}_{\mathbf{Q}}\{\mathbf{1}_{[X_n > K]}\} = \mathbf{Q}\{X_n > K\} < \epsilon,$$

the last inequality by Lemma 13.18. This proves that $(X_n, n \geq 0)$ is UI; the martingale convergence theorem then gives that $X_n \rightarrow X$ in $L^1(\Omega, \mathcal{F}, \mathbf{P})$.

It follows that for all $F \in \bigcup_{n \geq 0} \mathcal{F}_n$,

$$\mathbf{E}_{\mathbf{P}}\{X \mathbf{1}_{[F]}\} = \lim_{n \rightarrow \infty} \mathbf{E}_{\mathbf{P}}\{X_n \mathbf{1}_{[F]}\} = \mathbf{Q}(F)$$

In other words, the measures $X\mathbf{P}$ and \mathbf{Q} agree on $\bigcup_{n \geq 0} \mathcal{F}_n$. Since $\bigcup_{n \geq 0} \mathcal{F}_n$ is a π -system generating \mathcal{F} , it follows that $X\mathbf{P} = \mathbf{Q}$ as claimed. □

The next theorem describes how the dichotomy between absolute continuity and mutual singularity of measures manifests when observed along a filtration.

Theorem 13.19. *Fix an increasing sequence of sub- σ -field $(\mathcal{F}_n)_{n \geq 1}$ with $\sigma(\bigcup_n \mathcal{F}_n) = \mathcal{F}$. Write $\mathbf{P}_n := \mathbf{P}|_{\mathcal{F}_n}$ and $\mathbf{Q}_n := \mathbf{Q}|_{\mathcal{F}_n}$. Suppose that $\mathbf{Q}_n \ll \mathbf{P}_n$ for all n , and write $X_n = d\mathbf{Q}_n/d\mathbf{P}_n : \Omega \rightarrow [0, \infty)$ for the corresponding Radon-Nikodym derivatives. Then setting $X = \limsup_{n \rightarrow \infty} X_n$, it holds that*

$$\mathbf{Q} = X\mathbf{P} + \mathbf{Q}\mathbf{1}_{[X = \infty]}. \tag{13.3}$$

Exercise 13.9. *In the notation of Theorem 13.19, show that $(X_n, n \geq 1)$ is an \mathcal{F}_n -martingale for \mathbf{P} .*

Remark. Since the X_n are non-negative, Exercise 13.9 implies that X_n converges \mathbf{P} -almost surely, so we must have $\mathbf{P} \{ \lim_{n \rightarrow \infty} X_n = X \} = 1$. But it is standard that if $Z_n \xrightarrow{\mathbf{P}} Z_\infty$ and $\mathbf{E}|Z_n| < \infty$ for all n , then $\mathbf{E}|Z_n| \rightarrow \mathbf{E}|Z|$ if and only if $(Z_n, n \geq 1)$ is uniformly integrable. (See Exercise 12.8.) This means that there is another equivalent property which may be added to (1) in Theorem 13.19: that $(X_n, n \geq 1)$ is \mathbf{P} -uniformly integrable.

Lemma 13.20. *In the setting of Theorem 13.19, if $\mathbf{Q} \ll \mathbf{P}$ then $\mathbf{Q} = X\mathbf{P}$.*

Recall that $\mathbf{Q} = X\mathbf{P}$ is shorthand for the statement that for all $E \in \mathcal{F}$,

$$\mathbf{Q}(E) = \int_E d\mathbf{Q} = \int_E X d\mathbf{P} = \mathbf{E}_{\mathbf{P}} \{ X \mathbf{1}_{[E]} \} .$$

Proof. First suppose \mathbf{Q} is absolutely continuous with respect to \mathbf{P} . Then the Radon-Nikodým derivative $\tilde{X} = d\mathbf{Q}/d\mathbf{P}$ exists and satisfies $\mathbf{Q}(\tilde{X} = \infty) = 0$, so we just want to show that for all $E \in \mathcal{F}$,

$$\mathbf{Q}(E) = \mathbf{E}_{\mathbf{P}} \{ X \mathbf{1}_{[E]} \} .$$

For all $E \in \mathcal{F}$, by the definition of the Radon-Nikodým derivative,

$$\mathbf{E}_{\mathbf{Q}} \{ \mathbf{1}_{[E]} \} = \int_E 1 d\mathbf{Q} = \int_E \tilde{X} d\mathbf{P} = \mathbf{E}_{\mathbf{P}} \{ \tilde{X} \mathbf{1}_{[E]} \} . \quad (13.4)$$

For all $E \in \mathcal{F}_n$, we also have

$$\begin{aligned} \mathbf{E}_{\mathbf{P}} \{ X_n \mathbf{1}_{[E]} \} &= \int_E X_n d\mathbf{P} \\ &= \int_E X_n d\mathbf{P}_n && \text{(Homework)} \\ &= \int_E 1 d\mathbf{Q}_n && \text{(Since } X_n = d\mathbf{Q}_n/d\mathbf{P}_n) \\ &= \int_E 1 d\mathbf{Q} \\ &= \mathbf{E}_{\mathbf{Q}} \{ \mathbf{1}_{[E]} \} , \end{aligned}$$

so X_n is a version of $\mathbf{E} [\tilde{X} | \mathcal{F}_n]$. Since $\mathcal{F}_\infty := \sigma(\bigcup_{n \rightarrow \infty} \mathcal{F}_n) = \mathcal{F}$, the non-negative martingale convergence theorem then gives that \mathbf{P} -almost surely

$$X_n \rightarrow \mathbf{E} [\tilde{X} | \mathcal{F}_\infty] = \mathbf{E} [\tilde{X} | \mathcal{F}] = \tilde{X} .$$

But $X = \limsup_{n \rightarrow \infty} X_n$ by definition, so \mathbf{P} -almost surely $X = \tilde{X}$. Thus $\mathbf{E}_{\mathbf{P}} \{ \tilde{X} \mathbf{1}_{[E]} \} = \mathbf{E}_{\mathbf{P}} \{ X \mathbf{1}_{[E]} \}$, and the result follows from (13.4). \square

Proof of Theorem 13.19. Let π be the average of \mathbf{P} and \mathbf{Q} , so $\pi(E) = (\mathbf{P}(E) + \mathbf{Q}(E))/2$ for $E \in \mathcal{F}$. For $n \geq 1$ let $\pi_n = \pi|_{\mathcal{F}_n} = (\mathbf{P}_n + \mathbf{Q}_n)/2$. Then both \mathbf{P} and \mathbf{Q} are absolutely continuous with respect to π , and likewise \mathbf{P}_n and \mathbf{Q}_n are absolutely continuous with respect to π_n for all $n \geq 1$.

Write $U_n = d\mathbf{Q}_n/d\pi_n$ and $V_n = d\mathbf{P}_n/d\pi_n$ and let $U = \limsup_{n \rightarrow \infty} U_n$ and $V = \limsup_{n \rightarrow \infty} V_n$. Since $\mathbf{Q} \ll \pi$ it follows by Lemma 13.20 (applied with π in place of \mathbf{P}) that π -almost surely $U_n \rightarrow U$ and that $\mathbf{Q} = U\pi$. Likewise, applying Lemma 13.20 with π in place of \mathbf{P} and \mathbf{P} in place of \mathbf{Q} , it follows that π -almost surely $V_n \rightarrow V$ and that $\mathbf{P} = V\pi$.

Next, π -almost surely we have

$$U_n + V_n = \frac{d\mathbf{Q}_n}{d\pi_n} + \frac{d\mathbf{P}_n}{d\pi_n} = 2 \frac{d\pi_n}{d\pi_n} = 2 .$$

It follows that

$$\pi(U + V = 0) = \pi(\limsup_n (U_n + V_n) = 0) = 0 ,$$

so π -almost surely, U/V is well-defined (and equal to ∞ if $U = \infty$ and $V = 0$), and

$$\begin{aligned} \frac{U}{V} &= \frac{\lim_{n \rightarrow \infty} U_n}{\lim_{n \rightarrow \infty} V_n} \\ &= \lim_{n \rightarrow \infty} \frac{U_n}{V_n} \\ &= \lim_{n \rightarrow \infty} X_n && \text{(chain rule)} \\ &= X. \end{aligned}$$

Finally, we already know $\mathbf{Q} = U\pi$ and $\mathbf{P} = V\pi$. We may also write $U = XV + U\mathbf{1}_{[V=0]} = XV + U\mathbf{1}_{[X=\infty]}$, so

$$\mathbf{Q} = U\pi = XV\pi + \mathbf{1}_{[X=\infty]}U\pi = X\mathbf{P} + \mathbf{1}_{[X=\infty]}\mathbf{Q},$$

as claimed. □

Corollary 13.21. *In the setting of Theorem 13.19, we have the following.*

- (1) $\mathbf{Q} \ll \mathbf{P} \Leftrightarrow \mathbf{Q}(X = \infty) = 0 \Leftrightarrow \mathbf{E}_{\mathbf{P}} X = 1.$
- (2) $\mathbf{Q} \perp \mathbf{P} \Leftrightarrow \mathbf{Q}(X = \infty) = 1 \Leftrightarrow \mathbf{E}_{\mathbf{P}} X = 0.$

Proof. If $\mathbf{Q} \ll \mathbf{P}$ then by Lemma 13.20 we have $\mathbf{Q} = X\mathbf{P}$ so clearly $\mathbf{Q}(X = \infty) = 0$. We now repeatedly use (13.3). If $\mathbf{Q}(X = \infty) = 0$ then by (13.3) we have

$$\mathbf{E}_{\mathbf{P}} \{X\} = \mathbf{E}_{\mathbf{Q}} \{1\} - \mathbf{E}_{\mathbf{Q}} \{\mathbf{1}_{[X=\infty]}\} = 1.$$

If $\mathbf{E}_{\mathbf{P}} \{X\} = 1$ then again by (13.3), $\mathbf{Q}(X = \infty) = 0$ so $\mathbf{Q} = X\mathbf{P}$ and thus $\mathbf{Q} \ll \mathbf{P}$. This proves the first line of equivalences of the theorem.

Note that by Exercise 13.9, X_n is an \mathcal{F}_n -martingale for \mathbf{P} so $\mathbf{E}_{\mathbf{P}} \{X\} \leq \liminf_{n \rightarrow \infty} \mathbf{E}_{\mathbf{P}} \{X_n\} < \infty$. It follows that $\mathbf{P}(X = \infty) = 0$.

If $\mathbf{Q} \perp \mathbf{P}$ then \mathbf{Q} has no absolutely continuous part with respect to \mathbf{P} . On the other hand, $X\mathbf{P} \ll \mathbf{P}$, so by (13.3) we must have $\mathbf{Q} = \mathbf{1}_{[X=\infty]}\mathbf{Q}$; this in turn implies that $\mathbf{Q}(X = \infty) = 1$.

If $\mathbf{Q}(X = \infty) = 1$ then by (13.3), $\mathbf{E}_{\mathbf{P}} \{X\} = \int X d\mathbf{P} = X\mathbf{P} = 0$. Finally, if $X\mathbf{P} = 0$ then by (13.3) we have $\mathbf{Q}(X = \infty) = 1$. But $\mathbf{P}X = \infty = 0$, which implies $\mathbf{Q} \perp \mathbf{P}$. □

14. Branching process limits

14.1. Branching process recap. It's useful to have a fixed way to label the nodes of our trees. We do so using the Ulam-Harris tree \mathcal{U} , which has nodes labelled by $\bigcup_{n \geq 0} \mathbb{N}^n$, where $\mathbb{N}^0 := \{\emptyset\}$. The node \emptyset is the root. In general, a node at level n is labeled by a string $v = v_1 v_2 \dots v_n$; it has parent $\text{par}(v) = v_1 v_2 \dots v_{n-1}$ and children $(vi, i \geq 1) = (v_1 \dots v_n i, i \geq 1)$. We think of the children of v as being born one-at-a-time: first v_1 , then v_2 and so on. If $i < j$ we say vi is an *older sibling* of vj .

We write $\mathcal{U} = \bigcup_{n \geq 0} \mathbb{N}^n$, identifying \mathcal{U} with the set of labels of its nodes. (This is a bit sloppy, since the Ulam-Harris tree is not the only graph with these node labels, but in this context it shouldn't cause any confusion.)

A *subtree* of \mathcal{U} is a set $t \subset \mathcal{U}$ with the following properties:

- (a) $\emptyset \in t$.
- (b) If $v \in t$ then $\text{par}(v) \in t$; the ancestors of v are all in t as well.
- (c) If $v \in t$, $v = wi$ then $wj \in t$ for all $j \leq i$; the older siblings of v are all in t as well.

Given a subtree t of \mathcal{U} , for $v \in t$ we write $c(v; t) = \max\{i : vi \in t\}$; this is the *outdegree*, or *number of children* of v in t , and it may be infinite. We also write $t_n := t \cap \mathbb{N}^n$, and $t_{\leq n} = \bigcup_{m=0}^n t_m$ and the like.

A subtree $t \subset \mathcal{U}$ is *finite* if $|t| < \infty$. It is *locally finite* if $t_n := t \cap \mathbb{N}^n$ is finite for all n . Its *height* is $\text{ht}(t) := \max\{n : t_n \neq \emptyset\}$.

From now on, the word “tree” means “subtree of \mathcal{U} ”, and we write \mathcal{T} for the set of locally finite trees. We wish to consider random trees, and for this we need to turn the set of trees into a measurable space.

Definition 14.1. For a tree t and an integer $n \geq 0$, let $[t]_{\leq n} = \{\text{trees } t' : t'_{\leq n} = t_{\leq n}\}$.

It’s useful to also introduce the notation $[\]_{< n} := [\]_{\leq n-1}$. The equivalence relation $[\]_{\leq n}$ partitions the set of trees into countably many equivalence classes; we let $\mathcal{F}_n = \sigma(\{[t]_{\leq n} : t \in \mathcal{T}\})$, and let $\mathcal{F} = \sigma(\bigcup_{n \geq 0} \mathcal{F}_n)$. Note that $[\]_{\leq n+1}$ refines $[\]_{\leq n}$, which implies that $(\mathcal{F}_n, n \geq 0)$ is a filtration. Note that since $[\]_{\leq n}$ is an equivalence relation, the sets $[t]_{\leq n}$ are all atoms of \mathcal{F}_n .

Fix an *offspring distribution* μ ; this is a probability distribution with $\mu(\mathbb{Z}_+) = 1$. A random subtree T of \mathcal{U} is Galton-Watson(μ)-distributed (or B_μ -distributed for short) if for all $n \geq 1$ and all locally finite subtrees t of \mathcal{U} ,

$$\mathbf{P}\{T_{\leq n} = t_{\leq n}\} = B_\mu([t]_{\leq n}) := \prod_{v \in t_{< n}} \mu(\{c(v; t)\}) = \prod_{m=0}^{n-1} \prod_{v \in t_m} \mu(\{c(v; t)\}).$$

One “concrete” way to build such a tree T is as follows. Let $(X_v, v \in \mathcal{U})$ be independent with law μ . Then let T be the subtree of \mathcal{U} in which the root \emptyset has X_\emptyset children and more generally, inductively, if $v \in T$ then $c(v, T) := X_v$.

Let T be B_μ -distributed. In what follows we’ll write $\mu(i)$ instead of the more cumbersome $\mu(\{i\})$. Let $\alpha := \sum_{i \geq 1} i\mu(i) = |\mu|_1$ and $\sigma^2 := \sum_{i \geq 1} i(i - \alpha)\mu(i) = |\mu|_2^2 - |\mu|_1^2$ be the mean and variance of the offspring distribution, respectively. The fundamental theorem of branching processes states that $\mathbf{P}\{|T| = \infty\} > 0$ if and only if either $\alpha > 1$ or $\mu(1) = 1$.

Those who took Math 587 in Fall 2018 saw (in particular on the final exam) the following. For $t \in \mathcal{T}$ write $Z_n = Z_n(t) = |t_n|$, set $M_n(t) := Z_n(t)/\alpha^n$, and let $M(t) = \limsup_{n \rightarrow \infty} M_n(t)$. Then (M_n) is a \mathbf{P} -martingale with respect to the filtration (\mathcal{F}_n°) , where $\mathcal{F}_n^\circ := \sigma(Z_m, 0 \leq m \leq n)$, so M_n converges almost surely to M . Moreover, by Fatou’s lemma $\mathbf{E}[M] \leq 1$.

Exercise 14.1. (M_n) is also a \mathbf{P} -martingale with respect to (\mathcal{F}'_n) , where $\mathcal{F}'_n = \sigma(X_v, v \in \mathcal{U}_{< n})$. Equivalently, (M_n) is a B_μ -martingale with respect to (\mathcal{F}_n) , where $\mathcal{F}_n = \sigma(\{[t]_{< n} : t \in \mathcal{T}\})$.

Theorem 14.2 (Fundamental theorem of branching processes). *Let B be a non-negative random variable integer random variable with distribution μ , and let T be B_μ -distributed. Then $\mathbf{P}\{|T| = \infty\} > 0$ if and only if one of the following two conditions holds.*

- $\mathbf{P}\{B = 1\} = 1$
- $\mathbf{E}[B] > 1$.

As a warm up, we prove the following lemma.

Lemma 14.3. *Let B be μ -distributed. Then for all n , $\mathbf{E}[Z_n] = [\mathbf{E}B]^n$.*

Proof. This is obviously true for $n = 0$. Supposing the equality holds for a given n , we write

$$\mathbf{E}[Z_{n+1}] = \sum_{i=0}^{\infty} \mathbf{E}[Z_{n+1} \mid B_\emptyset = i] \mathbf{P}\{B_\emptyset = i\}.$$

Given that $B_\emptyset = i$, the children $1, \dots, i$ of \emptyset are each the root of an independent copy of the whole process, so

$$\mathbf{E}[Z_{n+1} \mid B_\emptyset = i] = i\mathbf{E}[Z_n].$$

We thus have

$$\mathbf{E}[Z_{n+1}] = \sum_{i=0}^{\infty} i\mathbf{E}[Z_n] \mathbf{P}\{B_\emptyset = i\} = \mathbf{E}[Z_n] \cdot \mathbf{E}[B] = [\mathbf{E}[B]]^{n+1},$$

the last step by induction. □

Corollary 14.4. *If $\mathbf{E}[B] < 1$ then $\mathbf{E}|T| < \infty$, so $\mathbf{P}\{|T| = \infty\} = 0$.*

Proof. If $\mathbf{E}[B] < 1$ then

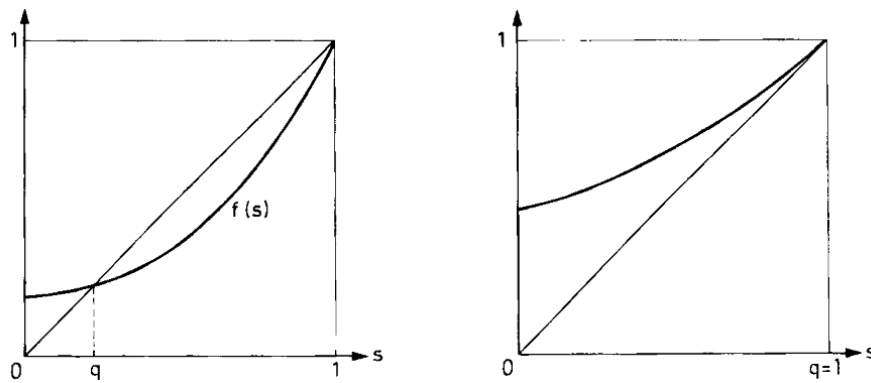
$$\mathbf{E}|T| = \sum_{n=0}^{\infty} \mathbf{E}[Z_n] = \sum_{n=0}^{\infty} (\mathbf{E}B)^n = \frac{1}{1 - \mathbf{E}[B]} < \infty.$$

It follows by Markov's inequality that $\mathbf{P}\{|T| = \infty\} = 0$. □

Now let $F(z) = \mathbf{E}[z^B] = \sum_{k=0}^{\infty} \mathbf{P}\{B = k\} z^k$.

Proposition 14.5 (Fundamental theorem of branching processes). *If $\mathbf{P}\{B = 1\} < 1$ then*

$$\mathbf{P}\{|T| < \infty\} = \min_{x \geq 0} \{F(x) = x\}.$$



Proof. Write $p = \mathbf{P}\{|T| < \infty\}$. We prove the proposition in two parts: first we show that $F(p) = p$, and second we show that p is the *smallest* non-negative solution of $F(x) = x$.

The proof of the first part is similar to that of the proof of the lemma. We begin by noting that

$$|T| < \infty \Leftrightarrow Z_n = 0 \text{ for some } n,$$

so

$$p = \mathbf{P}\left\{\bigcup_{n=0}^{\infty} Z_n = 0\right\}.$$

The events on the right are increasing (if $Z_n = 0$ then $Z_{n+1} = 0$) so it follows that

$$p = \lim_{n \rightarrow \infty} \mathbf{P}\{Z_n = 0\}.$$

Now write $F_1(x) = F(x)$ and for $n > 1$ write $F_n(x) = F(F_{n-1}(x))$, so $F_n(x)$ is the result of applying F to x n times.

We claim that for all $n \geq 1$, $\mathbf{P}\{Z_n = 0\} = F_n(0)$. When $n = 1$, we have $F_1(0) = F(0) = \mathbf{P}\{B = 0\} = \mathbf{P}\{Z_1 = 0\}$. For larger n , we apply the same inductive technique as in Lemma 1.

$$\begin{aligned} \mathbf{P}\{Z_n = 0\} &= \sum_{i=0}^{\infty} \mathbf{P}\{Z_n = 0 \mid Z_1 = i\} \mathbf{P}\{Z_1 = i\} \\ &= \sum_{i=0}^{\infty} \mathbf{P}\{Z_{n-1} = 0\}^i \mathbf{P}\{B = i\} \\ &= \sum_{i=0}^{\infty} F_{n-1}(0)^i \mathbf{P}\{B = i\} \\ &= F(F_{n-1}(0)) \\ &= F_n(0). \end{aligned}$$

We now have

$$p = \lim_{n \rightarrow \infty} F_n(0).$$

Since $F_n(0) \rightarrow p$ and F is continuous, we also have $F(F_n(0)) \rightarrow F(p)$. But $F(F_n(0)) \rightarrow p$, so we must have $p = F(p)$.

For the second part, suppose q is any other non-negative solution of $F(x) = x$. By differentiation we see that F is non-decreasing and so since $q \geq 0$ we must have $q = F(q) \geq F(0)$. Repeatedly applying F we see that we must have $q \geq F_n(0)$ for all n , and so $q \geq \lim_{n \rightarrow \infty} F_n(0) = p$. \square

Proof of Fundamental Theorem. We already saw that if $\mathbf{E}[B] < 1$ then extinction is certain, so we assume that $\mathbf{E}[B] \geq 1$. Case (a) is also obvious so we assume that $\mathbf{P}\{B = 1\} < 1$. Note that $F(0) = \mathbf{P}\{B = 0\} \geq 0$ and that $F''(x) > 0$ for all $x > 0$. Also,

$$F'(z) = \left(\sum_{n=0}^{\infty} \mathbf{P}\{B = n\} z^n \right)' = \sum_{n=1}^{\infty} n \mathbf{P}\{B = n\} z^{n-1},$$

so $F'(1) = \sum_{n=1}^{\infty} n \mathbf{P}\{B = n\} = \mathbf{E}[B]$. If $\mathbf{E}[B] > 1$ then by continuity there is $x < 1$ such that $F(x) < x$, so by the intermediate value theorem, there is $0 \leq y < x$ with $F(y) = y$, and we must have $p < 1$. On the other hand, if $\mathbf{E}[B] = 1$ then since $\mathbf{P}\{B = 1\} < 1$ there must be $k > 1$ such that $\mathbf{P}\{B = k\} > 0$. It follows that $F''(x) > 0$ for all $x > 0$, so we must have $F(x) > x$ for all $0 \leq x < 1$, and so $p = 1$. \square

Exercise 14.2. Let $(Z_n, n \geq 0)$ be the generation sizes in a Galton-Watson process with offspring distribution μ . Let B be μ -distributed and write $\alpha = \mathbf{E}B$ and $\sigma^2 = \mathbf{Var}(B)$. We suppose in this question that $\sigma^2 \in (0, \infty)$ and that $\alpha > 1$. Also, write $M_n = Z_n / (\mathbf{E}B)^n$ and let M be the a.s. martingale limit of M_n .

(a) Prove that for every $n \geq 0$,

$$\mathbf{E}\{Z_{n+1}^2 \mid \mathcal{F}_n\} = (\mathbf{E}B)^2 Z_n^2 + \sigma^2 Z_n.$$

(b) Prove that for every $n \geq 0$,

$$\mathbf{E}[Z_n^2] = \alpha^{2n} + \frac{\sigma^2(\alpha^n - \alpha^{2n})}{\alpha(1 - \alpha)}$$

(c) Prove that $M_n \rightarrow M$ in L_2 and that $\mathbf{Var}(M) = \frac{\sigma^2}{\alpha(\alpha - 1)}$.

14.2. Branching processes with immigration. These are very natural extensions of branching processes where at each generation a random number of individuals “immigrate”, joining the current population. The generation size process $(U_n)_{n \geq 0}$ of a branching process with offspring distribution μ and immigration distribution ν may be constructed as follows. Let $(X_{n,k}, n, k \geq 1)$ be IID with law μ , and independently let $(Y_n, n \geq 0)$ be IID with law ν . Then set $U_0 = Y_0$ and, for $n \geq 0$ let $U_{n+1} = Y_{n+1} + X_{n,1} + \dots + X_{n,U_n}$. Note that this construction makes perfect sense with $(Y_n, n \geq 0)$ replaced by a deterministic vector $y = (y_n, n \geq 0)$ of non-negative integers; this will be useful below.

The next theorem characterizes when immigration leads to super-exponential population growth.

Theorem 14.6 (Seneta, 1970). *Let Y have law ν . If $\mathbf{E}[\max(\log Y, 0)] < \infty$ then $\lim_{n \rightarrow \infty} U_n / \alpha^n$ exists and is almost surely finite. If $\mathbf{E}[\max(\log Y, 0)] = \infty$ then $\lim_{n \rightarrow \infty} U_n / c^n$ is almost surely infinite for all $c > 0$.*

Lemma 14.7. *Let $(R_n, n \geq 1)$ be IID and non-negative.*

(a) *If $\mathbf{E}R_1 < \infty$ then almost surely $\limsup_{n \rightarrow \infty} \frac{R_n}{n} = 0$ and $\sum_{n \geq 1} e^{R_n} c^n < \infty$ for all $c \in (0, 1)$.*

(b) *If $\mathbf{E}R_1 = \infty$ then almost surely $\limsup_{n \rightarrow \infty} \frac{R_n}{n} = \infty$ and $\sum_{n \geq 1} e^{R_n} c^n = \infty$ for all $c \in (0, 1)$.*

Proof. Suppose $\mathbf{E}R_1 < \infty$ and fix any $\epsilon > 0$. Then

$$\sum_{n > 0} \mathbf{P}\{R_n \geq \epsilon n\} = \sum_{n > 0} \mathbf{P}\{R_1 \geq \epsilon n\} \leq \frac{1}{\epsilon} \sum_{n \geq 0} \mathbf{P}\{R_1 \geq n\} = \frac{\mathbf{E}R_1}{n} < \infty,$$

so by the first Borel-Cantelli lemma, $\limsup_{n \rightarrow \infty} R_n/n \leq \epsilon$ almost surely, and since $\log(1 - a) < -a$ for $a \in (0, 1)$, letting $N_0 = \sup\{n : R_n \geq \epsilon n\}$, which is almost surely finite, for all $c \in (0, 1 - 2\epsilon)$ we have

$$\begin{aligned} \sum_{n \geq 1} e^{R_n} c^n &< \sum_{n \geq 1} e^{R_n + n \log(1-2\epsilon)} \\ &\leq \sum_{n \geq 1} e^{R_n - 2\epsilon n} \\ &\leq \sum_{n \leq N_0} e^{R_n - 2\epsilon n} + \sum_{n > N_0} e^{-\epsilon n} \\ &< \infty. \end{aligned}$$

Since $\epsilon > 0$ was arbitrary, the first result follows.

Next suppose $\mathbf{E}R_1 = \infty$ and fix any $C > 1$. Then

$$\sum_{n > 0} \mathbf{P}\{R_n \geq Cn\} = \sum_{n > 0} \mathbf{P}\{R_1 \geq Cn\} \geq \frac{1}{C} \sum_{n \geq C} \mathbf{P}\{R_1 \geq n\} \geq \frac{\mathbf{E}R_1 - C}{C} = \infty,$$

so by the second Borel-Cantelli lemma, almost surely $R_n/n \geq C$ infinitely often. It follows that almost surely $\limsup_{n \rightarrow \infty} R_n/n \geq C$, and for any $c > 1/C$,

$$\sum_{n > 0} e^{R_n} c^n \geq \sup_{n > 0} e^{R_n} c^n \geq \sup_{n \geq 0} (Cc)^n = \infty.$$

Since $C > 1$ was arbitrary, the second result follows. □

Proof of Theorem 14.6. First suppose that $\mathbf{E}[\max(\log Y_1, 0)] = \infty$. Then for all $c > 0$,

$$\limsup_n \frac{U_n}{c^n} \geq \limsup_n Y_n c^n = \infty,$$

the last inequality holding almost surely by Lemma 14.7.

Next suppose that $\mathbf{E}[\max(\log Y_1, 0)] < \infty$. Let $U_{n,k}$ be the number of generation- n descendants of generation- k immigrants. Conditionally given Y_k , $U_{n,k}$ is just distributed as the number of individuals in generation $n - k$ of a branching process started with Y_k individuals. Moreover, $U_{n,k}$ is independent of $(Y_j, j \neq k)$, so if $\mathcal{G} := \sigma(Y_n, n \geq 1)$ then

$$\mathbf{E}\{U_{n,k} \mid \mathcal{G}\} = \mathbf{E}\{U_{n,k} \mid Y_k\} = Y_k \alpha^{n-k}.$$

Since $U_n = \sum_{k=0}^n U_{n,k}$ it follows that

$$\begin{aligned} \mathbf{E}\{\alpha^{-n} U_n \mid \mathcal{G}\} &= \sum_{k \leq n} \mathbf{E}\{\alpha^{-n} U_{n,k} \mid \mathcal{G}\} \\ &= \sum_{k \leq n} \frac{Y_k}{\alpha^k}. \end{aligned} \tag{14.1}$$

Lemma 14.7 gives that $\sum_{k \leq n} \frac{Y_k}{\alpha^k} \rightarrow \sum_{n \leq 0} \frac{Y_n}{\alpha^n} < \infty$ almost surely, so by the conditional Fatou lemma, almost surely

$$\mathbf{E}\left\{\liminf_{n \rightarrow \infty} \alpha^{-n} U_n \mid \mathcal{G}\right\} \leq \liminf_{n \rightarrow \infty} \mathbf{E}\{\alpha^{-n} U_n \mid \mathcal{G}\} = \sum_{n \leq 0} \frac{Y_n}{\alpha^n} < \infty.$$

Thus, $\mathbf{P}\{\liminf_{n \rightarrow \infty} \alpha^{-n} U_n = \infty \mid \mathcal{G}\} = 0$ almost surely. But then

$$\mathbf{P}\left\{\liminf_{n \rightarrow \infty} \alpha^{-n} U_n = \infty\right\} = \mathbf{E}\left[\mathbf{P}\left\{\liminf_{n \rightarrow \infty} \alpha^{-n} U_n = \infty \mid \mathcal{G}\right\}\right] = 0,$$

so almost surely

$$\liminf_{n \rightarrow \infty} \alpha^{-n} U_n < \infty.$$

It remains to show that $\alpha^{-n}U_n$ converges almost surely. For this we use the submartingale convergence theorem, which states that a submartingale which is bounded in expectation converges almost surely. We have

$$\mathbf{E}\{U_{n+1} \mid U_1, \dots, U_n, \mathcal{G}\} = \alpha U_n + Y_{n+1},$$

so $(\alpha^{-n}U_n, n \geq 0)$ is a submartingale with respect to its natural filtration given \mathcal{G} ; the fact that it is bounded in expectation (given \mathcal{G}) follows from (14.1). \square

There is a nice construction of branching processes with immigration within the Ulam-Harris tree. A *spinal* tree is a pair (t, p) , where $t \in \mathcal{T}$ and $p = p_0, p_1, \dots$ is a finite or infinite path in t , starting from the root. We write $p_{\leq n}$ for the truncation of p at level n , so if p has at most n nodes then $p = p_{\leq n}$, and otherwise $p_{\leq n} = p_0, p_1 \dots p_n$.

Let $X = (X_v, v \in \mathcal{U})$ be IID with distribution μ . Fix a vector $y = (y_n, n > 0)$ of non-negative integers, and another vector $i = (i_n, n > 0)$ of integers with $1 \leq i_n \leq y_n + 1$ for all $n > 0$. Let $P_n = P_n(i) := i_1 i_2 \dots i_n$, so $P_{n+1} = P_n i_{n+1}$ for $n \geq 0$, and let $P = P_0, P_1, \dots$. Then define a random tree $T = T(X, y, i)$ containing $P(i)$, as follows.

- (1) Let $\emptyset \in T$ and let $p_0 = \emptyset$.
- (2) For $n \geq 0$, given $T_{\leq n}$:
 - Let $c(P_n; T) = y_{n+1} + 1$. (Note that $i_{n+1} \leq c(P_n; T)$ so $P_{n+1} \in T_{n+1}$.)
 - For $v \in T_n$ with $v \neq P_n$, let $c(v; T) = X_v$.

Exercise 14.3. *The process $(|T_{n+1}| - 1, n \geq 0)$ is distributed as a branching process with immigration with offspring distribution μ and immigration vector y .*

We next introduce a sigma-field on spinal trees, much the same as we did for the set of trees \mathcal{T} . The set of spinal trees is

$$\mathcal{T}^* = \{(t, p) : t \in \mathcal{T}, p \text{ a path in } t \text{ starting at the root}\}.$$

For each $n \geq 0$, for each pair (t, v) where $t \in \mathcal{T}$ and $v \in t_n$, we define an equivalence class

$$[(t, v)]_{\leq n} = \{(t', p') \in \mathcal{T}^* : t'_{\leq n} = t_{\leq n}, p'_{\leq n} \text{ passes through } v\}.$$

Let $\mathcal{F}_n^* = \sigma(\bigcup_{m=0}^n \{[(t, p)]_{\leq m} : (t, p) \in \mathcal{T}^*\})$, and let $\mathcal{F}^* = \sigma(\bigcup_{n \geq 0} \mathcal{F}_n^*)$. The reason the definition of \mathcal{F}_n^* has a union over $m \leq n$ is that we allow for finite paths, which may end at some level $m \leq n$. Again, $(\mathcal{F}_n^*, n \geq 0)$ is a filtration, and it is easy to see that \mathcal{F}_n^* refines \mathcal{F}_n for each n .

14.3. The Kesten-Stigum theorem. Recall that $M_n = Z_n/\alpha^n$, and that $M = \limsup_{n \rightarrow \infty} M_n$ is the a.s. martingale limit of M_n . The goal of this section is to prove the *Kesten-Stigum* theorem, which provides necessary and sufficient conditions for M_n to converge to M in L_1 .

Theorem 14.8 (Kesten-Stigum Theorem). *Fix an offspring distribution μ with $\alpha = \sum_{i \geq 1} i\mu(i) > 1$. Let T be B_μ -distributed, and let M_n and M be defined as above. Then the following are equivalent.*

- (i) $\mathbf{P}\{M = 0\} = \mathbf{P}\{|T| < \infty\}$
- (ii) $\mathbf{E}M = 1$
- (iii) $\sum_{i \geq 1} \mu(i) \cdot i \log i < \infty$.

Remarks.

- Note that if ω is such that $|T(\omega)| < \infty$ then $M_n(\omega) = 0$ for all n large, so $M(\omega) = 0$. It follows that $\mathbf{P}\{M = 0\} \geq \mathbf{P}\{|T| < \infty\}$.
- We have seen (thm/ex references) that $\mathbf{E}M_n \rightarrow \mathbf{E}M$ if and only if (M_n) is uniformly integrable (in which case $M_n \xrightarrow{L_1} M$), so a fourth equivalent condition which can be added to the Kesten-Stigum theorem is that (M_n) is UI.

To prove the Kesten-Stigum theorem we use a beautiful method called a “spinal change of measure”. Recall that the *size-biasing* $\hat{\mu}$ of μ is the probability distribution with $\hat{\mu}(i) = i\mu(i)/\alpha$. Note that if B as law $\hat{\mu}$ then $\mathbf{P}\{B \geq 1\} = 1$.

Let ν be the probability measure on \mathbb{Z}^+ defined by setting $\nu(i) = \hat{\mu}(i - 1)$ for all i . Let $X = (X_v, v \in \mathcal{U})$ are independent with law μ , let $Y = (Y_n, n > 0)$ be independent with law ν , and let $U = (U_n, n > 0)$ be independent Uniform $[0, 1]$ random variables, with X, Y and U mutually independent. For $n > 0$ let $I_n = \lceil (Y_n + 1)U_n \rceil$, so that I_n is a uniformly random element of $\{1, \dots, Y_n + 1\}$. Then write BPI_μ^* for the law of the pair $(T, P) = (T(X, Y, I), P(I))$, and let BPI_μ be the law of the tree $T = T(X, Y, I)$ obtained from (T, P) by “ignoring the spine”.

Proposition 14.9. *For any offspring distribution μ with $\mu(0) < 1$, and any spinal tree (t, p) , for all $n \geq 0$,*

$$\text{BPI}_\mu^*(t_{\leq n}, p_{\leq n}) = \frac{1}{\alpha^n} \text{B}_\mu(t_{\leq n}).$$

Proof. Let (T, P) be constructed as above, so that

$$\text{BPI}_\mu^*(t_{\leq n}, p_{\leq n}) = \mathbf{P}\{(T_{\leq n}, P_{\leq n}) = (t_{\leq n}, p_{\leq n})\}.$$

Then write

$$\mathbf{P}\{(T_{\leq n}, P_{\leq n}) = (t_{\leq n}, p_{\leq n})\} = \prod_{i=0}^{n-1} \mathbf{P}\{T_{i+1} = t_{i+1}, P_{i+1} = p_{i+1} \mid (T_{\leq i}, P_{\leq i}) = (t_{\leq i}, p_{\leq i})\}.$$

Now, given that $(t_{\leq i}, p_{\leq i})$, in order to have $T_{i+1} = t_{i+1}$ and $P_{i+1} = p_{i+1}$, the following must occur: p_i must have the right number of children; the correct extension of $p_{\leq i}$ must be chosen; and all the other nodes in t_i must also have the right number of children. The probability of all these occurring is

$$\begin{aligned} & \mathbf{P}\{T_{i+1} = t_{i+1}, P_{i+1} = p_{i+1} \mid (T_{\leq i}, P_{\leq i}) = (t_{\leq i}, p_{\leq i})\} \\ &= \hat{\mu}(c(p_i; t)) \cdot \frac{1}{c(p_i; t)} \cdot \prod_{v \in t_i, v \neq p_i} \mu(c(v; t)) \\ &= \frac{c(p_i; t)\mu(c(p_i; t))}{\alpha} \cdot \frac{1}{c(p_i; t)} \cdot \prod_{v \in t_i, v \neq p_i} \mu(c(v; t)) \\ &= \frac{1}{\alpha} \prod_{v \in t_i} \mu(c(v; t)), \end{aligned}$$

which combined with the two previous equations gives

$$\text{BPI}_\mu^*(t_{\leq n}, p_{\leq n}) = \prod_{i=0}^{n-1} \left(\frac{1}{\alpha} \prod_{v \in t_i} \mu(c(v; t)) \right) = \frac{1}{\alpha^n} \text{B}_\mu(t_{\leq n}). \quad \square$$

Corollary 14.10. *For all $n \geq 0$,*

$$\frac{d\text{BPI}_\mu|_{\mathcal{F}_n}}{d\text{B}_\mu|_{\mathcal{F}_n}} = M_n.$$

Proof. For any subtree t of \mathcal{U} , by definition,

$$\text{BPI}_\mu(t_{\leq n}) = \sum_p \text{BPI}_\mu^*(t_{\leq n}, p),$$

where the sum is over paths p from the root to generation n in $t_{\leq n}$. But the number of such paths is just $|t_n|$. Using Proposition 14.9 and the fact that $M_n(t) = |t_n|/\alpha^n$, we thus have

$$\text{BPI}_\mu(t_{\leq n}) = \frac{|t_n|}{\alpha^n} \text{B}_\mu(t_{\leq n}) = M_n(t) \cdot \text{B}_\mu(t_{\leq n}).$$

and the result follows. □

Before proving the Kesten-Stigum theorem, we need one further lemma.

Lemma 14.11. *Either $\mathbf{P}\{M = 0\} = \mathbf{P}\{|T| < \infty\}$ or $\mathbf{P}\{M = 0\} = 1$.*

Proof. If $i \in T_1$ then the subtree of T rooted at i is itself a B_μ -branching process. Writing

$$M_n^{(i)} = \frac{1}{\alpha^{n-1}} \#\{v \in T_n : 1 \text{ is an ancestor of } v\},$$

then $M_n^{(i)}$ is a martingale; writing $M^{(i)}$ for its almost sure limit, we may decompose M as

$$M = \frac{1}{\alpha} \left(M_n^{(1)} + \dots + M_n^{(X_\emptyset)} \right).$$

Conditionally given that $X = k$, the limits $M_n^{(1)}, \dots, M_n^{(k)}$ are independent copies of M , and $M = 0$ if and only if each of $M_n^{(1)}, \dots, M_n^{(k)}$ equals zero. Thus

$$p := \mathbf{P}\{M = 0\} = \sum_{k \geq 0} \mathbf{P}\{X = k\} \mathbf{P}\{M = 0\}^k = \mathbf{E}[p^X].$$

The only roots the equation $s = \mathbf{E}[s^X]$ are $\mathbf{P}\{|T| < \infty\}$ and 1, so the lemma follows. \square

Proof of Theorem 14.8. Let X have law μ , let Y have law ν where $\nu(i) = \hat{\mu}(i+1)$, and let $L = \log(Y+1)$. It is easy to verify that $\mathbf{E}L < \infty$ if and only if $\mathbf{E}[\log^+ Y] < \infty$, and

$$\mathbf{E}[X \log^+ X] = \sum_{i > 0} (i \log i) \mu(i) = \sum_{i > 0} \log(i) \hat{\mu}(i) = \mathbf{E}L,$$

so by Theorem 14.6 $\text{BPI}_\mu(M < \infty) = 1$ if and only if $\mathbf{E}L < \infty$, i.e. if and only if $\mathbf{E}[X \log^+ X] < \infty$.

We now use that

$$M = \limsup_n M_n = \limsup_n \frac{d\hat{B}_\mu|_{\mathcal{F}_n}}{dB_\mu|_{\mathcal{F}_n}}.$$

by Corollary 14.10. Since

$$\mathbf{E}M = \int M(t) B_\mu(dt) = B_\mu(M),$$

It follows by Theorem 13.19 that $\mathbf{E}M = 1$ if and only if $\text{BPI}_\mu(M < \infty) = 1$, which occurs if and only if $\mathbf{E}[X \log^+ X] < \infty$.

Now, if $\mathbf{E}M = 1$ we must have $\mathbf{P}\{M = 0\} < 1$, in which case $\mathbf{P}\{M = 0\} = \mathbf{P}\{|T| < \infty\}$ by Lemma 14.11.

Finally, if $\mathbf{E}[X \log^+ X] = \infty$ then $\text{BPI}_\mu(M = \infty) = 1$, and by Theorem 13.19 this implies that $B_\mu(M = 0) = 1$, or in other words, that $\mathbf{P}\{M = 0\} = 1$; we then have $\mathbf{E}M = 0 < 1$. \square

15. Transforms 1: Moment-generating functions

15.1. **Introduction.** Here is a high-level description of many arguments involving transforms in probability. A transform τ takes as input a random variable X , and outputs some function τ_X , the transform of X . All the transforms τ we study will *factor* through the set of probability distributions, in the sense that if $X \stackrel{d}{=} Y$ then $\tau_X = \tau_Y$. The first property that makes a transform useful is *uniqueness*.

- (1) **[Uniqueness.]** There exists a suitably rich collection \mathcal{C} of probability measures such that τ is injective relative to \mathcal{C} : if $\mathcal{L}_X \in \mathcal{C}$ and $\mathcal{L}_Y \in \mathcal{C}$ and $\mathcal{L}_X \neq \mathcal{L}_Y$ then $\tau_X \neq \tau_Y$.

The cumulative distribution function is an example of a transform, and has the uniqueness property relative to the collection of all probability measures on \mathbb{R} .

The second property that makes a transform useful is *stability*.

- (1) **[Stability.]** There exists a suitably rich collection \mathcal{C} of probability measures such that if $(X_n, 1 \leq n \leq \infty)$ are random variables and $\mathcal{L}_{X_n} \in \mathcal{C}$ for all $n \leq \infty$, then $X_n \xrightarrow{d} X_\infty$ if and only if $\tau_{X_n} \rightarrow \tau_{X_\infty}$, for an appropriately defined notion of convergence of transforms.

Again using the cumulative distribution function as an example, we have that $X_n \xrightarrow{d} X_\infty$ if and only if $F_n(x) \rightarrow F_\infty(x)$ for all x which are points of continuity of F_∞ . Often it takes a little thought to find the right notion of convergence for the transform. One way to organize this thought is to start with a naive guess of what notion of convergence to use, then look for counterexamples to stability, and use the counterexamples to improve the guess of what notion of convergence is reasonable.

15.2. The moment generating function. Given a random variable X with distribution $\mathcal{L}_X = \mu$, the *moment generating function* of X is

$$G_X(s) := \mathbf{E} [e^{sX}] = \int_{\mathbb{R}} e^{sx} \mu(dx) \in (0, \infty].$$

Theorem 15.1. *If G_X is finite in a neighbourhood of s then for all $k \geq 0$, $X^k e^{sX}$ is integrable and*

$$G_X^{(k)}(s) = \mathbf{E} [X^k e^{sX}].$$

Moreover, if G_X is finite in a neighbourhood of 0 then for $|s|$ sufficiently small,

$$G_X(s) = \sum_{n=0}^{\infty} \frac{s^n}{n!} \mathbf{E} [X^n].$$

We begin by proving the theorem when G_X is finite in neighbourhood of 0.

Proposition 15.2. *If $G_X(s) < \infty$ and $G_X(-s) < \infty$ then X^n is integrable for all n and $G(s) = \sum_{n=0}^{\infty} \frac{s^n}{n!} \mathbf{E} [X^n]$.*

Lemma 15.3. *If random variables $(X_n, n \geq 0)$ are such that $\sum_{n \geq 0} \mathbf{E} |X_n| < \infty$ then $\sum_{n \geq 0} X_n$ converges absolutely almost everywhere, its a.e. limit is integrable, and*

$$\mathbf{E} \sum_{n \geq 0} X_n = \sum_{n \geq 0} \mathbf{E} X_n.$$

Proof. Write $Y = \sum_{n \geq 0} |X_n|$. Then by the monotone convergence theorem,

$$\mathbf{E} [Y] = \sum_{n \geq 0} \mathbf{E} |X_n| < \infty,$$

so Y is finite almost everywhere. It follows that $\sum_{n \geq 0} |X_n|$ converges almost everywhere, and thus $\sum_{n \geq 0} X_n$ also converges almost everywhere. Moreover, $|\sum_{n=0}^m X_n| \leq Y$ for all m , so $|\sum_{n \geq 0} X_n| \leq Y$ and is therefore integrable, and by the dominated convergence theorem

$$\mathbf{E} \sum_{n \geq 0} X_n = \sum_{n \geq 0} \mathbf{E} X_n. \quad \square$$

Proof of Proposition 15.2. We take $X_n = (sX)^n/n!$, and obtain

$$\sum_{n \geq 0} \mathbf{E} |X_n| = \sum_{n \geq 0} \mathbf{E} \frac{|sX|^n}{n!} = \mathbf{E} \sum_{n \geq 0} \frac{|sX|^n}{n!},$$

by the monotone convergence theorem. The last sum is just the Taylor expansion of $e^{|sX|}$, so

$$\sum_{n \geq 0} \mathbf{E} |X_n| = \mathbf{E} e^{|sX|} \leq \mathbf{E} [e^{sX} + e^{-sX}] = G(s) + G(-s) < \infty,$$

the last bound by assumption. It follows by Lemma 15.3 that

$$e^{sX} = \sum_{n \geq 0} \frac{(sX)^n}{n!} < \infty$$

Defined earlier, first few lines here probably redundant. However, there is a sign change for earlier that needs to be propagated!

almost everywhere and that

$$G(s) = \sum_{n \geq 0} \mathbf{E} \frac{(sX)^n}{n!} = \sum_{n \geq 0} \frac{s^n}{n!} \mathbf{E} [X^n].$$

□

If $f : \mathbb{R} \rightarrow \mathbb{R}$ which has a Taylor series expansion $f(s) = \sum_{n \geq 0} c_n s^n$ around 0, we write $[s^n]f(s) := c_n$. This allows us to refer to coefficients of such an expansion without naming them explicitly.

Proposition 15.4. *If G_X is finite in a neighbourhood of 0 then $G_X^{(k)}(0) = \mathbf{E} [X^k]$ for all $k \geq 0$.*

Proof. Recall²¹ that if $f : \mathbb{R} \rightarrow \mathbb{R}$ has a Taylor series expansion on $(-r, r)$, say $f(s) = \sum_{n \geq 0} c_n s^n$, then f is differentiable on $(-r, r)$, and its derivative has a Taylor series expansion around 0 on $(-r, r)$ as

$$f'(s) = \sum_{n \geq 1} n c_n s^{n-1}.$$

In the shorthand introduced above, this is summarized by the statement that

$$[s^{n-1}]f'(s) = n \cdot [s^n]f(s)$$

on $(-r, r)$. In particular, $f'(0) = [s^1]f(s)$. By induction, for all $n \geq k \geq 1$ we have

$$[s^{n-k}]f^{(k)}(s) = n \cdot (n-1) \cdots (n-k+1) [s^n]f(s),$$

and taking $n = k$ gives $f^{(k)}(0) = k! [s^k]f(s)$. Applying this to $G_X(s) = \sum_{n \geq 0} (\mathbf{E} [X^n] / n!) s^n$, we obtain that $G_X^{(k)}(0) = k! \cdot (\mathbf{E} [X^k] / k!) = \mathbf{E} [X^k]$, as claimed. □

To extend from the case $s = 0$ to the general case, we use a technique which is also a fundamental tool in Monte Carlo estimation, in particular for importance sampling: *exponential tilting*.

Definition 15.5 (Exponential Tilting). *Fix $s \in \mathbb{R}$ with $G_X(s)$ finite, and let μ_s be the probability measure on \mathbb{R} defined by $\mu_s(A) = \mathbf{E} [e^{sX} \mathbf{1}_{[X \in A]}] / G_X(s)$.²² Then the density (Radon-Nikodym derivative) of μ_s with respect to μ is $e^{sX} / G_X(s)$, and we call μ_s the tilting of μ by e^{sX} .*

If \hat{X} has distribution μ_s then we may also refer to \hat{X} as a tilting of X by e^{sX} .

Proof of Theorem 15.1. It only remains to prove that the statements in the first sentence of the theorem are true. Suppose that G_X is finite near s , and let \hat{X} be the tilting of X by e^{sX} . Then

$$G_{\hat{X}}(t) = \int_{\mathbb{R}} e^{tx} \mu_s(dx) = \int_{\mathbb{R}} e^{tx} \frac{e^{sx}}{G_X(s)} \mu(dx) = \frac{G_X(s+t)}{G_X(s)},$$

so $G_{\hat{X}}$ is finite in a neighbourhood of 0. Proposition 15.2 yields that \hat{X}^k is integrable; since \hat{X} has density $e^{sX} / G_X(s)$ with respect to \mathcal{L}_X , it follows that $X^k e^{sX}$ is integrable and, using Proposition 15.4, that

$$G_{\hat{X}}^{(k)}(0) = \mathbf{E} [\hat{X}^k] = \int_{\mathbb{R}} x^k \mu_s(dx) = \int_{\mathbb{R}} x^k \frac{e^{sx}}{G_X(s)} \mu(dx) = \frac{1}{G_X(s)} \mathbf{E} [X^k e^{sX}].$$

Also, $G_{\hat{X}}^{(k)}(0) = G_X^{(k)}(s) / G_X(s)$, which together with the above gives

$$G_X^{(k)}(s) = \mathbf{E} [X^k e^{sX}].$$

□

²¹Or prove it for yourself.

²²Exercise: prove it is a probability measure.

Example 15.6. Let N be standard normal. Then

$$G_N(s) = \mathbf{E}e^{sN} = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{sx} e^{-x^2/2} dx = e^{s^2/2} \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-(x-s)^2/2} dx = e^{s^2/2}.$$

The Taylor expansion of $e^{s^2/2}$ around zero is

$$e^{s^2/2} = \sum_{n \geq 0} \frac{1}{n!} (s^2/2)^n = \sum_{n \geq 0} \frac{s^{2n}}{2^n n!},$$

so it follows from Theorem 15.1 that the odd moments of N vanish, and that for $k \geq 0$,

$$\mathbf{E} [N^{2k}] = G_N^{(2k)}(0) = \frac{(2k)!}{2^k k!}.$$

Example 15.7. Let P be Poisson(λ), so $\mathbf{P}\{P = k\} = \lambda^k e^{-\lambda}/k!$. Then

$$G_P(s) = \sum_{k \geq 0} e^{ks} \frac{\lambda^k e^{-\lambda}}{k!} = e^{-\lambda} \sum_{k \geq 0} \frac{(\lambda e^s)^k}{k!} = e^{\lambda(e^s - 1)}.$$

We wish to calculate

$$\mathbf{E} [P^k] = G_P^{(k)}(0) = k! [s^k] G_P(s) = k! [s^k] e^{\lambda(e^s - 1)}.$$

Since

$$e^{\lambda(e^s - 1)} = \sum_{n \geq 0} \frac{1}{n!} (\lambda(e^s - 1))^n,$$

it follows that

$$\mathbf{E} [P^k] = k! \sum_{n \geq 0} [s^k] (\lambda(e^s - 1))^n = \sum_{n \geq 0} \lambda^n [s^k] (e^s - 1)^n.$$

To analyze this equation, we again use the Taylor expansion of $e^s - 1 = \sum_{m \geq 1} s^m/m!$, to obtain

$$\mathbf{E} [P^k] = k! \sum_{n \geq 0} \frac{\lambda^n}{n!} [s^k] \left(\sum_{m \geq 1} \frac{s^m}{m!} \right)^n.$$

Expanding out the n 'th power and gathering like powers of s , the smallest order term will be s^n , so if $n > k$ then the expression is zero. If $n \leq k$ then

$$[s^k] \left(\sum_{m \geq 1} \frac{s^m}{m!} \right)^n = [s^k] \sum_{\substack{m_1, \dots, m_n \geq 1 \\ m_1 + \dots + m_n = k}} \prod_{i=1}^n \frac{s^{m_i}}{m_i!} = \sum_{\substack{m_1, \dots, m_n \geq 1 \\ m_1 + \dots + m_n = k}} \frac{1}{m_i!} = \frac{1}{k!} \sum_{\substack{m_1, \dots, m_n \geq 1 \\ m_1 + \dots + m_n = k}} \binom{k}{m_1, \dots, m_n}.$$

We now use this in the previous formula (recall that we only need to consider $n \leq k$). The $k!$ terms cancel, and we obtain

$$\mathbf{E} [P^k] = \sum_{n=1}^k \lambda^n \sum_{\substack{m_1, \dots, m_n \geq 1 \\ m_1 + \dots + m_n = k}} \frac{1}{n!} \binom{k}{m_1, \dots, m_n}.$$

Now, $\binom{k}{m_1, \dots, m_n}$ is the number of ways to partition the set $\{1, \dots, k\}$ into an ordered sequence P_1, \dots, P_n of non-empty parts such that $|P_i| = m_i$ for $i \leq i \leq n$; the term $1/n!$ may be interpreted as saying that the order of the parts is unimportant. So the inner sum is simply the total number of ways of partitioning the set $\{1, 2, \dots, k\}$ into n nonempty parts. This is a Stirling number of the second kind, denoted $\left\{ \begin{smallmatrix} k \\ n \end{smallmatrix} \right\}$. We thus obtain

$$\mathbf{E} [P^k] = \sum_{n=1}^k \lambda^n \left\{ \begin{smallmatrix} k \\ n \end{smallmatrix} \right\}.$$

The transform which sends a variable to the sequence of its integer moments is quite natural, and the question of when it satisfies uniqueness and stability deserving of study. We will return to this later in the course and in the exercises.

15.3. The moment generating function and uniqueness. In this section we show that the moment generating function satisfies uniqueness provided we restrict our attention to non-negative random variables.

Theorem 15.8. *If $\mathbf{P}\{X \geq 0\} = 1$ then \mathcal{L}_X can be recovered from G_X .*

Proof. Since $X \geq 0$ almost surely, the moment generating function G_X is finite on $(-\infty, 0]$. For any $s > 0$, by Theorem 15.1 we have

$$G_X^{(k)}(-s) = \mathbf{E} \left[X^k e^{-sX} \right].$$

For $t > 0$ we then have

$$\sum_{k=0}^{\lfloor st \rfloor} \frac{s^k}{k!} G_X^{(k)}(-s) = \sum_{k=0}^{\lfloor st \rfloor} \mathbf{E} \left[\frac{(sX)^k}{k!} e^{-sX} \right] = \mathbf{E} \left[\sum_{k=0}^{\lfloor st \rfloor} \frac{(sX)^k}{k!} e^{-sX} \right].$$

Now, for any $m \in \mathbb{N}$ and $\lambda \geq 0$,

$$\sum_{k=0}^m \frac{\lambda^k}{k!} e^{-\lambda} = \mathbf{P} \{ \text{Poisson}(\lambda) \leq m \},$$

so we may rewrite the above as

$$\sum_{k=0}^{\lfloor st \rfloor} \frac{s^k}{k!} G_X^{(k)}(-s) = \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sX) \leq \lfloor st \rfloor \}].$$

The $\text{Poisson}(\lambda)$ distribution has²³ mean and variance λ , so for all $x < t$, by Chebyshev's inequality

$$\liminf_{s \rightarrow \infty} \mathbf{P} \{ \text{Poisson}(sx) \leq \lfloor st \rfloor \} \geq \liminf_{s \rightarrow \infty} \left(1 - \frac{(\lfloor st \rfloor - sx)^2}{sy} \right) = 1 \quad (15.1)$$

and likewise, for all $x > t$,

$$\limsup_{s \rightarrow \infty} \mathbf{P} \{ \text{Poisson}(sx) \leq \lfloor st \rfloor \} \leq \limsup_{s \rightarrow \infty} \left(\frac{(sx - \lfloor st \rfloor)^2}{sy} \right) = 0 \quad (15.2)$$

For any $x < t$, we then have

$$\begin{aligned} \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sX) \leq \lfloor st \rfloor \}] &\geq \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sX) \leq \lfloor st \rfloor \} \mathbf{1}_{[X \leq x]}] \\ &\geq \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sx) \leq \lfloor st \rfloor \} \mathbf{1}_{[X \leq x]}] \\ &= \mathbf{P} \{ \text{Poisson}(sx) \leq \lfloor st \rfloor \} \mathbf{P} \{ X \leq x \}. \end{aligned}$$

Since this holds for all s and for all $x < t$, using (15.1) we obtain

$$\liminf_{s \rightarrow \infty} \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sX) \leq \lfloor st \rfloor \}] \geq \sup_{x < t} \mathbf{P} \{ X \leq x \} = \mathbf{P} \{ X < t \}.$$

Similarly, for any $x > t$,

$$\begin{aligned} \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sX) \leq \lfloor st \rfloor \}] &\leq \mathbf{P} \{ X < x \} + \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sX) \leq \lfloor st \rfloor \} \mathbf{1}_{[X \geq x]}] \\ &\leq \mathbf{P} \{ X < x \} + \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sx) \leq \lfloor st \rfloor \} \mathbf{1}_{[X \geq x]}] \\ &\leq \mathbf{P} \{ X < x \} + \mathbf{P} \{ \text{Poisson}(sx) \leq \lfloor st \rfloor \} \mathbf{P} \{ X \geq x \}. \end{aligned}$$

Since this holds for all s and for all $x > t$, using (15.2) we obtain

$$\limsup_{s \rightarrow \infty} \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sX) \leq \lfloor st \rfloor \}] \leq \inf_{x > t} \mathbf{P} \{ X < x \} = \mathbf{P} \{ X \leq t \}.$$

²³Prove it.

It follows that

$$\lim_{s \rightarrow \infty} \sum_{k=0}^{\lfloor st \rfloor} \frac{s^k}{k!} G_X^{(k)}(-s) = \lim_{s \rightarrow \infty} \mathbf{E} [\mathbf{P} \{ \text{Poisson}(sX) \leq \lfloor st \rfloor \}] = \mathbf{P} \{ X \leq t \} = \mathcal{L}_X((-\infty, t]).$$

for all t such that $\mathbf{P} \{ X = t \} = \mathcal{L}_X(\{t\}) = 0$. Finally, for any countable set $Q \subset \mathbb{R}$, the set $\{(-\infty, t], t \in \mathbb{R} \setminus Q\}$ is a π -system generating the Borel σ -field $\mathcal{B}(\mathbb{R})$. Since \mathcal{L}_X has at most countably many atoms, the result follows. \square

Note that $G_X^{(k)}(-s) = \mathbf{E} [X^k e^{-sX}]$, so the method of proof is to recover the law of X from considering the integer moments of large negative exponential tilts of X . However, the proof requires knowledge of all moments for all large negative tilts - a single value of s is not enough. (The *Stieltjes moment problem* asks for necessary and sufficient conditions for the existence and uniqueness of a non-negative random variable with given integer moments; see [2] for recent work on this question.)

16. Transforms 2: Characteristic functions

In a probabilistic context, Fourier transforms are known as characteristic functions; they capture a function by decomposing it according to its rotational symmetries. To get a feeling for what this means, consider the following, very simple way to decompose a function into symmetric pieces. Fix a function $f : \mathbb{R} \rightarrow \mathbb{R}$, and define $f_{\text{even}} : \mathbb{R} \rightarrow \mathbb{R}$ and $f_{\text{odd}} : \mathbb{R} \rightarrow \mathbb{R}$ as

$$f_{\text{even}}(x) = \frac{f(x) + f(-x)}{2}, \quad f_{\text{odd}}(x) = \frac{f(x) - f(-x)}{2}.$$

We clearly have $f \equiv f_{\text{even}} + f_{\text{odd}}$, and the functions f_{even} and f_{odd} each satisfy a natural symmetry: f_{even} is even, which means that $f_{\text{even}}(-x) = f_{\text{even}}(x)$, and f_{odd} is odd, which means that $f_{\text{odd}}(-x) = -f_{\text{odd}}(x)$.

The above decomposition is easy to understand but is not particularly useful. The perspective of Fourier analysis generalizes the even-odd decomposition by viewing $-1 = e^{i\pi}$ as a rotation by π in the complex plane. A first natural generalization is to decompose using the family of rotations $(e^{2\pi i/n}, 0 \leq i < n)$; above we had $n = 2$. For a function $f : \mathbb{C} \rightarrow \mathbb{C}$, say f is (n, j) -symmetric if $f(e^{2\pi i/n} z) = (e^{2\pi i/n})^j f(z)$.

Proposition 16.1. Fix any function $f : \mathbb{C} \rightarrow \mathbb{C}$. Fix $n \geq 1$, write $\omega = e^{2\pi i/n}$, and for $j \in \mathbb{Z}$ let

$$f_{n,j}(z) = \frac{1}{2n} \sum_{k=-n}^{n-1} \omega^{-jk} f(\omega^k z).$$

Then each function $f_{n,j}$ is (n, j) -symmetric; if f is (n, j) -symmetric then $f = f_{n,j}$; and

$$f \equiv \sum_{j=-n}^{n-1} f_{n,j}. \tag{16.1}$$

Proof. First,

$$\begin{aligned} f_{n,j}(\omega z) &= \frac{1}{2n} \sum_{k=-n}^{n-1} \omega^{-jk} f(\omega^k \cdot \omega z) \\ &= \omega^j \cdot \frac{1}{2n} \sum_{k=-n}^{n-1} \omega^{-j(k+1)} f(\omega^{k+1} z) \\ &= \omega^j \cdot f_{n,j}(z), \end{aligned}$$

so $f_{n,j}$ is (n, j) -symmetric. Next, if f is (n, j) -symmetric then $f(\omega^k z) = \omega^{jk} f(z)$, so

$$f_{n,j}(z) = \frac{1}{2n} \sum_{k=-n}^{n-1} \omega^{-jk} f(\omega^k z) = \frac{1}{2n} \sum_{k=-n}^{n-1} f(z) = f(z).$$

Finally, for all z ,

$$\begin{aligned} \sum_{j=0}^{n-1} f_{n,j}(z) &= \frac{1}{2n} \sum_{j=-n}^{n-1} \sum_{k=-n}^{n-1} \omega^{-jk} f(\omega^k z) \\ &= \frac{1}{2n} \sum_{k=-n}^{n-1} f(\omega^k \cdot z) \cdot \sum_{j=-n}^{n-1} \omega^{-jk} \end{aligned}$$

It is straightforward that²⁴ $\sum_{j=-n}^{n-1} \omega^{-jk} = 2n \mathbf{1}_{[k=0]}$, so the previous equality gives

$$\sum_{j=-n}^{n-1} f_j(z) = \frac{1}{2n} f(\omega^k \cdot z) \cdot 2n \mathbf{1}_{[k=0]} = f(z). \quad \square$$

Let's restrict attention to points on the unit circle $S^1 = \{z \in \mathbb{C} : |z| = 1\}$. Equation (16.1) gives that $f(1) = \sum_{j=-n}^{n-1} f_{n,j}(1)$. Since $f_{n,j}$ is (n, j) -harmonic, for $k \in \mathbb{Z}$ we then have

$$f(e^{2\pi i k/n}) = \sum_{j=-n}^{n-1} f_{n,j}(e^{2\pi i k/n}) = \sum_{j=-n}^{n-1} e^{2\pi i (k/n) \cdot j} f_{n,j}(1). \quad (16.2)$$

It also jumps out that the definition of $f_{n,j}(1)$ looks like a Riemann approximation in n :

$$f_{n,j}(1) = \frac{1}{2n} \sum_{k=-n}^{n-1} \omega^{-jk} f(\omega^k) \approx \frac{1}{4\pi} \int_{-2\pi}^{2\pi} e^{-j \cdot ix} f(e^{ix}) dx = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{ix \cdot j} f(e^{ix}) dx =: \hat{f}(j); \quad (16.3)$$

the last definition makes sense provided f is Riemann integrable. In this case, provided f is sufficiently well-behaved, one may take k and n to infinity jointly in such a way that $k/n \rightarrow x$ in (16.2), and obtain the following result.

Theorem 16.2. *Suppose that $f : S^1 \rightarrow \mathbb{C}$ is continuous and that $\sum_{j \in \mathbb{Z}} |\hat{f}(j)| < \infty$. Then*

$$\lim_{N \rightarrow \infty} \sum_{j=-N}^N e^{2\pi i x \cdot j} \hat{f}(j) = f(e^{2\pi i x}),$$

uniformly over $x \in [0, 1)$.

For a proof, see [4], Chapter 2. It is natural to try to further continuize, taking $j \in \mathbb{R}$ instead of $j \in \mathbb{Z}$ in (16.3); one then guesses at the Fourier inversion formula:

$$f(e^{2\pi i x}) = \int_{\mathbb{R}} \hat{f}(t) e^{2\pi i x \cdot t} dt. \quad (16.4)$$

We do not prove the Fourier inversion formula in this form, as the natural setting of probability is not functions on S^1 but probability measures. However, having seen the derivation of a Fourier inversion formula in a simpler setting will help with intuition in what follows. The key points to remember from the above development are as follows. First, f is recovered from a collection of its symmetrizations, each of which captures its symmetries with respect to a different oscillatory frequency. Second, because each symmetrization is symmetric, its behaviour can be recovered from its value at a single point.

²⁴Prove it

16.1. Characteristic functions: basic properties and examples. Let X be a real random variable and write $\mu = \mathcal{L}_X$ for the distribution of X . The characteristic function of X is the function $\varphi_X : \mathbb{R} \rightarrow \mathbb{C}$ given by

$$\varphi_X(t) = \mathbf{E} [e^{itX}] = \mathbf{E} [\cos(tX)] + i\mathbf{E} [\sin(tX)] .$$

Note that if X has density $f : \mathbb{R} \rightarrow [0, \infty)$ with respect to Lebesgue measure then by the change of variables formula,

$$\varphi_X(t) = \int_{\mathbb{R}} e^{itx} f(x) dx ;$$

this corresponds to (16.3) but with the range of integration \mathbb{R} instead of S^1 since f now has domain \mathbb{R} .

Proposition 16.3. *For any real random variable X , the characteristic function $\varphi = \varphi_X$ satisfies the following properties*

- (1) $\varphi(0) = 1$
- (2) $\varphi(-t) = \overline{\varphi(t)}$
- (3) $\|\varphi\|_{\infty} \leq 1$
- (4) For all $t, h \in \mathbb{R}$, $|\varphi(t+h) - \varphi(t)| \leq \mathbf{E}|e^{ihX} - 1|$, so $\varphi(t)$ is uniformly continuous on \mathbb{R} .

Proof. The first property is obvious. The second is an easy calculation using the fact that \cos is even and \sin is odd:

$$\begin{aligned} \varphi(-t) &= \mathbf{E} [e^{i(-t)X}] = \mathbf{E} [\cos(-tX)] + i\mathbf{E} [\sin(-tX)] \\ &= \mathbf{E} [\cos(tX)] - i\mathbf{E} [\sin(tX)] = \overline{\mathbf{E} [\cos(tX)] + i\mathbf{E} [\sin(tX)]} \\ &= \overline{\varphi(t)} . \end{aligned}$$

Next, for any $x \in \mathbb{R}$ we have $|e^{itx}| = 1$, so for all $t \in \mathbb{R}$,

$$|\mathbf{E} [e^{itX}]| \leq \mathbf{E} [|e^{itX}|] = 1 ,$$

establishing the third property. For the fourth assertion, we compute

$$\begin{aligned} |\varphi(t+h) - \varphi(t)| &= |\mathbf{E} [e^{itX} e^{ihX}] - \mathbf{E} [e^{itX}]| \\ &= |\mathbf{E} [e^{itX} (e^{ihX} - 1)]| \\ &\leq \mathbf{E} [|e^{itX} (e^{ihX} - 1)|] \\ &= \mathbf{E} [|e^{ihX} - 1|] . \end{aligned}$$

□

Fundamental to the utility of characteristic functions in probability is their connection to independence. The basic fact is recorded in the following proposition.

Proposition 16.4. *Let X, Y be random variables defined on a common space. If X and Y are independent then for all $a \in \mathbb{R}$ and all $t \in \mathbb{R}$,*

$$\varphi_{aX+Y}(t) = \varphi_X(at)\varphi_Y(t) .$$

Proof. We use that for bounded, complex-valued random variables U, V , if U and V are independent then $\mathbf{E}[UV] = \mathbf{E}[U]\mathbf{E}[V]$; this can be checked by considering the real and complex parts separately. We then have

$$\varphi_{aX+Y}(t) = \mathbf{E} [e^{it(aX+Y)}] = \mathbf{E} [e^{i(at)X} e^{itY}] = \mathbf{E} [e^{i(at)X}] \mathbf{E} [e^{itY}] = \varphi_X(at)\varphi_Y(t) . \quad \square$$

A special case of the above proposition is that for any $a, b \in \mathbb{R}$ and $t \in \mathbb{R}$, $\varphi_{aX+b}(t) = e^{itb}\varphi_X(at)$; simply take Y to be a random variable which is everywhere equal to b .

Examples.

- (1) **Poisson distributions.** If X is $\text{Poisson}(\lambda)$, then

$$\begin{aligned}\varphi_X(t) &= \mathbf{E} [e^{itX}] = \sum_{k=0}^{\infty} e^{itk} \mathbf{P} \{X = k\} \\ &= \sum_{k=0}^{\infty} e^{itk} \frac{\lambda^k e^{-\lambda}}{k!} \\ &= e^{-\lambda} \sum_{k=0}^{\infty} \frac{(\lambda e^{it})^k}{k!} \\ &= e^{-\lambda} e^{\lambda e^{it}} = e^{\lambda(e^{it}-1)}.\end{aligned}$$

- (2) **Normal distributions.** Let X be $\text{Normal}(0, 1)$. To compute $\varphi_X(t)$ we will use the Taylor expansion

$$e^{itX} = \sum_{k \geq 0} \frac{(itX)^k}{k!};$$

we will justify the use of this formula in expectation calculations shortly. Assuming we can interchange expectation and sum without issue, we then have

$$\begin{aligned}\varphi_X(t) &= \mathbf{E} [e^{itX}] = \sum_{k \geq 0} \frac{(it)^k}{k!} \mathbf{E} [X^k] \\ &= \sum_{k \geq 0} \frac{(it)^{2k}}{(2k)!} \frac{(2k)!}{2^k k!},\end{aligned}$$

where we have used the formula for the moments of the normal distribution from Example 15.6; recall that the odd moments of X vanish. We thus have

$$\varphi_X(t) = \sum_{k \geq 0} \frac{1}{k!} \left(\frac{(it)^2}{2} \right)^k = \sum_{k \geq 0} \frac{1}{k!} \left(\frac{-t^2}{2} \right)^k = e^{-t^2/2}.$$

Note that $\sigma X + \mu$ is $\text{Normal}(\mu, \sigma^2)$, so the characteristic function of any normal can be easily derived from that of a $\text{Normal}(0, 1)$.

- (3) **Uniform distributions.** Let X be $\text{Uniform}[-a, a]$. Then X has density $1/(2a)$ on the interval $[-a, a]$, so

$$\varphi_X(t) = \mathbf{E} [e^{itX}] = \int_{-a}^a \frac{e^{itx}}{2a} dx = \frac{1}{2a} \left[\frac{e^{itx}}{it} \right]_{-a}^a = \frac{e^{ita} - e^{-ita}}{2iat} = \frac{\sin(at)}{at}.$$

- (4) **Symmetric simple random walk.** Let $(X_i, i \geq 1)$ be independent fair coin tosses, i.e., $\mathbf{P} \{X_i = 1\} = \mathbf{P} \{X_i = -1\} = 1/2$, and let $S_n = n^{-1/2} \sum_{i=1}^n X_i$. Then

$$\varphi_{X_i}(t) = \frac{1}{2}(e^{it} + e^{-it}) = \cos(t).$$

Using the factorization formula for characteristic functions, this gives

$$\varphi_{S_n}(t) = \mathbf{E} [e^{itS_n}] = \left(\mathbf{E} [e^{itX_1/n^{1/2}}] \right)^n = (\cos(t/n^{1/2}))^n,$$

so

$$\lim_{n \rightarrow \infty} \varphi_{S_n}(t) = \lim_{n \rightarrow \infty} (\cos(t/n^{1/2}))^n = \lim_{n \rightarrow \infty} \left(1 - \frac{t^2}{n} + O\left(\frac{1}{n^2}\right) \right)^n = e^{-t^2/2}.$$

In other words, the characteristic function of S_n converges pointwise to that of a Normal(0, 1) random variable.

The last example above is particularly suggestive; it hints that perhaps the distribution of S_n approximates that of a Normal(0, 1) when n is large. To turn this from a suggestive observation to a persuasive mathematical argument, we need to relate convergence of characteristic functions to convergence in distribution. This is done in the next section.

16.2. The inversion and continuity theorems. In this section we establish the following two theorems.

Theorem 16.5 (Inversion theorem for characteristic functions). *Fix any random variable X and write μ for the distribution of X . Then for all $a, b \in \mathbb{R}$ with $a < b$,*

$$\frac{1}{2\pi} \lim_{T \rightarrow \infty} \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi_X(t) dt = \mu(a, b) + \frac{1}{2}(\mu(\{a\}) + \mu(\{b\})).$$

In particular, μ can be recovered from φ_X .

Theorem 16.6 (Continuity theorem for characteristic functions). *Fix random variables $(X_n, 1 \leq n \leq \infty)$. If $X_n \xrightarrow{d} X_\infty$ then $\varphi_{X_n} \rightarrow \varphi_{X_\infty}$ pointwise. Conversely, if $\varphi_{X_n} \rightarrow \varphi$ pointwise and φ is continuous at 0, then $(X_n, n \geq 1)$ converges in distribution to a random variable Y satisfying $\varphi_Y = \varphi$.*

The above inversion theorem is an “integrated” version of Fourier inversion, which is the best we can hope for if X doesn’t have a density. When X does have a density, we can obtain something which looks more like (16.4).

Exercise 16.1. *Let X be a real random variable. Show that if X has a density then $|\varphi_X(t)| \rightarrow 0$ as $t \rightarrow \infty$. Also show that if $\int_{\mathbb{R}} |\varphi_X(t)| < \infty$ then the function defined by*

$$f(x) = \frac{1}{2\pi} \int_{\mathbb{R}} e^{-itx} \varphi_X(t) dt$$

is a continuous density for X .

If X has a density, then it follows from the above exercise and the inversion theorem that (ignoring whether or not Fubini can be applied)

$$\begin{aligned} \mu(a, b) &= \int_a^b f(x) dx = \frac{1}{2\pi} \int_a^b \int_{\mathbb{R}} e^{-itx} \varphi_X(t) dt dx \\ &= \frac{1}{2\pi} \int_{\mathbb{R}} \varphi_X(t) \int_a^b e^{-itx} dx dt = \frac{1}{2\pi} \int_{\mathbb{R}} \varphi_X(t) \frac{e^{-ita} - e^{-itb}}{it} dt. \end{aligned}$$

Before starting the proofs of these theorems, we develop some preliminary estimates for Taylor expansions of complex exponentials, which also belatedly justify the calculations from examples 2 and 4, above.

Proposition 16.7. *For any $x \in \mathbb{R}$ and integer $n \geq 0$,*

$$e^{ix} = \sum_{k=0}^n \frac{(ix)^k}{k!} + \frac{i^{n+1}}{n!} \int_0^x (x-s)^n e^{is} ds = \sum_{k=0}^{n+1} \frac{(ix)^k}{k!} + \frac{i^{n+1}}{n!} \int_0^x (x-s)^n (e^{is} - 1) ds$$

Proof. We begin with the first identity, and first consider the case $n = 0$. In this case we have

$$\frac{(ix)^0}{0!} + \frac{i}{0!} \int_0^x (x-s)^0 e^{is} ds = 1 + i \left[\frac{e^{is}}{i} \right]_{s=0}^x = 1 + e^{ix} - 1 = e^{ix}$$

as required. For $n > 0$, suppose inductively that the first identity holds for smaller values of n . Using integration by parts with $u = e^{is}$ and $dv = (x-s)^{n-1} ds$, we have

$$\int_0^x (x-s)^{n-1} e^{is} ds = \frac{x^n}{n} + \frac{i}{n} \int_0^x (x-s)^n e^{is} ds.$$

It follows by induction that

$$\begin{aligned} e^{ix} &= \sum_{k=0}^{n-1} \frac{(ix)^k}{k!} + \frac{i^n}{(n-1)!} \int_0^x (x-s)^{n-1} e^{is} ds \\ &= \sum_{k=0}^{n-1} \frac{(ix)^k}{k!} + \frac{i^n}{(n-1)!} \frac{x^n}{n} + \frac{i^n}{(n-1)!} \frac{i}{n} \int_0^x (x-s)^n e^{is} ds \\ &= \sum_{k=0}^n \frac{(ix)^k}{k!} + \frac{i^{n+1}}{n!} \int_0^x (x-s)^n e^{is} ds, \end{aligned}$$

as required.

Next, rearranging the above integration by parts formula gives that

$$\int_0^x (x-s)^n e^{is} ds = \frac{n}{i} \int_0^x (x-s)^{n-1} e^{is} ds - \frac{x^n}{i} = \frac{n}{i} \int_0^x (x-s)^{n-1} e^{is} ds - \frac{n}{i} \int_0^x (x-s)^{n-1} ds.$$

Substituting this formula for $\int_0^x (x-s)^n e^{is} ds$ into the first identity, the second identity is immediate. \square

Here are some useful bounds which follow from the above identities. The first bounds the error terms in the Taylor expansion of e^{ix} ; the second states a probabilistic consequence of the first.

Corollary 16.8. *For any $x \in \mathbb{R}$ and integer $n \geq 0$,*

$$\left| e^{ix} - \sum_{k=0}^n \frac{(ix)^k}{k!} \right| \leq \min \left(\frac{|x|^{n+1}}{(n+1)!}, \frac{2|x|^n}{n!} \right)$$

Proof. For the first term, use the first identity from Proposition 16.7 with the bound

$$\left| \frac{i^{n+1}}{n!} \int_0^x (x-s)^n e^{is} ds \right| \leq \frac{1}{n!} \int_0^x |(x-s)^n e^{is}| ds = \frac{|x|^{n+1}}{(n+1)!}.$$

For the second term, use the second identity from Proposition 16.7 (note: we use the identity with the upper limit of the sum equal to n , not $n+1$) with the bound

$$\left| \frac{i^n}{(n-1)!} \int_0^x (x-s)^{n-1} (e^{is} - 1) ds \right| \leq \frac{1}{(n-1)!} \int_0^x |(x-s)^{n-1} |e^{is} - 1| ds \leq \frac{2|x|^n}{n!}. \quad \square$$

Corollary 16.9. *Let X be a real random variable. For all $n \geq 0$, if $\mathbf{E}[|X|^n] < \infty$ then*

$$\left| \varphi_X(t) - \sum_{k=0}^n \frac{(it)^k}{k!} \mathbf{E}[X^k] \right| \leq \mathbf{E} \min \left(\frac{|tX|^{n+1}}{(n+1)!}, \frac{2|tX|^n}{n!} \right).$$

Thus, if $\lim_{n \rightarrow \infty} \frac{t^n}{n!} \mathbf{E}[|X|^n] = 0$ then

$$\varphi_X(t) = \sum_{k=0}^{\infty} \frac{(it)^k}{k!} \mathbf{E}[X^k].$$

Here are the cases $n = 1, 2$ of the above bound, written out explicitly.

- (1) If $\mathbf{E}|X| < \infty$ then $|\varphi_X(t) - 1 - it\mathbf{E}[X]| \leq \mathbf{E} \min((tX)^2/2, 2|tX|)$,
- (2) If $\mathbf{E}[X^2] < \infty$ then $|\varphi_X(t) - 1 - it\mathbf{E}[X] + t^2\mathbf{E}[X^2]| \leq \mathbf{E} \min(|tX|^3/6, (tX)^2)$.

At one point in the proof of the inversion theorem we will also need to know the value of the following famous integral.

Lemma 16.10.

$$\int_0^{\infty} \frac{\sin x}{x} dx := \lim_{t \rightarrow \infty} \int_0^t \frac{\sin(x)}{x} dx = \frac{\pi}{2}$$

Proof. To compute the integral we introduce a second parameter: set

$$I(b) := \int_0^\infty \frac{\sin x}{x} e^{-bx} dx.$$

Our goal is to recover $I(0)$. The point of adding the parameter b is that $\frac{\partial}{\partial b} \frac{\sin x}{x} e^{-bx} = -\sin(x)e^{-bx}$, which is integrable. We then have

$$I'(b) = \frac{\partial}{\partial b} \int_0^\infty \frac{\sin x}{x} e^{-bx} dx = \int_0^\infty \frac{\sin x}{x} \frac{\partial}{\partial b} e^{-bx} dx = - \int_0^\infty \sin(x) e^{-bx} dx.$$

Integration by parts (applied twice) gives that

$$\int_0^\infty \sin(x) e^{-bx} dx = \left[-\cos(x) e^{-bx} - \sin(x) \cdot b e^{-bx} \right]_0^\infty - \int_0^\infty b^2 e^{-bx} \sin x dx,$$

which on rearrangement gives that

$$- \int_0^\infty \sin(x) e^{-bx} dx = \left[e^{-bx} \cdot \frac{\cos x + b \sin x}{b^2 + 1} \right]_0^\infty = \frac{-1}{b^2 + 1}$$

We have calculated that $I'(b) = -1/(b^2 + 1) = -\frac{d}{db} \tan^{-1}(b)$, so $I(b) = -\tan^{-1}(b) + I(0)$. Taking $b \rightarrow \infty$, it follows that

$$0 = \lim_{b \rightarrow \infty} \int_0^\infty \sin(x) e^{-bx} dx = \lim_{b \rightarrow \infty} (-\tan^{-1}(b) + I(0)) = -\frac{\pi}{2} + I(0);$$

so $I(0) = \pi/2$. □

Proof of Theorem 16.5. First, since $\varphi_X(t) = \int_{\mathbb{R}} e^{itx} \mu(dx)$, we have

$$\int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi_X(t) dt = \int_{-T}^T \int_{\mathbb{R}} \frac{e^{-ita} - e^{-itb}}{it} e^{itx} \mu(dx) dt$$

We would like to change the order of integration; to apply Fubini's theorem we need to verify that the integrand is absolutely integrable on the domain of integration. This isn't too hard: the case $n = 0$ of Corollary 16.8 says that $|e^{ix} - 1| \leq x$, so

$$|e^{-ita} - e^{-itb}| = |e^{-itb}(e^{-it(a-b)} - 1)| = |e^{-it(a-b)} - 1| \leq t|a - b|,$$

and therefore

$$\int_{-T}^T \int_{\mathbb{R}} \left| \frac{e^{-ita} - e^{-itb}}{it} e^{itx} \right| \mu(dx) dt \leq \int_{-T}^T \int_{\mathbb{R}} |a - b| \mu(dx) dt = 2T|a - b| < \infty.$$

So applying Fubini's theorem is justified, and we obtain

$$\int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi_X(t) dt = \int_{\mathbb{R}} \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} e^{itx} dt \mu(dx).$$

We next manipulate the inner integral to obtain a real-valued integrand. Using that $e^{ix} = \cos x + i \sin x$ we have

$$\begin{aligned} & \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} e^{itx} dt \\ &= \int_{-T}^T \frac{\cos(t(x-a)) - \cos(t(x-b))}{it} + \frac{\sin(t(x-a)) - \sin(t(x-b))}{t} dt \end{aligned}$$

We wish to split the latter integral in two; to apply linearity of integration we need to know that the terms in the integrand are absolutely integrable. This is true for the first fraction since $|\cos(t)/it| =$

$|\cos(t)/t| = O(t)$ as $t \downarrow 0$. For the second it is true by Lemma 16.10. We thus have

$$\begin{aligned} & \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} e^{itx} dt \\ &= \int_{-T}^T \frac{\cos(t(x-a)) - \cos(t(x-b))}{it} dt + \int_{-T}^T \frac{\sin(t(x-a)) - \sin(t(x-b))}{t} dt \\ &= 2 \int_0^T \frac{\sin(t(x-a)) - \sin(t(x-b))}{t} dt, \end{aligned}$$

where in the last step we have used that \cos is even and that \sin is odd.

This identity implies that

$$\int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi_X(t) dt = 2 \int_{\mathbb{R}} \int_0^T \frac{\sin(t(x-a)) - \sin(t(x-b))}{t} dt \mu(dx),$$

so

$$\begin{aligned} & \lim_{T \rightarrow \infty} \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi_X(t) dt \\ &= \lim_{T \rightarrow \infty} 2 \int_{\mathbb{R}} \left(\int_0^T \frac{\sin(t(x-a))}{t} dt - \int_0^T \frac{\sin(t(x-b))}{t} dt \right) \mu(dx) \\ &= 2 \int_{\mathbb{R}} \left(\lim_{T \rightarrow \infty} \int_0^T \frac{\sin(t(x-a))}{t} dt - \lim_{T \rightarrow \infty} \int_0^T \frac{\sin(t(x-b))}{t} dt \right) \mu(dx). \end{aligned}$$

Lemma 16.10 implies the the inner integrals are bounded, so the bounded convergence theorem justifies moving the limit under the integral sign in the last equality above.

Finally, we have

$$\lim_{T \rightarrow \infty} \int_0^T \frac{\sin(t(x-a))}{t} dt = \lim_{T \rightarrow \infty} \int_0^T \frac{\sin(t(x-a))}{(x-a)t} \cdot (x-a) dt = \frac{\pi}{2} \text{sign}(x-a)$$

and likewise $\lim_{T \rightarrow \infty} \int_0^T \frac{\sin(t(x-b))}{t} dt = (x-b) \cdot \frac{\pi}{2} \text{sign}(x-b)$, so

$$\begin{aligned} \lim_{T \rightarrow \infty} \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi_X(t) dt &= 2 \int_{\mathbb{R}} \frac{\pi}{2} \text{sign}(x-a) - \frac{\pi}{2} \text{sign}(x-b) \mu(dx) \\ &= \pi \int_{\mathbb{R}} \text{sign}(x-a) - \text{sign}(x-b) \mu(dx). \end{aligned}$$

The integrand is 1 if $x \in \{a, b\}$, is 2 if $x \in (a, b)$, and is 0 otherwise; so the last expression is precisely $\pi(\mu(\{a\}) + \mu(\{b\})) + 2\pi\mu(a, b)$. \square

We now turn our attention to the inversion theorem. Its proof has two steps, one “hard” (more quantitative) and one “soft” (more qualitative). The hard step is to show that convergence of characteristic functions implies tightness of the family of random variables. The soft step is to show that tightness together with pointwise convergence of characteristic functions implies convergence in distribution; this step uses the fact that the map $X \mapsto \varphi_X$ is injective, which is implied by the inversion theorem.

We begin with the more quantitative step.

Lemma 16.11. *For any real random variable X , for all $u > 0$,*

$$\mathbf{P}\{|X| \geq 2/\mu\} \leq \frac{1}{u} \int_{-u}^u (1 - \varphi_X(t)) dt$$

Proof. Write μ for the distribution of X . If $|x| \geq 2/\mu$ then

$$2 \left(1 - \frac{\sin(ux)}{ux} \right) \geq 2 \left(1 - \frac{1}{|ux|} \right) \geq 1,$$

so

$$\mathbf{P} \left\{ |X| \geq \frac{2}{u} \right\} = \int_{|x| \geq 2/\mu} 1 \mu(dx) \leq 2 \int_{|x| \geq 2/\mu} \left(1 - \frac{1}{|ux|} \right) \mu(dx) \leq 2 \int_{|x| \geq 2/\mu} \left(1 - \frac{\sin(ux)}{ux} \right) \mu(dx).$$

Integrating over \mathbb{R} rather than over $\{x : |x| \geq 2/\mu\}$ can only increase the integral; using this and the fact that $2 \sin(ux)/x = \int_{-u}^u e^{itx} dt$, we obtain that

$$\begin{aligned} \mathbf{P} \left\{ |X| \geq \frac{2}{u} \right\} &\leq 2 \int_{\mathbb{R}} \left(1 - \frac{\sin(ux)}{ux} \right) \mu(dx) \\ &= 2 - \frac{1}{u} \int_{\mathbb{R}} \int_{-u}^u e^{itx} dt \mu(dx) \\ &= \frac{1}{u} \int_{-u}^u (1 - \varphi_X(t)) dt, \end{aligned}$$

where in the last step we have used Fubini's theorem. □

Here is the key consequence of the preceding bound. We say a collection $(X_i, i \in I)$ of random variables is tight if the corresponding family of probability measures is tight.

Corollary 16.12. *Fix real random variables $(X_n, n \geq 1)$. If φ_{X_n} converges pointwise to some function φ and φ is continuous at zero then $(X_n, n \geq 1)$ is a tight family.*

Proof. Fix $\epsilon > 0$. Since $\varphi_{X_n} \rightarrow \varphi$ pointwise, $\varphi(0) = 1$. Since φ is continuous at zero, we may thus choose $u > 0$ such that $\varphi(t) > 1 - \epsilon/2$ for $|t| < u$. Then

$$\frac{1}{u} \int_{-u}^u (1 - \varphi(t)) dt \leq \frac{1}{u} \int_{-u}^u \frac{\epsilon}{2} = \epsilon.$$

By the bounded convergence theorem,

$$\lim_{n \rightarrow \infty} \frac{1}{u} \int_{-u}^u (1 - \varphi_n(t)) dt = \frac{1}{u} \int_{-u}^u (1 - \varphi(t)) dt,$$

so there is $n_0 \in \mathbb{N}$ such that $\frac{1}{u} \int_{-u}^u (1 - \varphi_{X_n}(t)) dt < \epsilon$ for all $n > n_0$; by the lemma it follows that

$$\sup_{n > n_0} \mathbf{P} \{ |X_n| \geq 1/u \} < \epsilon.$$

Now choose $K > 1/u$ large enough that $\max_{n \geq n_0} \mathbf{P} \{ |X_n| \geq 1/u \} < \epsilon$; then by the previous bound we have $\sup_{n \geq 1} \mathbf{P} \{ |X_n| \geq K \} < \epsilon$. Since $\epsilon > 0$ was arbitrary, it follows that $(X_n, n \geq 1)$ is a tight family. □

We now turn to the more qualitative arguments. The first is a lemma which is broadly applicable and has nothing to do with characteristic functions.

Lemma 16.13. *If $(X_n, n \geq 1)$ is a tight family of random variables, then there exists an increasing sequence $(n_k, k \geq 1)$ such that $(X_{n_k}, k \geq 1)$ converges in distribution.*

Proof. Let F_n be the cumulative distribution function of X_n . By a diagonalization argument, there exists an increasing sequence $(n_k, k \geq 1)$ such that $F_{n_k}(q)$ converges for all $q \in \mathbb{Q}$. Call the subsequential limit $G(q)$, and define a function $F : \mathbb{R} \rightarrow [0, 1]$ by $F(r) = \inf \{ G(q) : q > r, q \in \mathbb{Q} \}$.

The function F is clearly non-decreasing. It is also right-continuous, because for any $\epsilon > 0$ there is $q > r$ such that $G(q) < F(r) + \epsilon$, so $F(s) < F(r) + \epsilon$ for $s \in (r, q)$. Moreover, for all $\epsilon > 0$ there is $M \in \mathbb{N}$ such that $\mathbf{P} \{ |X_n| \geq M \} < \epsilon$, so

$$F(-M - 1) \leq \sup_{k \geq 1} F_{n_k}(-M) < \epsilon \quad \text{and} \quad F(M) \geq \inf_{k \geq 1} F_{n_k}(M) > 1 - \epsilon.$$

It follows that $\lim_{M \rightarrow -\infty} F(M) = 0$ and $\lim_{M \rightarrow \infty} F(M) = 1$, so F is the cumulative distribution function of some random variable X .

Explain this when tightness is first introduced?

Finally, if F is continuous at x then

$$\liminf_{k \rightarrow \infty} F_{n_k}(x) \geq \sup_{s \in \mathbb{Q}, s < x} \lim_{k \rightarrow \infty} F_{n_k}(s) = \sup_{s \in \mathbb{Q}, s < x} G(s) = F(x)$$

and

$$\limsup_{k \rightarrow \infty} F_{n_k}(x) \leq \inf_{s \in \mathbb{Q}, s > x} \lim_{k \rightarrow \infty} F_{n_k}(s) = \inf_{s \in \mathbb{Q}, s > x} G(s) = F(x).$$

Thus $F_{n_k} \rightarrow F$ at continuity points of F ; so $X_{n_k} \xrightarrow{d} X$ (find lemma reference from earlier in the text for the last implication). \square

So is the second.

Lemma 16.14 (Subsubsequence principle). *Let $(X_n, 1 \leq n \leq \infty)$ be a tight family. Then $X_n \xrightarrow{d} X_\infty$ if and only if for any increasing sequence $(n_k, k \geq 1)$, if $X_{n_k} \xrightarrow{d} Y$ for some random variable Y , then $Y \stackrel{d}{=} X_\infty$.*

Proof. First suppose that $X_n \xrightarrow{d} X_\infty$. Then for any increasing sequence $(n_k, k \geq 1)$, $X_{n_k} \xrightarrow{d} X_\infty$; this shows that one of the two implications holds.

For the other implication, suppose that $X_n \not\xrightarrow{d} X_\infty$. Then by definition there is some $x \in \mathbb{R}$ such that $\mathbf{P}\{X_\infty = x\} = 0$ and such that $\mathbf{P}\{X_n \leq x\} \not\rightarrow \mathbf{P}\{X_\infty \leq x\}$. We may thus find $\epsilon > 0$ and an increasing sequence $(m_k, k \geq 1)$ such that

$$\inf_{k \geq 1} |\mathbf{P}\{X_{m_k} \leq x\} - \mathbf{P}\{X_\infty \leq x\}| > \epsilon.$$

Since $(X_n, n \geq 1)$ is tight, $(X_{m_k}, k \geq 1)$ is also tight, so by Lemma 16.13 there is an increasing subsequence $(n_k, k \geq 1)$ of $(m_k, k \geq 1)$ such that $X_{n_k} \xrightarrow{d} Y$ for some random variable Y . If $\mathbf{P}\{Y = x\} > 0$ then clearly $Y \not\stackrel{d}{=} X_\infty$, and if $\mathbf{P}\{Y = x\} = 0$ then

$$\mathbf{P}\{X_{n_k} \leq x\} \rightarrow \mathbf{P}\{Y \leq x\}$$

so $|\mathbf{P}\{Y \leq x\} - \mathbf{P}\{X_\infty \leq x\}| \geq \epsilon$ and again $X_\infty \not\stackrel{d}{=} Y$. This establishes the second implication. \square

We'll now put the two preceding lemmas together with information about convergence of characteristic functions to prove the continuity theorem.

Proof of Theorem 16.6. First, if $X_n \xrightarrow{d} X_\infty$ then for all $t \in \mathbb{R}$,

$$\mathbf{E}[\sin(tX_n)] \rightarrow \mathbf{E}[\sin(tX_\infty)] \quad \text{and} \quad \mathbf{E}[\cos(tX_n)] \rightarrow \mathbf{E}[\cos(tX_\infty)],$$

so $\varphi_{X_n}(t) = \mathbf{E}[e^{itX_n}] \rightarrow \mathbf{E}[e^{itX_\infty}] = \varphi_{X_\infty}(t)$.

Next, if $\varphi_{X_n} \rightarrow \varphi$ pointwise and φ is continuous at zero, then by Corollary 16.12, $(X_n, n \geq 1)$ is a tight family. Let $(n_k, k \geq 1)$ be any increasing sequence along which $X_{n_k} \xrightarrow{d} Y$ for some random variable Y . Then by the first part of the theorem, $\varphi_{X_{n_k}} \rightarrow \varphi_Y$ pointwise, so it must be that $\varphi_Y = \varphi$. In particular, φ is a characteristic function. Moreover, by the inversion theorem, the distribution of Y is uniquely determined by φ ; so if $(m_k, k \geq 1)$ is any other increasing sequence such that $(X_{m_k}, k \geq 1)$ converges in distribution, then it must be that $X_{m_k} \xrightarrow{d} Y$ as well. It follows by Lemma 16.14 that $(X_n, n \geq 1)$ converges in distribution and that the distributional limit has characteristic function φ . \square

16.3. The central limit theorem. Probably the most important use of characteristic functions is to prove the central limit theorem for iid random variables with finite second moment. Having proved the inversion and continuity theorems, we are almost ready to do so; there is just one more easy analytic estimate we will need in during the proof, which we state in advance: for all $a, b \in \mathbb{C}$ with $|a - b| \leq 1$, and all $n \in \mathbb{N}$,

$$|a^n - b^n| \leq n|a - b| \tag{16.5}$$

This is easily proved by induction. The case $n = 1$ is obvious; for the inductive step write $a^{n+1} - b^{n+1} = (a - b)a^n + b(a^n - b^n)$ and apply the triangle inequality.

Theorem 16.15 (Lindberg-Lévy Central Limit Theorem). *Let $(X_n, n \geq 1)$ be independent, identically distributed random variables with $X_1 \in L_2(\Omega, \mathcal{F}, \mathbf{P})$. Write $c = \mathbf{E}[X_1]$ and $\sigma^2 = \mathbf{Var}(X_1)$, and set $S_n = X_1 + \dots + X_n$ for $n \geq 1$. Then*

$$\frac{S_n - cn}{\sigma\sqrt{n}} \xrightarrow{d} N$$

where N is a $\text{Normal}(0, 1)$ random variable.

Proof. By an affine transformation (replacing X_i by $X'_i = X_i - c$) we may assume $c = 0$, in which case $\sigma^2 = \mathbf{E}[X_1^2]$.

Let $\varphi_N(t) = e^{-t^2/2}$ be the characteristic function of a $\text{Normal}(0, 1)$ random variable. By the continuity theorem it suffices to show that $\varphi_{S_n/\sigma\sqrt{n}} \rightarrow \varphi_N$ pointwise.

Writing $\varphi = \varphi_{X_1}$, then

$$\varphi_{S_n/\sigma\sqrt{n}}(t) = \mathbf{E}\left[e^{itS_n/\sqrt{n}}\right] = (\varphi(t/\sqrt{n}))^n,$$

so to prove convergence of $\varphi_{S_n/\sigma\sqrt{n}}(t)$ to $\varphi_N(t) = e^{-t^2/2}$, we need to control the behaviour of $\varphi(x)$ near $x = 0$.

By the case $n = 2$ of (16.9) we have

$$\left| \varphi(t) - \left(1 - \frac{t^2\sigma^2}{2}\right) \right| = \left| \varphi(t) - \left(1 - it\mathbf{E}X - \frac{t^2}{2}\mathbf{E}[X^2]\right) \right| \leq t^2 \mathbf{E} \min\left(\frac{|t|}{6}|X|^3, X^2\right).$$

Writing $\beta(t)$ for the final term on the right, by the dominated convergence theorem, $\beta(t)/t^2 = \mathbf{E} \min\left(\frac{|t|}{6}|X|^3, X^2\right) \rightarrow 0$ as $t \rightarrow 0$.

Using (16.5), we thus have

$$\begin{aligned} \left| \varphi_{S_n/\sigma\sqrt{n}}(t) - \left(1 - \frac{t^2}{2n}\right)^n \right| &= \left| \left(\varphi\left(\frac{t}{\sigma\sqrt{n}}\right)\right)^n - \left(1 - \frac{t^2}{2n}\right)^n \right| \\ &\leq n \left| \varphi\left(\frac{t}{\sigma\sqrt{n}}\right) - \left(1 - \frac{t^2}{2n}\right) \right| \\ &= n\beta\left(\frac{t}{\sigma\sqrt{n}}\right) \rightarrow 0 \end{aligned}$$

the convergence holding as $n \rightarrow \infty$ by the previously established control on the behaviour of β near zero.

Since $(1 - t^2/(2n))^n \rightarrow e^{-t^2/2}$ as $n \rightarrow \infty$, it follows that $\varphi_{S_n/\sigma\sqrt{n}}(t) \rightarrow e^{-t^2/2} = \varphi_N(t)$, as required. \square

16.4. Characteristic functions and moments. It was mentioned in Section 15.3 that knowledge of the moments of a random variable X is in general not sufficient to determine the distribution of X . This is even true for distributions that arise in “real-world” situations: if N is a $\text{Normal}(\mu, \sigma^2)$ distributed then e^N has all moments finite but its distribution is not determined by its moments. There are even uncountably many purely discrete distributions with the same moments as e^N ; for more on this see [1, 3].

The paragraph and the three after it can safely be skipped without detracting from the readability of the rest. Plan for this paragraph: describe the situation for log-normal distributions in more detail.

Note that if μ and ν are two distributions with the same moments then $p\mu + (1 - p)\nu$ also has the same moments and, more general, the set of distributions with given moments is always convex. There is a theorem of Riesz (see [3], Theorem 2.14) which says that the extremal points of the set

of distributions with given moments are precisely the distributions for which there is a Parseval's identity. More precisely, fix a sequence of real values $\mathbf{m} = (m_n, n \geq 0)$ and let

$$\mathcal{D}_{\mathbf{m}} = \left\{ \mu : \mu \text{ is a probability measure on } (\mathbb{R}, \mathcal{B}(\mathbb{R})); \forall n, \int_{\mathbb{R}} t^n \mu(dt) = m_n \right\}.$$

be the set of probability distributions on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ with moments given by \mathbf{m} .

If the set $\mathcal{D}_{\mathbf{m}}$ is non-empty then the moments \mathbf{m} uniquely determine a set of *orthogonal polynomials* $(p_n(x), n \geq 0)$ with p_n a degree- n polynomial, such that $\int_{\mathbb{R}} p_n(x)p_m(x)\mu(dx) = 0$ for all $0 \leq m \neq n$ and $\int_{\mathbb{R}} p_n(x)^2\mu(dx) = 1$ for all $n \geq 0$. The definition of the set $\mathcal{D}_{\mathbf{m}}$ means it doesn't matter which measure $\mu \in \mathcal{D}_{\mathbf{m}}$ is used for the preceding integrals; if the identities hold for one measure in $\mathcal{D}_{\mathbf{m}}$ then they hold for all measures in $\mathcal{D}_{\mathbf{m}}$.

Write $L_2(\mu)$ for the set of Borel functions $f : \mathbb{R} \rightarrow \mathbb{C}$ such that $\int_{\mathbb{R}} |f(t)|^2\mu(dt) < \infty$. Think of $(p_n)_{n \geq 0}$ as a sort of orthonormal basis for $L_2(\mu)$. (So maybe the p_n should be complex-valued) For $f \in L_2(\mu)$ write

$$\hat{f}(n, \mu) = \int_{\mathbb{R}} f(t)p_n(t)\mu(dt)$$

for the "Fourier coefficient" of f relative to $(p_n)_{n \geq 0}$. Then μ is an extreme point of $\mathcal{D}_{\mathbf{m}}$ if and only if for all $f \in L_2(\mu)$,

$$\int_{\mathbb{R}} |f(t)|^2\mu(dt) = \sum_{n=0}^{\infty} |\hat{f}(n, \mu)|^2.$$

We see that the connection of moment problems to complex analysis runs deep. Though we won't exhaustively explore the connection, we will at least prove the following result, whose proof proceed via characteristic functions.

Theorem 16.16. *For a real random variable X , if G_X has a positive radius of convergence then \mathcal{L}_X is determined by the moments of X .*

We will prove Theorem 16.16 by showing that if G_X has a positive radius of convergence then φ_X is determined by the moments of X ; we can then apply the inversion theorem to recover \mathcal{L}_X from φ_X . To carry this out, we need the following lemma, a complex-valued analogue of the first part of Theorem 15.1. Throughout what follows, we'll work with a fixed random variable X and write $\varphi = \varphi_X$.

Lemma 16.17. *For all $k \in \mathbb{N}$, if $\mathbf{E}[|X|^k] < \infty$ then for all $t \in \mathbb{R}$ we have $\varphi^{(k)}(t) = \mathbf{E}[(iX)^k e^{itX}]$.*

Proof. We proceed by induction. For all $t \in \mathbb{R}$ and $h > 0$, we have

$$\begin{aligned} \frac{\varphi(t+h) - \varphi(t)}{h} - iX e^{itX} &= \frac{e^{i(t+h)X} - e^{itX}}{h} - iX e^{itX} \\ &= e^{itX} \frac{e^{ihX} - 1 - ihX}{h}, \end{aligned}$$

so

$$\left| \frac{\varphi(t+h) - \varphi(t)}{h} - iX e^{itX} \right| \leq \left| e^{itX} \frac{e^{ihX} - 1 - ihX}{h} \right| \leq \min \left(\frac{h|X|^2}{2}, 2|X| \right),$$

where the inequality follows from the case $n = 1$ of Corollary 16.8. If $\mathbf{E}|X| < \infty$ then the second of the two bounds implies that $\left| \frac{\varphi(t+h) - \varphi(t)}{h} - iX e^{itX} \right|$ is dominated by an integrable random variable, so using the first of the two bounds and the dominated convergence theorem we obtain that

$$\limsup_{h \rightarrow 0} \left| \frac{\varphi(t+h) - \varphi(t)}{h} - iX e^{itX} \right| \leq \mathbf{E} \left[\limsup_{h \rightarrow 0} \frac{h}{2} |X|^2 \right] = 0.$$

This establishes the assertion of the lemma when $k = 1$.

Now fix $k \geq 1$, suppose $\mathbf{E}[|X|^{k+1}] < \infty$, and assume by induction that $\varphi^{(k)}(t) = \mathbf{E}[(iX)^k e^{itX}]$. Then

$$\begin{aligned} \left| \frac{\varphi^{(k)}(t+h) - \varphi^{(k)}(t)}{h} - (iX)^{k+1} e^{itX} \right| &= \left| \frac{(iX)^k e^{i(t+h)X} - (iX)^k e^{itX}}{h} - (iX)^{k+1} e^{itX} \right| \\ &= \left| (iX)^k e^{itX} \frac{e^{ihX} - 1 - ihX}{h} \right| \\ &\leq \left| (iX)^k \min\left(\frac{h|X|^2}{2}, 2|X|\right) \right| \\ &= \min\left(\frac{h}{2}|X|^{k+1}, 2|X|^k\right). \end{aligned}$$

From this bound, arguing using the dominated convergence theorem as before shows that $\varphi^{(k+1)}(t) = \mathbf{E}[(iX)^{k+1} e^{itX}]$. \square

We will also use the following lemma which provides a bound on the growth rate of the absolute moments of a random variable whose moment generating function is finite in a neighbourhood of the origin.

Lemma 16.18. *If $G_X(s) < \infty$ and $G(-s) < \infty$ then for all $r \in (0, s)$,*

$$\frac{r^n \mathbf{E}[|X|^n]}{n!} \rightarrow 0,$$

as $n \rightarrow \infty$.

Proof. Under the assumptions of the lemma, for all $r \in (0, s)$ we have $G_X(r) < \infty$ and $G_X(-r) < \infty$, so by Proposition 15.2 we have

$$G_X(r) = \sum_{n \geq 0} \mathbf{E}[X^n] \frac{r^n}{n!};$$

so $|\mathbf{E}[X^n] \frac{r^n}{n!}| \rightarrow 0$ as $n \rightarrow \infty$. It follows that $\mathbf{E}[|X|^{2k}] \frac{r^{2k}}{(2k)!} \rightarrow 0$ as $k \rightarrow \infty$, since for even values the moments and absolute moments agree.

To handle odd values, note that for $k \geq 1$ we have $|rX|^{2k-1} \leq 1 + |rX|^{2k}$. Moreover, for k sufficiently large, $r^{2k} < s^{2k}/2k$, and for such k we have

$$r^{2k-1} \mathbf{E}[|X|^{2k-1}] = \mathbf{E}[|rX|^{2k-1}] \leq 1 + \mathbf{E}[|rX|^{2k}] \leq 1 + \frac{s^k}{2k} \mathbf{E}[|X|^{2k}].$$

It follows that

$$\limsup_{k \rightarrow \infty} \frac{r^{2k-1} \mathbf{E}[|X|^{2k-1}]}{(2k-1)!} \leq \limsup_{k \rightarrow \infty} \frac{1}{(2k-1)!} \left(1 + \frac{s^k}{2k} \mathbf{E}[|X|^{2k}] \right) = \limsup_{k \rightarrow \infty} \frac{s^k \mathbf{E}[|X|^{2k}]}{(2k)!}$$

and the last limit is zero as noted above. \square

Proof of Theorem 16.16. Fix $r > 0$ strictly inside the radius of convergence of G_X . By Corollary 16.8, for all $n \in \mathbb{N}$ and $x \in \mathbb{R}$ we have

$$\left| e^{ihx} - \sum_{k=0}^n \frac{(ihx)^k}{k!} \right| \leq \frac{|hx|^{n+1}}{(n+1)!},$$

so for $t, h \in \mathbb{R}$ with $|h| \leq r$,

$$\begin{aligned} \left| \varphi(t+h) - \sum_{k \leq n} \frac{\varphi^{(k)}(t)}{k!} h^k \right| &= \left| \mathbf{E} \left[e^{i(t+h)X} \right] - \sum_{k \leq n} \frac{h^k}{k!} \mathbf{E} \left[(iX)^k e^{itX} \right] \right| \\ &= \left| \mathbf{E} \left[e^{i(t+h)X} \cdot \left(e^{ihX} - \sum_{k \leq n} \frac{h^k}{k!} (iX)^k \right) \right] \right| \\ &\leq \mathbf{E} \left| e^{ihX} - \sum_{k \leq n} \frac{h^k}{k!} (iX)^k \right| \\ &\leq \mathbf{E} \frac{|hX|^{n+1}}{(n+1)!} = \frac{h^{n+1}}{(n+1)!} \mathbf{E} [|X|^{n+1}]. \end{aligned}$$

The final quantity tends to zero as $n \rightarrow \infty$ by Lemma 16.18, and it follows that for such t and h ,

$$\varphi(t+h) = \sum_{k=0}^{\infty} \frac{\varphi^{(k)}(t)}{k!} h^k. \quad (16.6)$$

We now argue by induction that for all $a \in \mathbb{N}$, the moments of X determine $\varphi(x)$ for all $x \in [-ar, ar]$. For $a = 1$, taking $t = 0$ in (16.6) and using the formula for the derivatives of φ given in Lemma 16.17, we obtain that for $|h| \leq r$,

$$\varphi(h) = \sum_{k=0}^{\infty} \frac{\varphi^{(k)}(0)}{k!} h^k = \sum_{k \geq 0} \frac{\mathbf{E} [(iX)^k]}{k!} h^k,$$

so $\varphi(h)$ is determined by the moments of X . Supposing the claim is true for a given value of a , fix x with $|x| \in (ar, (a+1)r]$. Then we may write $x = t+h$ with $|t| \leq ar$ and $|h| \leq r$, and by (16.6) we obtain that

$$\varphi(x) = \varphi(t+h) = \sum_{k=0}^{\infty} \frac{\varphi^{(k)}(t)}{k!} h^k.$$

By induction, the values $\varphi^{(k)}(t)$ are determined by the moments of X ; it follows that $\varphi(x)$ is determined by the moments of X as well. \square

That's it for characteristic functions in this section. There are two goals for the rest of the section. The first is to relate convergence of moments to convergence in distribution. The second is to use that relation to provide a second proof of a special case of the central limit theorem.

Theorem 16.19. Fix random variables $(X_n, 1 \leq n \leq \infty)$ such that $\mathbf{E} [|X_n|^r] < \infty$ for all $r > 0$. If \mathcal{L}_{X_∞} is determined by the moments of X_∞ , and $\mathbf{E} [X_n^r] \rightarrow \mathbf{E} [X_\infty^r]$ for all $r \in \mathbb{N}$, then $X_n \xrightarrow{d} X_\infty$.

Proof. Set $K := \sup_{1 \leq n < \infty} \mathbf{E} [X_n^2]$. Since $\mathbf{E} [X_n^2] \rightarrow \mathbf{E} [X_\infty^2] < \infty$, we have $K < \infty$, so for any M ,

$$\sup_{n \in \mathbb{N}} \mathbf{P} \{|X_n| > M\} \leq \sup_{n \in \mathbb{N}} \frac{\mathbf{E} [X_n^2]}{M^2} \leq \frac{K}{M^2}.$$

For any $\epsilon > 0$, it follows that $\sup_{n \in \mathbb{N}} \mathbf{P} \{|X_n| > K/\epsilon^{1/2}\} \leq \epsilon$, so $(X_n, 1 \leq n < \infty)$ is tight.

A similar argument establishes uniform integrability of any fixed powers of the random variables in the sequence. Recall that for a sequence $(Y_n, n \geq 1)$ of random variables, if $\sup_{n \geq 1} \mathbf{E} [Y_n^t] < \infty$ for some $t > 1$ then $(Y_n, n \geq 1)$ is uniformly integrable; we saw this in the course of proving Theorem 13.15. Now fix integer $p \geq 1$; then $\sup_{n \geq 1} \mathbf{E} [X_n^{2p}] < \infty$, so $(X_n^p, n \geq 1)$ is uniformly integrable. Since also $X_n^p \xrightarrow{d} X_\infty^p$, it follows that

$$\mathbf{E} [X_n^p] \rightarrow \mathbf{E} [X_\infty^p].$$

We now apply the Lemma 16.14 to conclude that $X_n \xrightarrow{d} X_\infty$. Fix any increasing sequence $(n_k, k \geq 1)$ along which $X_n \xrightarrow{d} Z$ for some random variable Z . Then for any integer $p > 1$, by the uniform integrability of $(X_{n_k}^p, k \geq 1)$ it follows that $\mathbf{E}[X_{n_k}^p] \rightarrow \mathbf{E}[Z^p]$. Thus $\mathbf{E}[Z^p] = \mathbf{E}[X_\infty^p]$ for all p ; since \mathcal{L}_{X_∞} is determined by the moments of X_∞ , it follows that $\mathcal{L}_Z = \mathcal{L}_{X_\infty}$. We have shown that any subsequential distributional limit of X_n has the same distribution as X_∞ , so Lemma 16.14 then implies that $X_n \xrightarrow{d} X_\infty$. \square

Theorem 16.20 (Central limit theorem via moments). *Let $(X_n, n \geq 1)$ be independent identically distributed random variables with $\mathbf{E}[X_1] = 0$, $\mathbf{E}[X_1^2] = 1$, and set $\bar{S}_n := n^{-1/2}S_n = (X_1 + \dots + X_n)/n^{1/2}$. If $\mathbf{E}[|X_1|^p] < \infty$ for all $p > 0$ then \bar{S}_n converges in distribution to a Normal(0, 1) as $n \rightarrow \infty$.*

The key step in the proof is the following proposition, which is independently interesting.

Proposition 16.21. *Let $(X_n, n \geq 1)$ be independent identically distributed random variables with $\mathbf{E}[X_1] = 0$, $\mathbf{E}[X_1^2] = 1$, and set $\bar{S}_n := n^{-1/2}S_n = (X_1 + \dots + X_n)/n^{1/2}$. For any $p \in \mathbb{N}$, if $\mathbf{E}[|X_1|^p] < \infty$ then $L_p := \lim_{n \rightarrow \infty} \mathbf{E}[\bar{S}_n^p]$ exists, and*

$$L_p = \begin{cases} \frac{p!}{2^{p/2}(p/2)!} & \text{if } p \text{ is even} \\ 0 & \text{if } p \text{ is odd.} \end{cases}$$

Proof of Theorem . Let N be a Normal(0, 1)-distributed. Then $\mathbf{E}N^p = L_p$ for all $p \in \mathbb{N}$ (cite example from earlier in text). Moreover, G_N has a positive radius of convergence (in fact it converges everywhere), so \mathcal{L}_N is determined by its moments by Theorem 16.16. Proposition 16.21 says that all moments of $(\bar{S}_n, n \geq 1)$ converge to those of N , and it follows by Theorem 16.19 that $\bar{S}_n \xrightarrow{d} N$. \square

Proof of Proposition 16.21. The claim is obvious for $p = 1$ since $\mathbf{E}[\bar{S}_n] = 0$ for all n . It is also easy for $p = 2$ since $\mathbf{E}[S_n^2] = \sum_{i=1}^n \mathbf{E}[X_i^2] = n$ so $\mathbf{E}[\bar{S}_n^2] = 1 = L_2$.

Now fix $p \geq 2$, and suppose inductively that the claim is true for all $1 \leq q \leq p$. If $\mathbf{E}[|X_1|^{p+1}] < \infty$ then for any $n \in \mathbb{N}$,

$$\mathbf{E}[|S_n^{p+1}|] \leq \sum_{i_1, \dots, i_{p+1}=1}^n \mathbf{E}\left[\prod_{j=1}^p +1X_j^{i_j}\right] < \infty,$$

the last identity holding by the factorization formula since in each term in the sum, any single random variable X_i shows up at most $p + 1$ times.

We now write

$$S_n^{p+1} = S_n^p(X_1 + \dots + X_n);$$

by linearity of expectation and since (X_1, \dots, X_n) are identically distributed we then have

$$\mathbf{E}[S_n^{p+1}] = \sum_{i=1}^n \mathbf{E}[S_n^p X_i] = n\mathbf{E}[S_n^p X_n] = n\mathbf{E}[(S_{n-1} + X_n)^p X_n].$$

Applying the binomial expansion to $(S_{n-1} + X_n)^p$ it follows that

$$\mathbf{E}[S_n^{p+1}] = n \sum_{j=0}^p \binom{p}{j} \mathbf{E}[S_{n-1}^j X_n^{p-j} X_n] = n \sum_{j=0}^p \binom{p}{j} \mathbf{E}[S_{n-1}^j] \mathbf{E}[X_n^{p+1-j}].$$

The $j = p$ term of the sum is zero since $\mathbf{E}[X_n] = 0$. The $j = p - 1$ term of the sum is $p\mathbf{E}[S_{n-1}^{p-1}] \mathbf{E}[X_n^2] = \mathbf{E}[S_{n-1}^{p-1}]$. So the above identity may be rewritten as

$$\mathbf{E}[S_n^{p+1}] = n p \mathbf{E}[S_n^{p-1}] + n \sum_{j=0}^{p-2} \binom{p}{j} \mathbf{E}[S_{n-1}^j] \mathbf{E}[X_n^{p+1-j}].$$

To obtain an identity for $\mathbf{E} \left[\overline{S}_n^{p+1} \right]$ we divide through by $n^{(p+1)/2}$, which yields

$$\begin{aligned} \mathbf{E} \left[\overline{S}_n^{p+1} \right] &= p \left(\frac{n-1}{n} \right)^{\frac{p-1}{2}} \mathbf{E} \left[\overline{S}_n^{p-1} \right] + \sum_{j=0}^{p-2} \binom{p}{j} \frac{1}{n^{\frac{p-1-j}{2}}} \left(\frac{n-1}{n} \right)^{\frac{j}{2}} \mathbf{E} \left[X_n^{p+1-j} \right] \mathbf{E} \left[\overline{S}_n^j \right] \\ &\rightarrow pL_{p-1} = L_{p+1}. \end{aligned}$$

□

17. Weak convergence

17.1. Measures on metric spaces; the Portmanteau theorem. In this section we develop the theory of convergence of probability measures on metric spaces. Throughout what follows, we fix a metric space $M = (M, d)$, and write \mathcal{B}_M for the Borel σ -field on M . By a probability measure on M we mean a probability measure on (M, \mathcal{B}_M) .

$M = (M, d)$ metric space
 \mathcal{B}_M Borel σ -field

An M -valued random variable is a $(\mathcal{F}/\mathcal{B}_M)$ -measurable map $X : \Omega \rightarrow M$, where $(\Omega, \mathcal{F}, \mathbf{P})$ is a probability space. The law \mathcal{L}_X of X is the push-forward $X_*\mathbf{P}$; in other words, $\mathcal{L}_X(B) = \mathbf{P}\{X \in B\} = \mathbf{P}\{X^{-1}(B)\}$ for $B \in \mathcal{B}_M$. Note that if $I : M \rightarrow M$ is the identity map then I has law \mathcal{L}_X on the probability space $(M, \mathcal{B}_M, \mathcal{L}_X)$; this is the basis of the change of variables formula which says that if $f : M \rightarrow \mathbb{R}$ is measurable and $f(X)$ is non-negative or integrable, then

$$\mathbf{E}[f(X)] = \int_M f d\mathcal{L}_X.$$

The following is a basic property of probability measures on metric spaces.

Proposition 17.1. *Every probability measure μ on (M, \mathcal{B}_M) is regular: that is to say, for all $A \in \mathcal{B}_M$ and all $\epsilon > 0$ there exist $F \subseteq A \subseteq G$ with F closed and G open such that $\mu(G \setminus F) < \epsilon$.*

Proof. Let \mathcal{G} be the collection of regular sets in \mathcal{B}_M . We show that $\mathcal{G} = \mathcal{B}_M$ by proving that \mathcal{G} is a σ -field containing the closed sets. Note that \mathcal{G} is clearly closed under complements.

Fix $A \in \mathcal{B}_M$ closed, let $F = A$, and let $G_n = B(A, 1/n) = \{x \in M : d(x, A) < 1/n\}$. Then G_n is open, and since $G_n \downarrow A$, it follows by dominated convergence that $\mu(G_n \setminus F) = \mu(G_n \setminus A) \downarrow 0$. Thus $A \in \mathcal{G}$ so \mathcal{G} contains the closed sets.

Finally, fix any sequence $(A_n, n \geq 1)$ of elements of \mathcal{G} and let $A = \bigcup_{n \geq 1} A_n$. Given $\epsilon > 0$, for each $n \geq 1$ choose $F_n \subseteq A_n \subseteq G_n$ such that $\mu(G_n \setminus F_n) \leq \epsilon/2^{n+1}$. Set $F' = \bigcup_{n \geq 1} F_n$ and $G = \bigcup_{n \geq 1} G_n$, and choose n_0 large enough that $\mu(F' \setminus \bigcup_{n \leq n_0} F_n) < \epsilon/2$. Taking $F = \bigcup_{n \leq n_0} F_n$, it follows that $F \subset A \subset G$, and

$$\mu(G \setminus F) \leq \mu(F' \setminus F) + \mu(G \setminus F') \leq \frac{\epsilon}{2} + \sum_{n \geq 1} \frac{\epsilon}{2^{n+1}} = \epsilon \quad \square$$

We say that $\mathcal{A} \subset \mathcal{B}_M$ is a *separating class* if for any two probability measures P, Q on M , if $P(A) = Q(A)$ for all $A \in \mathcal{A}$ then $P = Q$. The preceding proposition implies that the collection of closed sets in M is a separating class (exercise). We saw earlier in the term that the collection $\{(-\infty, q], q \in \mathbb{Q}\}$ is a separating class for $(\mathbb{R}, \mathcal{B}_{\mathbb{R}})$.

separating class

The next proposition shows that expectations of bounded continuous test functions characterize probability measures on metric spaces. We write $C_b(M)$ for the set of bounded continuous functions $f : M \rightarrow \mathbb{R}$.

Proposition 17.2. *Fix P, Q two probability measures on M . If $\int f dP = \int f dQ$ for all $f \in C_b(M)$ then $P = Q$.*

Proof. By Proposition 17.1, it suffices to show that $P(F) = Q(F)$ for all closed sets $F \subset M$. Fix such F and for $n \geq 1$ define $f_n : M \rightarrow [0, 1]$ by

$$f_n(x) = \begin{cases} 1 - nd(x, F) & \text{if } d(x, F) \leq 1/n \\ 0 & \text{otherwise.} \end{cases}$$

Then the functions $(f_n, n \geq 1)$ bounded and continuous and are pointwise decreasing to the function $\mathbf{1}_{[F]}$, so by dominated convergence and the hypothesis of the proposition,

$$P(F) = \int \mathbf{1}_{[F]} dP = \lim_{n \rightarrow \infty} \int f_n dP = \int f_n dQ = \int IF dQ = Q(F). \quad \square$$

We pause to introduce the notation $Pf := \int f dP$, where P is a probability measure and f is a non-negative or P -integrable function. Here is the fundamental definition of the section.

Definition 17.3. Given probability measures $(P_n, n \geq 1)$ and P on M , say P_n converges in distribution to P , and write $P_n \Rightarrow P$, if $P_n f \rightarrow P f$ for all $f \in C_b(M)$.

Exercise 17.1. Show that if $P_n \Rightarrow P$ and $P_n \Rightarrow Q$ then $P = Q$.

Theorem 17.4 (Portmanteau theorem). Given probability measures $(P_n, n \geq 1)$ and P on M , the following are equivalent.

- (1) $P_n f \rightarrow P f$ for all $f \in C_b(M)$.
- (2) $P_n f \rightarrow P f$ for all uniformly continuous f in $C_b(M)$.
- (3) $\limsup_{n \rightarrow \infty} P_n(F) \leq P(F)$ for all $F \subset M$ closed.
- (4) $\liminf_{n \rightarrow \infty} P_n(G) \geq P(G)$ for all $G \subset M$ open.
- (5) $\lim_{n \rightarrow \infty} P_n(A) = P(A)$ for all $A \in \mathcal{B}_M$ with $P(\partial A) = 0$.

Proof. We begin with some easy implications. First note that (1) is equivalent to the following condition.

$$(1') \quad P_n f \rightarrow P f \text{ for all continuous } f : M \rightarrow [0, 1].$$

Clearly (1) implies (1'); for the reverse implication, replace f by $(f/\text{range}(f)) - \inf(f)$. Next, clearly (1) implies (2). Also, (3) and (4) are clearly equivalent (take complements). To conclude the proof we show that (2) implies (3), that (3) and (4) together imply (5), and that (5) implies (1').

Suppose that (2) holds, and fix any closed set $F \subset M$ and $\delta > 0$. Let m be large enough that $G := B(F, 1/m)$ has $P(G) < P(F) + \delta$, and let

$$f(x) = \begin{cases} 1 - md(x, F) & \text{if } d(x, F) \leq 1/m \\ 0 & \text{otherwise.} \end{cases}$$

Then f is uniformly continuous and $\mathbf{1}_{[F]} \leq f \leq \mathbf{1}_{[G]}$, so

$$\limsup_{n \rightarrow \infty} P_n(F) = \limsup_{n \rightarrow \infty} P_n \mathbf{1}_{[F]} \leq \lim_{n \rightarrow \infty} P_n f = P f \leq P(G) = P(F) + \delta;$$

since $\delta > 0$ was arbitrary, (3) follows.

Suppose (3) and (4) hold. Fix $A \in \mathcal{B}_M$ and write $A^\circ = A \setminus \partial A$ and $A^+ = A \cup \partial A$ for the interior and closure of A , respectively. Then

$$P(A^\circ) \leq \liminf_{n \rightarrow \infty} P_n(A^\circ) \leq \liminf_{n \rightarrow \infty} P_n(A) \leq \limsup_{n \rightarrow \infty} P_n(A) \leq \limsup_{n \rightarrow \infty} P_n(A^+) \leq P(A^+),$$

where the first and last inequalities follow from (4) and (3), respectively. If $P(\partial A) = 0$ then $P(A^\circ) = P(A^+) = P(A)$, and the above bounds give $\lim_{n \rightarrow \infty} P_n(A) = P(A)$.

Finally, suppose (5) holds and fix $f : M \rightarrow [0, 1]$ continuous. Then $\partial\{f > t\} \subset \{f = t\}$, so if $P(\partial\{f > t\}) > 0$ then t is an atom of $P_* f$. Finite probability measures have at most countably many atoms, so it follows that $P(\{f > t\}) = 0$ for Lebesgue-a.e. $t \in [0, 1]$. By (5), it follows that $P_n(\{f > t\}) \rightarrow P(\{f > t\})$ for Lebesgue-a.e. $t \in [0, 1]$, so using Fubini's theorem,

$$P_n f = \int_{[0,1]} P_n(\{f > t\}) dt \rightarrow \int_{[0,1]} P(\{f > t\}) dt = P f,$$

as $n \rightarrow \infty$. □

Exercise 17.2 (Subsubsequence principle). Show that $P_n \Rightarrow P$ if and only if for all $(n_k, k \geq 1)$ increasing sequences of positive integers, there exists a subsequence $(m_k, k \geq 1)$ such that $P_{m_k} \Rightarrow P$ as $k \rightarrow \infty$.

$Pf := \int f dP$

A°, A^+

17.2. Weights and measures. A metric space $M = (M, d)$ is called a *Polish space* if it is complete and separable. One of the most important themes of the next part of the notes is how choosing the right topology for Polish spaces yields a nice theory of measure.

Polish space

Throughout this section, unless otherwise stated, $M = (M, d)$ is a Polish space, (and $x = \{x_i, i \geq 1\} \subset M$ is a countable dense set?) The set of probability measures on (M, \mathcal{M}) is denoted $\text{prob}(M)$.

 $x = \{x_i, i \geq 1\}$
 countable, dense set
 $\text{prob}(M)$

We write $K \subset\subset M$ to mean that $K \subset M$ and K is compact. Recall that if $F \subset M$ is closed and $K \subset\subset M$ then $K \cap F$ is compact.

The relationship between compactness and measures on Polish spaces is crucial for what follows. A family $\mathcal{C} \subset \text{prob}(M)$ is *tight* if for all $\epsilon > 0$ there is $K \subset\subset M$ such that $\inf_{P \in \mathcal{C}} P(K) > 1 - \epsilon$.

tight

Proposition 17.5. *Let $M = (M, d)$ be a Polish space. Then $\{P\}$ is tight for all $P \in \text{prob}(M)$.*

Proof. For each n we have $\bigcup_{i \geq 1} B(x_i, 2^{-n}) = M$ since x is dense, so there is $i(n) \in \mathbb{N}$ such that

$$\mathbf{P} \left\{ \bigcup_{i > i(n)} B(x_i, 2^{-(n+1)}) \right\} < \frac{\epsilon}{2^n}.$$

Let $K = (\bigcap_{n \geq 1} \bigcup_{i \leq i(n)} B(x_i, 2^{-(n+1)}))^+$; recall that the superscript $+$ denotes closure. Then since the complement of $\bigcup_{i \leq i(n)} B(x_i, 2^{-(n+1)})$ is contained in $\bigcup_{i > i(n)} B(x_i, 2^{-(n+1)})$, we have

$$\mathbf{P} \{K^c\} \leq \mathbf{P} \left\{ \bigcup_{n \geq 1} \bigcup_{i > i(n)} B(x_i, 2^{-(n+1)}) \right\} \leq \sum_{n \geq 1} \frac{\epsilon}{2^n} = \epsilon,$$

so to prove the proposition it suffices to show that K is compact.

We prove this by contradiction; suppose $G = (G_j, j \in J)$ is an open cover of K with no finite subcover. Now, $K \subseteq \bigcup_{i \leq i(1)} B(x_i, 2^{-1})$, so there must be $j(1) \leq i(1)$ such that $K \cap B(x_{j(1)}, 2^{-1})$ has no finite subcover from G . However, $K \cap B(x_{j(1)}, 2^{-1}) \subset \bigcup_{i \leq i(2)} B(x_i, 2^{-2})$ so there must be $j(2) \leq i(2)$ such that

$$K \cap B(x_{j(1)}, 2^{-1}) \cap B(x_{j(2)}, 2^{-2})$$

has no finite subcover from G . Proceeding inductively, for each $n \geq 1$ we may choose $j(n) \leq i(n)$ such that

$$K \cap B(x_{j(n)}, 2^{-n}) \cap B(x_{j(n+1)}, 2^{-(n+1)})$$

has no finite subcover from G .

Now choose $x_n \in B(x_{j(n)}, 2^{-n}) \cap B(x_{j(n+1)}, 2^{-(n+1)})$ for each $n \geq 1$. Then x_n and x_{n+1} are both in $B(x_{j(n+1)}, 2^{-(n+1)})$, so $d(x_n, x_{n+1}) \leq 2^{-n}$. Thus $(x_n)_{n \geq 1}$ is Cauchy by the triangle inequality. Since K is complete, it follows that $x_n \rightarrow x \in K$ as $n \rightarrow \infty$. Since G covers K , it follows that $x \in G_j$ for some $j \in J$. The set G_j is open, so there is m such that $B(x, 2^{-m}) \subset G_j$. But $d(x_n, x) \leq 2^{1-n}$, so it follows that $B(x_{j(n)}, 2^{-n}) \subset G_j$ for all $n \geq m+2$, a contradiction. \square

Note that the latter part of the preceding proof is really just the proof that a closed and totally bounded set is compact (a set is totally bounded if for all $r > 0$ it can be covered with finitely many balls of radius at most r). Also, it follows immediately from the preceding proposition that any finite collection of probability measures on a Polish space is tight.

Proposition 17.6. *Fix $\mathcal{C} \subset \text{prob}(M)$ countable and list the elements of \mathcal{C} as $(P_n, n \geq 1)$. Suppose that for all $r > 0$,*

$$\lim_{m \rightarrow \infty} \liminf_{n \geq 1} P_n(\bigcup_{i=1}^m B(x_i, r)) = 1.$$

Then \mathcal{C} is tight.

Proof. For $r > 0, \epsilon > 0$ and $m \in \mathbb{N}$ define

$$F(r, m) := \bigcup_{i=1}^m B[x_i, r].$$

Then let $m_1 = m_1(r, \epsilon)$ be such that

$$\liminf_{n \geq 1} P_n(\bigcup_{i=1}^{m_1} B(x_i, r)) \geq 1 - \epsilon/2;$$

such m_1 exists by hypothesis. Then there must exist $n_0 = n_0(r, \epsilon)$ such that for all $n \geq n_0$,

$$P_n(\bigcup_{i=1}^{m_1} B(x_i, r)) \geq 1 - \epsilon. \tag{17.1}$$

Now let $m_2 = m_2(r, \epsilon)$ be such that for $1 \leq n \leq n_0$,

$$P_n(\bigcup_{i=1}^{m_2} B(x_i, r)) \geq 1 - \epsilon; \tag{17.2}$$

such m_2 exists since $(B(x_i, r), i \geq 1)$ is a cover and since $(P_i, 1 \leq i \leq n_0)$ is a finite family.

Take $m = m(r, \epsilon) = m_1 \vee m_2$; then for all $n \geq 1$, by either (17.1) or (17.2),

$$P_n(F(r, m(r, \epsilon))) \geq 1 - \epsilon.$$

Now let

$$K = \bigcap_{j \geq 1} F\left(\frac{1}{2^j}, m\left(\frac{1}{2^j}, \frac{\epsilon}{2^j}\right)\right).$$

Then $P_n(K) \geq 1 - \epsilon$ for all n , and K is closed (it is the intersection of closed sets) and is totally bounded so is compact. \square

Exercise 17.3. Let $(P_n, n \geq 1)$ and P be Borel probability measures on a metric space M such that $P_n \Rightarrow P$.

- (a) Prove that if M is separable then $(P_n, n \geq 1)$ is tight.
- (b) Prove or disprove: if P is tight then $(P_n, n \geq 1)$ is tight.

17.3. Aside: the existence of non-tight probability measures. Proposition 17.5 says that any single Borel probability measure on a Polish space is tight. It would be satisfying if I could give you an example showing that the condition that M be Polish is necessary, but that's not so easy to do. Even properly explaining why it's not easy isn't easy; it's the subject of this (optional) section, which veers into set theory and large cardinal axioms. I'm following the presentation from Fremlin (cite) here.

For sets η, ξ , write $\eta \leq \xi$ if $\eta = \xi$ or $\eta \in \xi$. An *ordinal* is a set ξ such that the following all hold.

- If $\eta \in \xi$ then η is a set and $\eta \notin \eta$.
- If $\eta \in \zeta \in \xi$ then $\eta \in \xi$.
- The partial order \leq is a well-ordering of η .

(A well-ordering of a set S is a total ordering of S such that every non-empty subset of S has a least element under the ordering.)

A *cardinal* is an initial ordinal: that is, an ordinal η such that for any $\xi \in \eta$, there is no bijection between ξ and η . The axiom of choice is equivalent to the statement that for every set S , there is a unique cardinal η such that there is a bijection between S and η ; the cardinal η is called the *cardinal* of S , or the *cardinality* of S .

For any set S we write 2^S for the power-set of S . A cardinal κ is *measure-free* if whenever μ is a probability measure on $(\kappa, 2^\kappa)$, there is $\xi < \kappa$ such that $\mu(\{\xi\}) > 0$.

This paragraph and the next are meant to give a bit of intuition for this definition. The assertion that \mathbb{N} is measure-free says that if μ is any probability measure on $(\mathbb{N}, 2^\mathbb{N})$, then there is $n \in \mathbb{N}$ such that $\mu(\{n\}) > 0$. This is true by countable additivity.

The assertion that \aleph_1 is measure-free says that if μ is any probability measure on $(\aleph_1, 2^{\aleph_1})$, then there is a (countable) ordinal $\xi \in \aleph_1$ such that $\mu(\{\xi\}) > 0$. If \mathbb{R} has cardinality \aleph_1 (the continuum hypothesis) then this in particular implies that for any probability measure μ on $(\mathbb{R}, 2^\mathbb{R})$, there is a countable set $S \subset \mathbb{R}$ such that $\mu(S) > 0$. This is not true of Lebesgue measure on \mathbb{R} , for example.

So it follows that, assuming the continuum hypothesis, if all subsets of \mathbb{R} are Lebesgue measurable then \aleph_1 is not measure-free.

Given a topological space (M, \mathcal{M}) , a *base* for \mathcal{M} is a set $B \subset \mathcal{M}$ such that every open set in \mathcal{M} is a union of sets in B . The *weight* of \mathcal{M} is the smallest cardinal of a base for \mathcal{M} . Any separable topological space has a finite or countable base, so its weight is finite or countable.

A Hausdorff space is a topological space (M, \mathcal{M}) where points are separated: for any $x, y \in M$ there are open sets U, V with $x \in U$ and $y \in V$ and $U \cap V = \emptyset$.

Theorem 17.7 (Fremlin Vol 4 page 244). *Let (M, \mathcal{M}) be a complete Hausdorff space. Then every finite Borel measure μ on (M, \mathcal{M}) is tight if and only if the weight of \mathcal{M} is measure-free.*

The assertion that there exist cardinals which are *not* measure free is believed to be independent of the axiom of choice and the continuum hypothesis (see Fremlin Vol 3 page 310). So one should not hope for an overly straightforward construction of a space on which there exist non-tight probability measures.

Billingsley (Convergence of probability measures, Second Edition, exercise 1.13), suggests the following construction. Write λ for Lebesgue measure on $[0, 1]$. For a set $S \subset [0, 1]$, the *Lebesgue outer measure* of S is defined as

$$\lambda^*(S) := \inf\{\lambda(G) : S \subset G, G \text{ open}\};$$

informally, this is the smallest total length of any collection of intervals whose union covers S . The *Lebesgue inner measure* of S is

$$\lambda_*(S) := \sup\{\lambda(K) : K \subset S, K \text{ closed}\}.$$

It is a theorem that the Lebesgue-measurable sets are precisely those for which $\lambda_*(S) = \lambda^*(S)$.

Now suppose that S is a subset of $[0, 1]$ with $\lambda^*(S) = 1$ and $\lambda_*(S) < 1$ (such sets exist assuming the axiom of choice). For x, y in S let $d(x, y) = |x - y|$; then (S, d) is a metric space. Now, for any compact K , we have $\lambda^*(K) = 0$, so λ^* can not be tight.

I don't understand this construction, because I'm not sure what measurable space λ^* is supposed to be a measure on. Presumably (S, \mathcal{F}) for some σ -field \mathcal{F} , but I'm not sure which. And, it's not clear to me how compactness relative to S relates to compactness relative to $[0, 1]$.

17.4. Bounded Lipschitz functions and $\text{prob}(M)$. Large parts of the coming material are drawn from notes written by my friend and collaborator Roberto Imbuzeiro Oliveira, notably the use of Bounded Lipschitz functions, the notion of consistent weightings, and both of their use in proving Prokhorov's theorem. Let

$$\text{BL} = \text{BL}(M) = \{f \in C_b(M) : \|f\|_{\text{BL}} < \infty\},$$

where

$$\|f\|_{\text{BL}} = \|f\|_{\infty} + \|f\|_{\text{Lip}},$$

and

$$\|f\|_{\text{Lip}} = \sup_{x, y \in M, x \neq y} \frac{|f(x) - f(y)|}{d(x, y)}.$$

Exercise 17.4. *The pair $(\text{BL}, \|\cdot\|_{\text{BL}})$ is a normed vector space; moreover, if $f, g \in \text{BL}$ then $f \wedge g$ and $f \vee g$ are both in BL .*

Exercise 17.5. *If $A \subset M$ is open then there exist non-negative BL functions $(f_n)_{n \geq 1}$ such that $f_n \uparrow \mathbf{1}_{[A]}$ as $n \rightarrow \infty$. (Consider the functions f_n defined by $f_n(x) = 1 \wedge (nd(x, A^c))$.)*

Exercise 17.6. *If $A, B \subset M$ then $\partial(A \cup B)$, $\partial(A \cap B)$ and $\partial(A \setminus B)$ are all contained in $\partial A \cup \partial B$.*

Given metric spaces $((M_i, d_i), i \geq 1)$, we define a product space $M = (M, d)$ as follows. First take $M = \prod_{i \geq 1} M_i$. Next, for $x, y \in M$, writing $x = (x_i, i \geq 1)$ and $y = (y_i, i \geq 1)$, set

$$d(x, y) = \sup_{i \geq 1} \frac{\min(1, d_i(x_i, y_i))}{2^i}.$$

base

\mathcal{B}_M

$\text{BL}, \|\cdot\|_{\text{BL}}$

$\|\cdot\|_{\text{Lip}}$

Product space

It's not hard to check the following properties.

- If the coordinate spaces $((M_i, d_i), i \geq 1)$ are all complete separable metric spaces then M is again complete and separable.
- if $(x^{(n)}, n \geq 1)$ is a sequence of elements of M , then $x^{(n)} \rightarrow x \in M$ if and only if $x_i^{(n)} \rightarrow x_i$ for all $i \geq 1$.

Due to the second property, the topology on M generated by d is sometimes called the *topology of pointwise convergence*.

Proposition 17.8. *If M is separable then the Borel σ -field \mathcal{B}_M on M is precisely the product σ -field.*

Proof. Recall that the product σ -field Π is the smallest σ -field which makes the projection maps $\pi_i : M \rightarrow M_i$ continuous. These maps are also continuous with respect to d , so the topology generated by d contains the product topology and thus \mathcal{B}_M contains Π .

Next, for any $x \in M$ and $r > 0$, we have

$$B[x, r] = \bigcap_{i \geq 1: r < 2^{-i}} \pi_i^{-1}(B_{M_i}[x_i, 2^i r]),$$

so the product σ -field on M contains the closed balls in M , so the open balls, so (since the space is separable) the Borel σ -field. \square

$Pf = \int_M f dP$
 $E_f(P) := Pf$

Given $P \in \text{prob}(M)$ and a P -integrable function $f : M \rightarrow \mathbb{R}$, write $Pf = \int_M f dP$. For bounded measurable $f : M \rightarrow \mathbb{R}$, define

$$E_f : \text{prob}(M) \rightarrow \mathbb{R}$$

$$P \mapsto Pf.$$

Write \mathcal{P}_M for the smallest σ -field which makes the maps

$$\{E_f, f : M \rightarrow \mathbb{R} \text{ bounded and measurable}\}$$

all themselves measurable. (This looks roughly like a dual space to $\text{prob}(M)$ but we haven't actually defined $\text{prob}(M)$ as a vector space.)

Proposition 17.9. \mathcal{P}_M is the smallest σ -field which makes the maps $\{E_f, f \in \text{BL}\}$ all measurable.

Proof. For the proof, write \mathcal{P}^* for the smallest σ -field which makes the maps $\{E_f, f \in \text{BL}\}$ all measurable. We must show that $\mathcal{P}^* = \mathcal{P}_M$. Functions in BL are all bounded and measurable, so it is immediate that $\mathcal{P}^* \subset \mathcal{P}_M$.

For the other inclusion, let

$$\mathcal{H} = \{f : M \rightarrow \mathbb{R} \text{ s.t. } E_f \text{ is defined and } (\mathcal{P}^*/\mathcal{B}_{\mathbb{R}})\text{-measurable}\}.$$

Then $\mathcal{H} \supseteq \text{BL}$, and to prove $\mathcal{P}^* \supset \mathcal{P}_M$ it suffices to show that \mathcal{H} contains all bounded measurable functions $f : M \rightarrow \mathbb{R}$.

Note that if $f, g \in \mathcal{H}$ then E_{cf+g} and $E_{cf+g} = cE_f + E_g$, so E_{cf+g} is $(\mathcal{P}^*/\mathcal{B}_{\mathbb{R}})$ -measurable and thus $cf + g \in \mathcal{H}$. Also, if $(f_n, n \geq 1)$ and f are in \mathcal{H} and $0 \leq f_n \uparrow f$ as $n \rightarrow \infty$, then for all $P \in \text{prob}(M)$, by the monotone convergence theorem

$$E_{f_n}P = Pf_n = \int f_n dP \rightarrow \int f dP = Pf = E_fP.$$

In other words, $E_{f_n} \rightarrow E_f$ pointwise on $\text{prob}(M)$. Measurability is preserved under pointwise limits so E_f is measurable so $f \in \mathcal{H}$. It follows that \mathcal{H} is a monotone class.

Finally, fix $A \subset M$ open. Then by Exercise 17.5, there exist functions $(f_n, n \geq 1)$ in BL (and so in \mathcal{H}) such that $0 \leq f_n \uparrow \mathbf{1}_{[A]}$, so $\mathbf{1}_{[A]} \in \mathcal{H}$. Thus $\mathcal{H} \supset \{\mathbf{1}_{[A]} : A \subset M \text{ open}\}$, so by the monotone class theorem \mathcal{H} contains all bounded measurable functions $f : M \rightarrow \mathbb{R}$, as required. \square

Exercise 17.7. Show that two measures $P, Q \in \text{prob}(M)$ are equal if and only if $Pf = Qf$ for all $f \in \text{BL}$.

Exercise 17.8. Let (S, \mathcal{S}) be a measurable space. Then a function $\Psi : S \rightarrow \text{prob}(M)$ is measurable if and only if $E_f \circ \Psi$ is $\mathcal{S}/\mathcal{B}_{\mathbb{R}}$ -measurable for all $f \in \text{BL}$.

17.5. Recursively partitioning Polish spaces. Fix any sequence $r = (r_n, n \geq 1)$ of real numbers with $r_n \downarrow 0$ as $n \rightarrow \infty$. In this section, using x and r , we define a nested sequence of measurable partitions of M , each refining the previous one.

Exercise 17.9. *The collection $(B(x_i, r_j), i, j \geq 1)$ forms a countable base for the topology of M .*

We use the notation of the Ulam-Harris tree $\mathcal{U} = \bigcup_{n \geq 0} \mathbb{N}^n$; recall that $\mathbb{N}^0 := \{\emptyset\}$ and \emptyset is the root of \mathcal{U} . We also occasionally write $\mathcal{U}_n = \mathbb{N}^n$.

For $i, n \geq 1$, let $B_{n,i} = B(x_i, r_n) \setminus (\bigcup_{j=1}^{i-1} B(x_j, r_n))$; the sets $(B_{n,i}, i \geq 1)$ partition M since $\bigcup_{i \geq 1} B(x_i, r_n) = M$. We have suppressed the dependence on x and r in the notation $B_{n,i}$ to keep things readable; if we need to make this dependence explicit we will write $B_{n,i}^{x,r}$. This dependence is also left out of the notation introduced next.

Set $A(\emptyset) = M$. For $n \geq 1$, suppose $(A(v), v \in \mathbb{N}^{n-1})$ is already defined and forms a measurable partition of M . Then for each $v \in \mathbb{N}^{n-1}$ and $i \geq 1$, let $A(vi) = A(v) \cap B_{n,i}$. Since $\bigcup_{i \geq 1} B_{n,i} = M$, it follows that $(A(v), v \in \mathbb{N}^n)$ is indeed a measurable partition refining $(A(v), v \in \mathbb{N}^{n-1})$.

Proposition 17.10. *For any $n \geq 1$, for any $K \subset\subset M$ there is $m = m(K, n)$ such that $K \subset \bigcup_{v \in \{1, \dots, m\}^n} A(v)$.*

Proof. For all $n \geq 0$, the collection $(B(x_i, r_n), i \geq 1)$ is an open cover of M , so for any $K \subset\subset M$ there is $m = m(K, n)$ such that

$$K \subset \bigcup_{i \leq m} B(x_i, r_n) = \bigcup_{i \leq m} B_{n,i}.$$

Next note that for any $m, n \geq 1$,

$$\begin{aligned} \bigcup_{v \in \{1, \dots, m\}^n} A(v) &= \bigcup_{u \in \{1, \dots, m\}^{n-1}} \bigcup_{i \leq m} A(ui) \\ &= \bigcup_{u \in \{1, \dots, m\}^{n-1}} \left(A(u) \cap \bigcup_{i \leq m} B(x_i, r_n) \right) \\ &= \bigcup_{i \leq m} B(x_i, r_n) \cap \bigcup_{u \in \{1, \dots, m\}^{n-1}} A(u), \end{aligned}$$

so by induction

$$\bigcup_{v \in \{1, \dots, m\}^n} A(v) = \bigcap_{j \leq n} \bigcup_{i \leq m} B(x_i, r_j) = \bigcup_{i \leq m} B(x_i, r_n).$$

In particular if $m = m(K, n)$ as above then $K \subset \bigcup_{v \in \{1, \dots, m\}^n} A(v)$. □

We record the following fact for later use.

Proposition 17.11. *For all $v \in \mathbb{N}^n$,*

$$\partial A(v) \subseteq \bigcup_{i \geq 1} \bigcup_{m \leq n} \partial B(x_i, r_m) = \bigcup_{i \geq 1} \bigcup_{m \leq n} \{x \in M : d(x, x_i) = r_m\}.$$

Proof. If $v = ui$ then $A(v) = A(u) \cap (B(x_i, r_n) \setminus \bigcup_{j < i} B(x_j, r_n))$, so by Exercise 17.6,

$$\partial A(v) \subset \partial A(u) \cup \bigcup_{j < i} \partial B(x_j, r_n).$$

The result follows by induction. □

In sum, from the sequence $r = (r_n, n \geq 1)$ and the countable dense set $x = \{x_i, i \geq 1\}$ we have created a countable base $(B(x_i, r_j), i, j \geq 1)$ and a collection of sets $(A(v), v \in \mathcal{U})$ such that for all $n \geq 0$, $(A(v), v \in \mathcal{U}, |v| = n)$ is a partition of \mathcal{U} . Moreover, for any $n \geq 1$, any compact $K \subset M$ is covered by a finite collection of sets from $(A(v), v \in \mathcal{U}, |v| = n)$.

Given $P \in \text{prob}(M)$, let $w_P : \mathcal{U} \rightarrow [0, 1]$ be defined by $w_P(v) = P(A(v))$. The function w_P defines a *consistent weighting* of \mathcal{U} ; this means w_P satisfies the following properties.

- w_P is non-negative and $w_P(\emptyset) = 1$.
- For all $v \in \mathcal{U}$, if $A(v) = \emptyset$ then $w_P(v) = 0$.
- For all $v \in \mathcal{U}$, $w_P(v) = \sum_{i \geq 1} w_P(v_i)$

Our next aim is to prove a theorem which provides a sort of converse to this.

The set-up for the theorem requires the axiom of countable choice. Let $\mathcal{U}^+ = \{v \in \mathcal{U} : A(v) \neq \emptyset\}$ and for $n \geq 0$ write $\mathcal{U}_n^+ = \mathcal{U}^+ \cap \mathcal{U}_n$. Then fix points $y = (y_v, v \in \mathcal{U}^+)$ such that for all $v \in \mathcal{U}^+$, we have $y_v \in A(v)$. We will use y_v as a representative of the $A(v)$.

Note that since $(A(v), v \in \mathcal{U}_n^+)$ partitions v and $d(x_i, y_v) \leq r_n$ if $v = w_i \in \mathcal{U}_n^+$, it follows that $(y_v, v \in \mathcal{U}^+)$ is again a countable dense set in M . Moreover, single points are measurable with respect to Borel σ -field, so the singleton sets $\{y_v\}$ are all measurable

Exercise 17.10. Endow the set $\mathcal{U}^* = \{\text{Functions } w : \mathcal{U} \rightarrow \mathbb{R}\} = \mathbb{R}^{\mathcal{U}}$ with the product topology (the topology of pointwise convergence). Let

$$\mathcal{W} = \{w \in \mathcal{U}^* : w \text{ is a consistent weighting}\}.$$

Then \mathcal{W} is measurable with respect to the Borel σ -field $\mathcal{B}_{\mathcal{U}^*}$.

Given $w \in \mathcal{W}$, for each n we can define a probability measure $\Psi_n(w)$ by

$$\Psi_n(w) = \sum_{v \in \mathcal{U}_n^+} w(v) \delta_{y_v};$$

this is an atomic measure assigning mass $w(v)$ to point y_v .

Proposition 17.12. Ψ_n is a measurable map from \mathcal{W} to $\text{prob}(M)$.

Proof. For any function $f \in \text{BL}$,

$$(E_f \circ \Psi_n)(w) = E_f(\Psi_n(w)) = \int f d\Psi_n(w) = \sum_{v \in \mathcal{U}_n^+} w(v) f(y_v).$$

Equivalently, $E_f \circ \Psi_n : \mathcal{W} \rightarrow \mathbb{R}$ is given by

$$E_f \circ \Psi_n \equiv \sum_{v \in \mathcal{U}_n^+} f(y_v) \pi_v,$$

where $\pi_v : \mathcal{U}^* \rightarrow \mathbb{R}$ is the projection map, $\pi_v(w) = w(v)$.

Since the projection maps are all measurable, it follows that $E_f \circ \Psi_n$ is $\mathcal{B}_{\mathcal{W}}/\mathcal{B}_{\mathbb{R}}$ -measurable for all $f \in \text{BL}$ and all $n \in \mathbb{N}$. By Exercise 17.8, it follows that Ψ_n is measurable. Alternately, argue as follows (this is really the proof of Exercise 17.8): for all $f \in \text{BL}$ and all $D \subset \mathbb{R}$ open,

$$\Psi_n^{-1}(E_f^{-1}(D)) \in \mathcal{B}_{\mathcal{W}}.$$

Since the sets $\{E_f^{-1}(D), D \subset \mathbb{R} \text{ open}, f \in \text{BL}\}$ generate \mathcal{P}_M by Proposition 17.9, it follows that Ψ_n is measurable. \square

Proposition 17.13. There exists a measurable map $\Psi = \Psi_{r,x,y} : \mathcal{W} \rightarrow \text{prob}(M)$ such that the following hold.

- (1) for all $w \in \mathcal{W}$, for all $n \geq 1$ and all $f \in \text{BL}$,

$$|\Psi_n(w)f - \Psi(w)f| \leq \|f\|_{\text{Lip}} \cdot 2r_n.$$

- (2) For all $P \in \text{prob}(M)$ we have $\Psi(w_P) = P$.

Proof. We define a sequence $(Z_n, n \geq 1)$ of random elements of M on the probability space $(\Omega, \mathcal{F}, \mathbf{P}) = ([0, 1), \mathcal{B}_{[0,1)}, \text{Leb})$ as follows.

Use \prec for the lexicographic order on \mathcal{U} . For each $n \geq 1$ and $v \in \mathcal{U}_n$, define $I_v = I_v(\mathbf{w}) = [g_v, d_v)$, where $g_v = \sum_{u \in \mathcal{U}_n, u \prec v} \mathbf{w}(u)$ and $d_v = \sum_{u \in \mathcal{U}_n, u \preceq v} \mathbf{w}(u)$, so that the intervals $(I_v, v \in \mathcal{U}_n)$ partition $[0, 1)$. The resulting partitions inherit the nested structure of the partitions of M , in that $(I_{v_j}, j \geq 1)$ forms a partition of I_v .

Let $Z_n : [0, 1) \rightarrow M$ be given by $Z_n = \sum_{v \in \mathcal{U}_n} y_v \delta_{I_v}$; i.e., if $\omega \in I_v$ then $Z_n(\omega) = y_v \in A(v)$. Then Z_n has distribution $\Psi_n(\mathbf{w})$:

$$\mathbf{P} \{Z_n = y_v\} = \mathbf{P} \{Z_n \in I_v\} = |I_v| = \mathbf{w}(v) = \Psi_n(\mathbf{w})(v).$$

Fix $n \leq m$. For all $\omega \in [0, 1)$, and let $u \in \mathcal{U}_n$ and $v \in \mathcal{U}_m$ be such that $\omega \in I_u$ and $\omega \in I_v$. then $y_u \in A(u) \subset A(v)$ and $y_v \in A(v)$, so $d(y_u, y_v) \leq \text{diam}(A(v)) \leq 2r_n$. It follows that $d(Z_m, Z_n) \leq r_n$, so $(Z_n, n \geq 1)$ is almost surely (in fact, surely) Cauchy and so convergent. Let $Z : \Omega \rightarrow M$ be the limiting random variable and let $\Psi(\mathbf{w})$ be its law; then $|Z_n - Z| \leq 2r_n$ almost surely, so for any $f \in \text{BL}$, by a change of variables,

$$|\Psi(\mathbf{w})f - \Psi_n(\mathbf{w})f| = |\mathbf{E}[f(Z)] - \mathbf{E}[f(Z_n)]| \leq \mathbf{E}[|f(Z) - f(Z_n)|] \leq 2r_n \cdot \|f\|_{\text{Lip}},$$

which is the first assertion of the proposition.

Before proving the second assertion, we verify the measurability of Ψ . Note that

$$\Psi(\mathbf{w})f = \int f d\Psi(\mathbf{w}) = E_f(\Psi(\mathbf{w})) = (E_f \circ \Psi)(\mathbf{w})$$

and likewise $\Psi_n(\mathbf{w})f = (E_f \circ \Psi_n)(\mathbf{w})$. The preceding bound then yields that $E_f \circ \Psi_n \rightarrow E_f \circ \Psi$ for all $f \in \text{BL}$. Thus $E_f \circ \Psi$ is measurable for all $f \in \text{BL}$, so Ψ is measurable by Exercise 17.8.

Now suppose that in fact $\mathbf{w} = \mathbf{w}_P$ for some probability measure $P \in \text{prob}(M)$. By Exercise 17.7, to show $\Psi(\mathbf{w}_P) = P$ it suffices to show that $Pf = \Psi(\mathbf{w}_P)f$ for all $f \in \text{BL}$. So fix $f \in \text{BL}$, and for $v \in \mathcal{U}$ define

$$F(v) = \begin{cases} \frac{\int_{A(v)} f dP}{P(A(v))} & \text{if } P(A(v)) > 0 \\ 0 & \text{otherwise.} \end{cases}$$

Think of this as the conditional expectation of Y given that $Y \in A(v)$, where Y is a random variable with law P .

We have

$$\Psi_n(\mathbf{w}_P)f = \sum_{v \in \mathcal{U}_n} P(A(v))f(y_v) = \sum_{v \in \mathcal{U}_n: P(A(v)) > 0} P(A(v))f(y_v),$$

so

$$|\Psi_n(\mathbf{w}_P)f - Pf| = \left| \sum_{v \in \mathcal{U}_n: P(A(v)) > 0} P(A(v))(f(y_v) - F(v)) \right|. \quad (17.3)$$

For $v \in \mathcal{U}_n$, if $P(A(v)) > 0$ then since $A(v)$ has diameter at most $2r_n$,

$$\inf(f(y) : y \in A(v)) \leq y_v \leq \sup(f(y) : y \in A(v)) \leq 2r_n \cdot \|f\|_{\text{Lip}};$$

since $F(v)$ is an average of f over $A(v)$, it is also bounded between $\inf(f(y) : y \in A(v))$ and $\sup(f(y) : y \in A(v))$, so (17.3) gives

$$|\Psi_n(\mathbf{w}_P)f - Pf| \leq \sum_{v \in \mathcal{U}_n: P(A(v)) > 0} P(A(v)) \cdot 2r_n \cdot \|f\|_{\text{Lip}} = 2r_n \cdot \|f\|_{\text{Lip}}.$$

Together with the first assertion of the proposition, it follows that $Pf = \Psi(\mathbf{w}_P)f$, as required. \square

17.6. **Metriizing weak convergence.** Metrize $\text{prob}(M)$ via the dual norm

$$d_{BL}^*(P, Q) = \|P - Q\|_{BL,*} := \sup_{f \in BL, \|f\|_{BL} \neq 0} \frac{|Pf - Qf|}{\|f\|_{BL}}.$$

Theorem 17.14. *Assume M is Polish. Fix probability measures $(P_n, n \geq 1)$ and P from $\text{prob}(M)$. Then $P_n \Rightarrow P$ if and only if $\|P_n - P\|_{BL,*} \rightarrow 0$.*

This theorem may seem stronger than we have a right to expect. Weak convergence means $P_n f \rightarrow P f$ for all $f \in BL$. Theorem 17.14 implies that when this occurs, in fact the convergence is uniform over functions with $\|f\|_{BL} \leq 1$.

Lemma 17.15. *Fix a countable dense set $x = \{x_i, i \geq 1\}$ and a sequence $r = (r_n, n \geq 1)$ with $r_n \downarrow 0$, and let $\Psi : \mathcal{W} \rightarrow \text{prob}(M)$ be as in Proposition 17.13. Next fix $P, Q \in \text{prob}(M)$, and consistent weightings p, q of \mathcal{U} such that $\Psi(p) = P$ and $\Psi(q) = Q$. Then*

$$\|P - Q\|_{BL,*} \leq \inf_{n \geq 1} \left(4r_n + \sum_{v \in \mathcal{U}_n} |p(v) - q(v)| \right).$$

Proof. Fix a function $f \in BL$ which is not identically zero, and $n \geq 1$. We must show that for all $n \geq 1$,

$$|Pf - Qf| \leq \|f\|_{BL} \cdot \left(4r_n + \sum_{v \in \mathcal{U}_n} |p(v) - q(v)| \right).$$

We use the triangle inequality to write

$$|Pf - Qf| = |\Psi(p)f - \Psi(q)f| \leq |\Psi_n(p)f - \Psi_n(q)f| + |\Psi(p)f - \Psi_n(p)f| + |\Psi(q)f - \Psi_n(q)f|.$$

By Proposition 17.13, we know that $|\Psi(p)f - \Psi_n(p)f| \leq \|f\|_{Lip} \cdot 2r_n$; the same bound holds for $|\Psi(q)f - \Psi_n(q)f|$. We thus have

$$\begin{aligned} |Pf - Qf| &\leq |\Psi_n(p)f - \Psi_n(q)f| + 4r_n \|f\|_{Lip} \\ &= \left| \sum_{v \in \mathcal{U}_n} (p(v) - q(v)) f(y_v) \right| + 4r_n \|f\|_{Lip} \\ &\leq \sum_{v \in \mathcal{U}_n} |p(v) - q(v)| \|f(y_v)\| + 4r_n \|f\|_{Lip} \\ &\leq \|f\|_{\infty} \sum_{v \in \mathcal{U}_n} |p(v) - q(v)| + 4r_n \|f\|_{Lip}. \end{aligned}$$

The bound follows since both $\|f\|_{Lip}$ and $\|f\|_{\infty}$ are at most $\|f\|_{BL}$. □

Lemma 17.16. *Fix a sequence $(Q_k, k \geq 1)$ of elements of $\text{prob}(M)$. Suppose that there exists a function $p : \mathcal{U} \rightarrow [0, 1]$ such that for all $n \geq 1$,*

$$\lim_{k \rightarrow \infty} \sum_{v \in \mathcal{U}_n} |p(v) - Q_k(A(v))| = 0.$$

Then $p \in \mathcal{W}$ and $\|Q_k - \Psi(p)\|_{BL,} \rightarrow 0$ as $k \rightarrow \infty$.*

Proof. For all $v \in \mathcal{U}$,

$$|p(v) - Q_k(A(v))| \rightarrow 0;$$

since $Q_k(A(v))$ is non-negative it follows that $p(v) \geq 0$; since $Q_k(\emptyset) = 1$ for all k it follows that $p(\emptyset) = 1$. Also, if $A(v) = \emptyset$ then $Q_k(A(v)) = 0$ for all k so $p(v) = 0$. Moreover,

$$\begin{aligned} \left| p(v) - \sum_{i \geq 1} p(vi) \right| &= \left| p(v) - Q_k(A(v)) - \sum_{i \geq 1} (p(vi) - Q_k(A(vi))) \right| \\ &\leq |p(v) - Q_k(A(v))| + \sum_{i \geq 1} |p(vi) - Q_k(A(vi))|. \end{aligned}$$

The final bound tends to zero as $k \rightarrow \infty$ by assumption, so $p(v) = \sum_{i \geq 1} p(vi)$; thus p is a consistent weighting, so $P := \Psi(p)$ is defined. Now note that for all $n \geq 1$, by Lemma 17.15 and the assumptions of the Lemma,

$$\limsup_{k \rightarrow \infty} \|Q_k - P\|_{\text{BL},*} \leq \limsup_{k \rightarrow \infty} \left(4r_n + \sum_{v \in \mathcal{U}_n} |p(v) - Q_k(A(v))| \right) = 4r_n.$$

Taking $n \rightarrow \infty$, the result follows. \square

Lemma 17.17. *Fix a tight sequence $(Q_k, k \geq 1)$ from $\text{prob}(M)$, and suppose that there is $p : \mathcal{U} \rightarrow [0, 1]$ such that $Q_k(A(v)) \rightarrow p(v)$ for all $v \in \mathcal{U}$. Then $p \in \mathcal{W}$ and $\|Q_k - \Psi(p)\|_{\text{BL},*} \rightarrow 0$ as $k \rightarrow \infty$.*

Proof. Fix $\epsilon > 0$ and let $K \subset\subset M$ be such that $\inf_{k \geq 1} Q_k(K) > 1 - \epsilon$; such K exists since $(Q_k, k \geq 1)$ is tight. Now fix $n \geq 1$ and let m be such that $K \subset \cup_{v \in [m]^n} A(v)$; such m exists by Proposition 17.10. Then

$$\sum_{v \in \mathcal{U}_n} |p(v) - Q_k(A(v))| = \sum_{v \in [m]^n} |p(v) - Q_k(A(v))| + \sum_{v \in \mathcal{U}_n \setminus [m]^n} |p(v) - Q_k(A(v))|$$

The first sum has a finite number of terms so the pointwise convergence of p to Q_k implies that

$$\lim_{k \rightarrow \infty} \sum_{v \in [m]^n} |p(v) - Q_k(A(v))| = 0.$$

To handle the second, notice that since $Q_k(A(v)) \rightarrow p(v)$ for all v , by Fatou's lemma we have

$$\sum_{v \in \mathcal{U}_n \setminus [m]^n} p(v) = \sum_{v \in \mathcal{U}_n \setminus [m]^n} \liminf_{k \rightarrow \infty} Q_k(A(v)) \leq \liminf_{k \rightarrow \infty} \sum_{v \in \mathcal{U}_n \setminus [m]^n} Q_k(A(v)) < \epsilon,$$

the last bound holding since the sets $(A(v), v \in \mathcal{U}_n \setminus [m]^n)$ are disjoint and their union is contained in K . It follows that

$$\sum_{v \in \mathcal{U}_n \setminus [m]^n} |p(v) - Q_k(A(v))| \leq \sum_{v \in \mathcal{U}_n \setminus [m]^n} p(v) + \sum_{v \in \mathcal{U}_n \setminus [m]^n} Q_k(A(v)) < 2\epsilon.$$

We thus have

$$\limsup_{k \rightarrow \infty} \sum_{v \in \mathcal{U}_n} |p(v) - Q_k(A(v))| < 2\epsilon;$$

since $\epsilon > 0$ was arbitrary the result follows by Lemma 17.16. \square

Proof of Theorem 17.14. Suppose $P_n \Rightarrow P$. We assume the radii $r = (r_n, n \geq 1)$ are chosen so that $P(\partial A(v)) = 0$ for all $v \in V$; this is possible by Proposition 17.11. By the Portmanteau theorem, it follows that $P_n(A(v)) \rightarrow P(A(v))$ for all $v \in \mathcal{U}$. Also, the fact that $P_n \Rightarrow P$ implies that $(P_n, n \geq 1)$ is tight (see Exercise 17.3 (b)), and Lemma 17.17 then implies that $\|P_n - P\|_{\text{BL},*} \rightarrow 0$ as $n \rightarrow \infty$.

Next suppose that $\|P_n - P\|_{\text{BL},*} \rightarrow 0$. Fix $G \subset M$ open, and let $(f_k, k \geq 1)$ be BL functions with $\|f_k\|_{\text{BL}} \neq 0$ and with $0 \leq f_k \uparrow \mathbf{1}_{[G]}$ as $k \rightarrow \infty$. Then for all $k \geq 1$ we have $|P_n f_k - P f_k| \rightarrow 0$ as $n \rightarrow \infty$, so

$$\liminf_{n \rightarrow \infty} P_n \mathbf{1}_{[G]} \geq \liminf_{n \rightarrow \infty} P_n f_k = P f_k.$$

Since $f_k \uparrow \mathbf{1}_{[G]}$ as $k \rightarrow \infty$, by the monotone convergence theorem it follows that

$$\liminf_{n \rightarrow \infty} P_n \mathbf{1}_{[G]} \geq \sup_{k \geq 1} P f_k = \lim_{k \rightarrow \infty} P f_k = P \mathbf{1}_{[G]}.$$

The Portmanteau theorem now implies that $P_n \Rightarrow P$. □

Theorem 17.18 (Prokhorov's theorem). *If $(P_n, n \geq 1)$ is a tight family of measures from $\text{prob}(M)$, then there exists $P \in \text{prob}(M)$ and an increasing sequence $(n_k, k \geq 1)$ such that $P_{n_k} \Rightarrow P$ as $k \rightarrow \infty$.*

Proof. For each $u \in \mathcal{U}$, the sequence $(P_n(A(u)), n \geq 1)$ is bounded so contains a convergent subsequence. Since \mathcal{U} is countable, a diagonal argument implies that there exists an increasing sequence $(n_k, k \geq 1)$ and a function $p : \mathcal{U} \rightarrow \mathbb{R}$ such that $P_{n_k} \rightarrow p$ pointwise as $k \rightarrow \infty$. Lemma 17.17 now implies that p is a consistent weighting and that $\|P_{n_k} - \Psi(p)\|_{\text{BL},*} \rightarrow 0$ as $k \rightarrow \infty$, and Theorem 17.14 then yields that $P_{n_k} \Rightarrow \Psi(p)$ as $k \rightarrow \infty$. □

17.7. Probability kernels and conditional probabilities. For a measurable space $S = (\mathcal{S}, \mathcal{G})$ we write $\text{prob}(S)$ for the collection of probability measures on S . Given measurable spaces (Ω, \mathcal{F}) and $(\mathcal{S}, \mathcal{G}) =: S$, a *transition kernel* from (Ω, \mathcal{F}) to $(\mathcal{S}, \mathcal{G})$ is a function

$$K : \Omega \rightarrow \text{prob}(S) \\ \omega \mapsto K^\omega$$

such that for all $B \in \mathcal{G}$, the function $\omega \mapsto K^\omega(B)$ is $(\mathcal{F}/\mathcal{B}_{\mathbb{R}})$ -measurable.

Aside. In the case that $(\mathcal{S}, \mathcal{G}) = (M, \mathcal{B}_M)$ for $M = (M, \mathcal{B}_M)$ a metric space, this is equivalent to saying that for all $f \in \text{BL}(M)$,

$$E_f \circ K : \Omega \rightarrow \mathbb{R} \\ \omega \mapsto K^\omega f$$

is $(\mathcal{G}/\mathcal{B}_{\mathbb{R}})$ -measurable, or in other words that K is $(\mathcal{F}/\mathcal{P}_M)$ -measurable; see Exercise 17.8.

Theorem 17.19. *Given measurable spaces (Ω, \mathcal{F}) and $(\mathcal{S}, \mathcal{G}) =: S$ and a transition kernel $K : \Omega \rightarrow \text{prob}(S)$, for any probability measure \mathbf{P} on (Ω, \mathcal{F}) there is a unique measure $\mathbf{P} \otimes K$ on $(\Omega \times \mathcal{S}, \mathcal{F} \otimes \mathcal{G})$ such that for all $A \in \mathcal{F}$ and $B \in \mathcal{G}$,*

$$(\mathbf{P} \otimes K)(A \times B) = \int_A K^\omega(B) \mathbf{P}(d\omega).$$

The measure $\mathbf{P} \otimes K$ has the following property. If $f : \Omega \times \mathcal{S} \rightarrow \mathbb{R}$ is $(\mathcal{F} \times \mathcal{G}/\mathcal{B}_{\mathbb{R}})$ -measurable and either f is non-negative or $f \in L_1(\Omega \times \mathcal{S}, \mathcal{F} \otimes \mathcal{G}, \mathbf{P} \otimes K)$ then

$$\omega \mapsto \int_S f(\omega, s) K^\omega(ds)$$

is $(\mathcal{F}/\mathcal{B}_{\mathbb{R}})$ -measurable and

$$(\mathbf{P} \otimes K)f = \int_\Omega \int_S f(\omega, s) K^\omega(ds) \mathbf{P}(d\omega).$$

Proof. First note that by assumption, $K^\omega(B)$ is \mathcal{F} -measurable as a function of ω , and is bounded, so we may write

$$\int_A K^\omega(B) \mathbf{P}(d\omega) = \int_A \int_B \mathbf{1} K^\omega(ds) \mathbf{P}(d\omega) = \int_\Omega \int_S \mathbf{1}_{[A \times B]}(\omega, s) K^\omega(ds) \mathbf{P}(d\omega)$$

□

Theorem 17.20. *Fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and a sub- σ -field \mathcal{G} of \mathcal{F} . Then for any Polish space $M = (M, d)$ and any $(\mathcal{F}/\mathcal{B}_M)$ -measurable random variable $Y : \Omega \rightarrow M$, there exists a probability kernel*

$$P : \Omega \rightarrow \text{prob}(M) \\ \omega \mapsto P^\omega$$

Transition kernel

such that the following holds. If $f : M \rightarrow \mathbb{R}$ is $(\mathcal{B}_M/\mathcal{B}_{\mathbb{R}})$ -measurable and $\mathbf{E}|f(Y)| < \infty$, then for almost all $\omega \in \Omega$,

$$P^\omega |f| < \infty \quad \text{and} \quad \mathbf{E} \{f(Y) \mid \mathcal{G}\}(\omega) = P^\omega f.$$

Conditional distribution

The last condition means it makes sense to view P^ω as the conditional distribution of Y given \mathcal{G} , and we write

$$\begin{aligned} \mathbf{P} \{Y \in A \mid \mathcal{G}\} : \Omega &\rightarrow [0, 1] \\ \omega &\mapsto P^\omega(A). \end{aligned}$$

We proceed much as in the development for Prokhorov's theorem, using approximation by atomic measures and taking a suitable limit. For $\omega \in \Omega$ and $u \in \mathcal{U}$, let

$$p(u, \omega) = \mathbf{E} \{ \mathbf{1}_{[Y \in A(u)]} \mid \mathcal{G} \}(\omega).$$

Sometimes we hold u fixed and let ω vary; we write $p_u(\omega) = p(u, \omega)$ to emphasize that this is taking place. Likewise, we sometimes hold ω fixed and think of u as varying, and in these cases write $p^\omega(u) = p(u, \omega)$.

Note that $p_u : \Omega \rightarrow \mathbb{R}$ is $(\mathcal{G}/\mathcal{B}_{\mathbb{R}})$ -measurable for each $u \in \mathcal{U}$. Since $\mathcal{B}_{\mathbb{R}^{\mathcal{U}}}$ is generated by the projection maps onto single coordinates, it follows that the map

$$\omega \mapsto (p^\omega(u), u \in \mathcal{U}) \tag{17.4}$$

from Ω to $\mathbb{R}^{\mathcal{U}}$ is $(\mathcal{G}/\mathcal{B}_{\mathbb{R}^{\mathcal{U}}})$ -measurable.

Claim 17.21. *It holds that $p^\omega \in \mathcal{W}$ for \mathbf{P} -almost all $\omega \in \Omega$.*

Proof. The idea of the proof is this: for a function $q : \mathcal{U} \rightarrow \mathbb{R}$, the property that $q \in \mathcal{U}$ may be written as the intersection of countably many ‘‘local’’ properties relating to the value of q at a node and perhaps at its children. A countable intersection of probability-one events still has probability one, so it suffices to check that each of these local properties holds almost surely.

More precisely, to prove the claim, it suffices to observe that each of the following holds with \mathbf{P} -probability one.

- $p^\omega(\emptyset) = \mathbf{E} \{ \mathbf{1}_{[Y \in M]} \mid \mathcal{G} \}(\omega) = \mathbf{E} \{ 1 \mid \mathcal{G} \}(\omega) = 1$. This holds with \mathbf{P} -probability one since 1 is a version of $\mathbf{E} \{ 1 \mid \mathcal{G} \}$.
- $p^\omega(u) = \mathbf{E} \{ \mathbf{1}_{[Y \in A(u)]} \mid \mathcal{G} \}(\omega)$ is non-negative for all $u \in \mathcal{U}$, with equality if $A(u) = \emptyset$. This holds with \mathbf{P} -probability one by monotonicity of conditional probability and since 0 is a version of $\mathbf{E} \{ \mathbf{1}_{[Y \in \emptyset]} \mid \mathcal{G} \}$.
- For all $u \in \mathcal{U}$,

$$\begin{aligned} \sum_{i \geq 1} p^\omega(ui) &= \sum_{i \geq 1} \mathbf{E} \{ \mathbf{1}_{[Y \in A(ui)]} \mid \mathcal{G} \}(\omega) \\ &= \mathbf{E} \left\{ \sum_{i \geq 1} \mathbf{1}_{[Y \in A(ui)]} \mid \mathcal{G} \right\}(\omega) \\ &= \mathbf{E} \{ \mathbf{1}_{[Y \in A(u)]} \mid \mathcal{G} \}(\omega) \\ &= p^\omega(u). \end{aligned}$$

This holds with \mathbf{P} -probability one by linearity of conditional expectation. \square

Proof of Theorem 17.20. Fix $Y : \Omega \rightarrow M$ and $f : M \rightarrow \mathbb{R}$ with $\mathbf{E}|f(Y)| < \infty$ as in the statement of the theorem. We first assume $f \in \text{BL}$.

Let $P^\omega = \Psi(p^\omega)$ for $\omega \in \Omega$. As Ψ is $(\mathcal{B}_{\mathcal{W}}/\mathcal{P}_M)$ -measurable, by (17.4) it follows that P^ω is $(\mathcal{G}/\mathcal{P}_M)$ -measurable. Since $f \in \text{BL}$, by Proposition 17.13 (a) we have

$$P^\omega f = \Psi(p^\omega) f = \lim_{n \rightarrow \infty} \Psi_n(p^\omega) f. \tag{17.5}$$

Since

$$p^\omega(u) = p_u(\omega) = \mathbf{E} \{ \mathbf{1}_{[Y \in A(u)]} \mid \mathcal{G} \} (\omega),$$

we have

$$\Psi_n(p^\omega) = \sum_{u \in \mathcal{U}_n} \mathbf{E} \{ \mathbf{1}_{[Y \in A(u)]} \mid \mathcal{G} \} (\omega) \cdot \delta_{y_u},$$

so

$$\begin{aligned} \lim_{n \rightarrow \infty} \Psi_n(p^\omega) f &= \lim_{n \rightarrow \infty} \sum_{u \in \mathcal{U}_n} \mathbf{E} \{ \mathbf{1}_{[Y \in A(u)]} \mid \mathcal{G} \} (\omega) \cdot f(y_u) \\ &= \lim_{n \rightarrow \infty} \mathbf{E} \left\{ \sum_{u \in \mathcal{U}_n} \mathbf{1}_{[Y \in A(u)]} \cdot f(y_u) \mid \mathcal{G} \right\} (\omega), \end{aligned} \quad (17.6)$$

the last equality holding for all ω outside of a set of measure 0.

For $u \in \mathcal{U}_n$, if $Y \in A(u)$ then $d(Y, y_u) \leq 2r_n$, so $|f(y_u) - f(Y)| \leq 2r_n \|f\|_{\text{Lip}}$. It follows that

$$\left| f(Y) - \sum_{u \in \mathcal{U}_n} \mathbf{1}_{[Y \in A(u)]} \cdot f(y_u) \right| \leq 2r_n \|f\|_{\text{Lip}}.$$

In particular $\sum_{u \in \mathcal{U}_n} \mathbf{1}_{[Y \in A(u)]} \cdot f(y_u) \rightarrow f(Y)$ and

$$\left| \sum_{u \in \mathcal{U}_n} \mathbf{1}_{[Y \in A(u)]} \cdot f(y_u) \right| \leq |f(Y)| + 2r_1 \|f\|_\infty,$$

so by the conditional dominated convergence theorem it follows that

$$\lim_{n \rightarrow \infty} \mathbf{E} \left\{ \sum_{u \in \mathcal{U}_n} \mathbf{1}_{[Y \in A(u)]} \cdot f(y_u) \mid \mathcal{G} \right\} \stackrel{\text{a.s.}}{=} \mathbf{E} \{ f(Y) \mid \mathcal{G} \}. \quad (17.7)$$

Combining (17.5), (17.6) and (17.7) we see that $\mathbf{E} \{ f(Y) \mid \mathcal{G} \} (\omega) = P^\omega f$ for \mathbf{P} -almost all ω when $f \in \text{BL}$.

Next, if $f = \mathbf{1}_{[G]}$ for $G \subset M$ open, then we may take $0 \leq f_k \uparrow f$ as $k \rightarrow \infty$, with $f_k \in \text{BL}$ for all k . Since $P^\omega \in \text{prob}(M)$, it follows that $P^\omega f = \lim_{k \rightarrow \infty} P^\omega f_k$. The conditional monotone convergence theorem gives that

$$\mathbf{E} \{ f(Y) \mid \mathcal{G} \} \stackrel{\text{a.s.}}{=} \lim_{k \rightarrow \infty} \mathbf{E} \{ f_k(Y) \mid \mathcal{G} \},$$

so \mathbf{P} -almost surely $P^\omega f = \mathbf{E} \{ f(Y) \mid \mathcal{G} \} (\omega)$ for f the indicator of an open set. The same argument shows that the set

$$\text{Good} := \{ \text{Measurable } f : M \rightarrow \mathbb{R}, \text{ f.s.t. } \mathbf{P}\text{-almost-surely } P^\omega f = \mathbf{E} \{ f(Y) \mid \mathcal{G} \} (\omega) \}$$

is closed under non-negative monotone limits; the closure of this set under affine combinations follows from linearity of integration/expectation. The monotone class theorem then implies that Good contains all bounded measurable $f : M \rightarrow \mathbb{R}$. Applying monotonicity again shows Good contains all non-negative functions, Finally, if $\mathbf{E}|f(Y)| < \infty$ then writing $f = f^+ - f^-$ where f^+ and f^- are the positive and negative parts of f , using linearity of integration completes the proof. \square

It's worthwhile to briefly discuss the case when $\mathcal{G} = \sigma(X)$, where $X : \Omega \rightarrow N$ takes values in another Polish space $\mathbf{N} = (N, d_N)$. In this case for all $u \in \mathcal{U}$,

$$p^\omega(u) = \mathbf{E} \{ \mathbf{1}_{[Y \in A(u)]} \mid \mathcal{G} \} (\omega)$$

is \mathbf{P} -almost surely constant on fibres $\{\omega : X(\omega) = x\}$ for all $x \in \mathbb{R}$. Since \mathcal{U} is countable it follows that p^ω is almost surely constant on fibres, so $P^\omega = \Psi(p^\omega)$ is as well. Thus P^ω factors through N : there is a $(\mathcal{B}_N/\mathcal{B}_M)$ -measurable function $Q : N \rightarrow \text{prob}(M)$ such that \mathbf{P} -almost surely

$$P^\omega = Q^{X(\omega)}. \quad (17.8)$$

It follows that Q^x is a probability kernel from N to M . In this case it's natural to use the notation

$$\mathbf{P} \{Y \in B \mid X = x\} := Q^x(B);$$

for all $A \in \mathcal{B}_N$ and $B \in \mathcal{B}_M$, we then have

$$\begin{aligned} \mathbf{P} \{X \in A, Y \in B\} &= \mathbf{E} [\mathbf{1}_{[X \in A]} \mathbf{1}_{[Y \in B]}] \\ &= \mathbf{E} [\mathbf{1}_{[X \in A]} \mathbf{E} \{ \mathbf{1}_{[Y \in B]} \mid X \}] \\ &= \mathbf{E} [\mathbf{1}_{[X \in A]} Q_X(B)] \\ &= \int_{\Omega} \mathbf{1}_{[A]}(\omega) Q_{X(\omega)}(B) \mathbf{P}(d\omega). \end{aligned}$$

The second equality holds by the tower law since $\mathbf{1}_{[X \in A]}$ is X -measurable; the third holds by (17.8) and the definition of P^ω . The change of variables formula then yields

$$\mathbf{P} \{X \in A, Y \in B\} = \int_A \mathbf{P} \{Y \in B \mid X = x\} \mathcal{L}_X(dx).$$

In other words, the joint distribution of (X, Y) is given by $\mathcal{L}_X \otimes Q$ as in Theorem 17.19.

List of notation and terminology

$\mathcal{A}(\mathbb{R})$	$\mathcal{A}(\mathbb{R}) = \{\cup_{i=1}^n (a_i, b_i] : n \geq 1, -\infty < a_1 \leq b_1 \leq a_2 \leq \dots \leq a_n \leq b_n < \infty\}$; finite unions of half-open intervals.	6
$\mathcal{B}(\mathbb{R})$	The Borel σ -field of M ; smallest σ -field containing all open sets of M	12
$\mathcal{B}(\mathbb{R})$	The Borel σ -field of \mathbb{R} ; equals $\sigma(\mathcal{A}(\mathbb{R}))$	12
CDF	Cumulative distribution function: A Stieltjes function with F with $\lim_{x \rightarrow -\infty} F(x) = 0$ and $\lim_{x \rightarrow \infty} F(x) = 1$	11
$d\mu/d\nu$	The Radon-Nikodým derivative of μ with respect to ν	80
$\mathbf{E}_{\mathbf{Q}} \{X\}$	“Expectation” notation meaning $\int X d\mathbf{Q}$; \mathbf{Q} need not be a probability measure.	80
Field	A collection of subsets of a ground set, closed under finite union and complement.	5
λ -system	A set $\mathcal{A} \subset 2^\Omega$ with $\Omega \in \mathcal{A}$, closed under monotone limits and relative complements.	9
\mathcal{L}_X	The distribution of random variable X	3
Measure	A countably additive function $\mu : \mathcal{F} \rightarrow [0, \infty]$ on a σ -field, with $\mu(\emptyset) = 0$	5
G_X	The moment generating function of X , $G_X(s) = \mathbf{E} [e^{-sX}]$	42
G_X	The moment generating function of X , $G_X(s) = \mathbf{E} [e^{sX}]$	91
μ_X	The distribution of random variable X	19
Outer measure	A monotone, countably subadditive function $\mu : 2^\Omega \rightarrow [0, \infty]$ with domain the power set of some set Ω , such that $\mu(\emptyset) = 0$	7
φ_X	The characteristic function of X , $\varphi_X = \mathbf{E} [e^{itX}]$	97
π -system	A collection of subsets of a ground set, closed under finite intersection.	5
Pre-measure	A σ -additive function $\mu : \mathcal{A} \rightarrow [0, \infty]$ on a ring \mathcal{A} with $\mu(\emptyset) = 0$	6
Pre-measure space	A triple $(\Omega, \mathcal{A}, \mu)$ where \mathcal{A} is a ring over Ω and μ is a pre-measure on \mathcal{A}	6
Probability space	A measure space $(\Omega, \mathcal{F}, \mathbf{P})$ with $\mathbf{P}(\Omega) = 1$	13
Ring	A collection of subsets of a ground set, closed under finite union and relative complement.	5
S^1	The unit circle in \mathbb{C}	96
$\sigma(\mathcal{A})$	The smallest σ -field containing \mathcal{A}	5
σ -field	A collection of subsets of a ground set, closed under countable union and complement.	5
Stieltjes Function	A non-decreasing function $F : \mathbb{R} \rightarrow \mathbb{R}$ which is right continuous with left limits.	11

\mathcal{T}	The set of subtrees of \mathcal{U}	84
\mathcal{U}	The Ulam-Harris tree: rooted; each node has countably many offspring. .	83

References

- [1] C. C. Heyde. On a property of the lognormal distribution. *J. Roy. Statist. Soc. Ser. B*, 25:392–393, 1963. ISSN 0035-9246. URL <https://mathscinet.ams.org/mathscinet-getitem?mr=171336>. 105
- [2] Karol A Penson, Pawel Blasiak, Andrzej Horzela, and Allan I Solomon. On certain non-unique solutions of the Stieltjes moment problem. *Discrete Mathematics and Theoretical Computer Science*, 12(1):295–306, 2010. 95
- [3] J. A. Shohat and J. D. Tamarkin. *The Problem of Moments*. American Mathematical Society Mathematical surveys, vol. I. American Mathematical Society, New York, 1943. URL <https://mathscinet.ams.org/mathscinet-getitem?mr=0008438>. 105
- [4] Elias M. Stein and Rami. Shakarchi. *Fourier analysis : an introduction*. Princeton University Press, 2003. ISBN 9780691113845. URL <https://press.princeton.edu/titles/7562.html>. 96

Department of Mathematics and Statistics, McGill University, Montréal, Canada.

Email address: louigi.addario@mcgill.ca

URL: <http://problab.ca/louigi/>