Integration

surable function:

Measurable sets and functions

Definition 1 (\sigma-algebras) Let Ω be a set and $\mathcal{A} \subseteq \mathcal{P}(\Omega)$ be a collection of subsets of Ω : i. A is an algebra if $\varnothing \in A$ and for $A, B \in A$, $A^c = \Omega \setminus A \in A$ and $A \cup B \in A$. ii. A is a σ -algebra if $\varnothing \in A$, for $A \in A$, $A^c \in A$, and for (A_n) a sequence of sets in A, $\bigcup_{n=1}^{\infty} A_n \in \mathcal{A}.$

A collection of sets is an algebra subject to being closed under finite applications of the basic operators. The σ -algebra concept extends this slightly to infinite ones. Consider where this distinction is relevant?

Note that if we have $\{\mathcal{F}_i : i \in I\}$ are σ -algebras, then

is a σ -algebra. This allows us to consider the notion of a smallest σ -algebra containing a set (the σ -algebra 'generated' by a set). We write the σ -algebra generated by a collection of collections of sets \mathfrak{A} as $\sigma(\mathfrak{A})$.

Definition 2 (Borel σ -algebra) Let (E, \mathcal{T}) be a topological space. The σ -algebra generated by the open sets in E is called the Borel σ -algebra on E and is denoted $\mathcal{B}(E) = \sigma(\mathcal{T})$.

Definition 3 Suppose $(\Omega_i, \mathcal{F}_i)_{i \in I}$ are measurable spaces. With $\Omega = \prod_{i \in I} \Omega_i$, \mathcal{F} the σ -algebra generated by $A = \prod_{i \in I} A_i$ where $A_i \in \mathcal{F}_i$ for all $i \in I$ and for all but finitely many $i \in I$, $A_i = \Omega_i$: (Ω, \mathcal{F}) is the product space.

This space is measurable, and \mathcal{F} is a σ -algebra.

Definition 4 (\pi and \lambda-systems) A collection of sets A is called a π -system if it is closed under intersections.

A collection of sets \mathcal{M} is called a λ -system if $\Omega \in \mathcal{M}$, if $A, B \in \mathcal{M}$, $A \subseteq B$, then $B \setminus A \in \mathcal{M}$, and if $(A_n) \subseteq \mathcal{M}$ with $A_n \subseteq A_{n+1}$ increasing then $\bigcup_{n>1} A_n \in \mathcal{M}$.

A collection of sets is a σ -algebra if and only if it is both a π -system and a λ -system.

Lemma 1 (π - λ systems lemma) Let \mathcal{A} be a π -system and \mathcal{M} a λ -system. Then if $\mathcal{A} \subseteq \mathcal{M}$ then $\sigma(\mathcal{A}) \subseteq \mathcal{M}$.

We can use this with a convenient π -system to show that our λ -system contains more than is immediately obvious.

Let $\lambda(\mathcal{A})$ be the smallest λ -system containing \mathcal{A} . This is a subset of \mathcal{M} and $\sigma(\mathcal{A})$, so we just need to show that $\lambda(\mathcal{A})$ is a σ -algebra (for which we just have to show that it is a π -system).

Definition 5 (Random variables) With measurable spaces (Ω, \mathcal{F}) , (E, \mathcal{E}) , a function $f: \Omega \to \mathcal{F}$ E is said to be an E-valued random variable (or a measurable function) if for all $A \in \mathcal{E}$, $f^{-1}(A) \in \mathcal{F}$.

We get immediately that random variables can be composed as one would expect. We can also use random variables to define new σ -algebras. Note that $(\Omega, \{f^{-1}(A) : A \in \mathcal{E}\})$ is a σ -algebra.

 σ -algebra on Ω for which all f_i are measurable.

Definition 6 With $\{f_i: i \in I\}$ a family of functions $\Omega \to E$, $\sigma(f_i: i \in I)$ is the smallest

This is initially a slightly intimidating definition, but the intuition is just that we need our $\sigma(f_i: i \in I) = \sigma(f_i^{-1}(A): A \in \mathcal{E}, i \in I).$

Theorem 2 (Monotone Class Theorem) Let \mathcal{H} be a class of bounded functions from $\Omega \to \mathbb{R}$ such that

• the constant function $1 \in \mathcal{H}$,

• \mathcal{H} is a vector space over \mathbb{R} ,

• if $(f_n) \subseteq \mathcal{H}$, $f_n \to f$ monotonically increasing, then $f \in \mathcal{H}$,

then if $\mathcal{C} \subseteq \mathcal{H}$, and \mathcal{C} is closed under pointwise multiplication, then all bounded $\sigma(\mathcal{C})$ measurable functions are in \mathcal{H} .

To get an intuition for this, note that any $f \in \mathcal{C}$ is necessarily bounded and $\sigma(\mathcal{C})$ -measurable, but the converse is not immediate. Thus we essentially get a statement of the λ - π systems lemma but for functions on analogous systems.

We can firstly see that \mathcal{H} is closed in $\mathcal{F}_b(\Omega)$. Then, we can prove the statement for the special case of $\mathcal{C} = \{\chi_A : A \in \mathcal{A}\}$ for a π -system \mathcal{A} , then adding 1 to \mathcal{C} without loss of generality we can make the proof more general (concretely, because $\sigma(\mathcal{C}) \subseteq \sigma(\mathcal{C} \cup \{1\})$).

It may allow this theorem to make more sense to note that λ -systems are sometimes referred to as 'monotone classes'. Thus the π - λ systems lemma can be seen as saying that for \mathcal{A} a π -system, the (smallest) monotone class generated by \mathcal{A} is $\sigma(\mathcal{A})$.

We can use the monotone class theorem to demonstrate that for $f:\Omega_1\times\Omega_2\to\mathbb{R}$ is measurable, then fixing $\omega_1 \in \Omega_1$, $\omega_2 \mapsto f(\omega_1, \omega_2)$ is measurable.

Conditional Probability

Up until presently, we've considered the notion of event A conditioned on event B as having a fixed probability. This doesn't entirely capture what a conditional is however – we're conditioning on the amount of information we have, and therefore we want the conditional probability to change as a function of our information. In particular, we want our conditional probability to be a function of $\omega \in \Omega$, and in order to reflect conditioning as a reflection of information, we want to condition over events in a σ -algebra, rather than individual events.

Doing more algebra, we see that expectation is a more fitting operator, leading us to the following **Definition 7 (Conditional expectation)** Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, $X \in \mathcal{L}^1(\Omega, \mathcal{F}, \mathbb{P})$, $\mathcal{G}\subseteq\mathcal{F}$ a σ -algebra. A random variable $Y\in\mathcal{L}^1(\Omega,\mathcal{G},\mathbb{P})$ is (a version of) the conditional expectation of X given \mathcal{G} if for $G \in \mathcal{G}$,

 $\mathbb{E}[Y\chi_G] = \mathbb{E}[X\chi_G].$ The key aspect of this statement can be rewritten as

 $\int_{\mathcal{C}} \mathbb{E}[X \mid \mathcal{G}] \, d\mathbb{P} = \int_{\mathcal{C}} X \, d\mathbb{P},$

Theorem 3 The conditional expectation of X given \mathcal{G} exists, denoted $\mathbb{E}[X | \mathcal{G}]$, and if Z is also the conditional expectation of X given \mathcal{G} , then $Z = \mathbb{E}[X \mid \mathcal{G}]$ a.s.

which allows us to carry over all of our normal integration properties to conditional expectations.

Come back to the proof of this.

Note importantly how conditional expectations behave with respect to measurability. If we have X a \mathcal{G} -measurable random variable, then

 $\mathbb{E}[X \mid \mathcal{G}] \stackrel{\text{a.s.}}{=} X.$ Meanwhile if $\sigma(X)$ and \mathcal{G} are independent, then $\mathbb{E}[X \mid \mathcal{G}] \stackrel{\mathrm{a.s.}}{=} \mathbb{E}[X].$

Lemma 4 (Tower property) Take $X \in \mathcal{L}^1(\Omega, \mathcal{F}, \mathbb{P})$, \mathcal{F}_1 , \mathcal{F}_2 both σ -algebras, satisfying $\mathcal{F}_1 \subseteq$ $\mathcal{F}_2 \subseteq \mathcal{F}$. Then

 $\mathbb{E}\Big[\mathbb{E}ig[X\,|\,\mathcal{F}_2ig]\,|\,\mathcal{F}_1\Big] \stackrel{ ext{a.s.}}{=} \mathbb{E}ig[X\,|\,\mathcal{F}_1ig].$

This should be relatively intuitive – $\mathbb{E}[X | \mathcal{F}_2]$ contains more information than can be represented in \mathcal{F}_1 , but is fundamentally still expressing a reduced form of X, which can be reduced more to give $\mathbb{E}[X | \mathcal{F}_1]$.

One can also consider this as a commutativity statement: as $\mathbb{E}[X | \mathcal{F}_1]$ is \mathcal{F}_2 -measurable, thus $\mathbb{E}[X | \mathcal{F}_1] = \mathbb{E}[\mathbb{E}[X | \mathcal{F}_1] | \mathcal{F}_2]$, so the tower property is stating that with $\mathcal{F}_1 \subseteq \mathcal{F}_2$: $\mathbb{E}\Big[\mathbb{E}ig[X\,|\,\mathcal{F}_1ig]\,|\,\mathcal{F}_2\Big] \stackrel{ ext{a.s.}}{=} \mathbb{E}\Big[\mathbb{E}ig[X\,|\,\mathcal{F}_2ig]\,|\,\mathcal{F}_1\Big].$

I'm tempted to claim the more general statement, that for \mathcal{F}_1 , \mathcal{F}_2 both σ -algebras in \mathcal{F} : $\mathbb{E}\left[\mathbb{E}\left[X\,|\,\mathcal{F}_1
ight]\,|\,\mathcal{F}_2
ight]\stackrel{\mathrm{a.s.}}{=}\mathbb{E}\left[X\,|\,\mathcal{F}_1\cap\mathcal{F}_2
ight].$

This statement is true if \mathcal{F}_1 and \mathcal{F}_2 are independent, because then both sides are equal to $\mathbb{E}[X]$ but it seems that there could be a 'middle-ground' between independence and containment for which commutativity stops holding.

On the other side, attempting to prove this statement, the tripping point is that it's unclear that the LHS is \mathcal{F}_1 -measurable (although clearly it is \mathcal{F}_2 -measurable).

Lemma 5 Take X, Y random variables on $(\Omega, \mathcal{F}, \mathbb{P})$ with X, Y, and XY integrable. Then $\mathbb{E}[XY \mid \sigma(Y)] \stackrel{\text{a.s.}}{=} Y \mathbb{E}[X \mid \sigma(Y)].$

Ensure for yourself that it's clear why this implies the same holding for $\mathcal{G} \supseteq \sigma(Y)$ instead of $\sigma(Y)$. **Theorem 6** Take $X \in \mathcal{L}^1(\Omega, \mathcal{F}, \mathbb{P})$, $\{\mathcal{F}_i : i \in I\}$ a family of σ -algebras in \mathcal{F} . Then $\{\mathbb{E}|X|\mathcal{F}_i|:i\in I\}$ is uniformly integrable.

Using conditional expectation, we can introduce an inner product to \mathcal{L}^2 , $\langle X, Y \rangle := \mathbb{E}[XY]$ ('introduce' is probably slightly strong, this inner product already exists for other purposes in functional analysis – although we usually use the Lebesgue measure). This gives us that \mathcal{L}^2 is a Hilbert space, and all of the corresponding results.

Measures on $\mathbb R$

Definition 8 A measure space is a triple $(\Omega, \mathcal{F}, \mu)$ such that Ω is a set, \mathcal{F} is a σ -algebra on Ω , and $\mu: \mathcal{F} \to [0, \infty]$ is countably additive (μ is then a measure on (Ω, \mathcal{F})).

Definition 9 Let μ be a probability measure on $\mathcal{B}(\mathbb{R})$. The distribution function of μ is $F_{\mu}(x) = \mu(-\infty, x]$, where we require that F_{μ} is non-decreasing, tends to 0 as $x \to -\infty$, to 1 as $x \to \infty$, and is right continuous.

Definition 10 We say that ν is absolutely continuous with respect to μ , $\nu \ll \mu$, if for any $A \in \mathcal{F}$, $\mu(A) = 0$ implies that $\nu(A) = 0$. Further, we say that μ and ν are equivalent, $\mu \sim \nu$, if $\mu \ll \nu$ and $\nu \ll \mu$.

Extensions

For the most part, it's difficult to characterise a measure explicitly, due to σ -algebras being incredibly large in all but countable Ω . We therefore wish to characterise them in terms of their value on algebras.

Theorem 7 (Uniqueness of extension) Let μ_1 and μ_2 be measures on a space (Ω, \mathcal{F}) , and $\mathcal{A} \subseteq \mathcal{F}$ is a π -system with $\sigma(\mathcal{A}) = \mathcal{F}$. Then if $\mu_1(\Omega) = \mu_2(\Omega) < \infty$ and $\mu_1|_{\mathcal{A}} = \mu_2|_{\mathcal{A}}$, then

This follows immediately via the λ - π systems lemma.

Theorem 8 (Carathéodory Extension theorem) Let Ω be a set and A an algebra on Ω , then with $\mu_0: \mathcal{A} \to [0,\infty]$ a countably additive set function, there exists a measure $\mu:$ $\sigma(\mathcal{A}) \to [0, \infty] \text{ such that } \mu|_{\mathcal{A}} = \mu_0.$

One can derive this from defining the outer measure μ^* in terms of μ_0 , and claiming that a set is measurable iff for all $E \subseteq \Omega$, $\mu^*(E) = \mu^*(E \cap B) + \mu^*(E \setminus B)$. We can then prove that this gives the smallest σ -algebra containing \mathcal{A} .

Definition 11 (Distribution function) If a function $F : \mathbb{R} \to [0, 1]$ satisfies:

i. F is non-decreasing; ii. $F(x) \to 0$ as $x \to -\infty$, $F(x) \to 1$ as $x \to \infty$; and iii. F is continuous from the right,

then F is a distribution function.

Theorem 9 Let F be a distribution function. Then there exists a unique Borel probability measure μ on \mathbb{R} such that $\mu(-\infty,x]=F(x)$. Further, every Borel probability measure on \mathbb{R} defines a distribution function.

A corollary of this is that there is a unique Borel measure such that for all $a < b \in \mathbb{R}$, $\mu(a, b) = b - a$. This result demonstrates that there is a bijection between measures on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ and distribution

functions. In particular, we call these measures the Lebesgue-Stieltjes measures. The proof of this theorem follows using both of the extension theorems. In particular, we use the algebra of left open right closed intervals.

Definition 12 (Pushforward measure) Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, $X : \Omega \to E$, then for $A \in \mathcal{E}$, $\mathbb{Q}(E) := \mathbb{P}(X^{-1}(E))$

Thus we have a (very non-injective) map from random variables in E to probability measures on E. This is useful, on the basis that many properties of random variables will just be properties of the pushforward measure rather than the random variable itself.

Theorem 10 For $\{(\Omega_i, \mathcal{F}_i, \mathbb{P}_i) : i \in \{1, \dots, n\}\}$ a set of probability measures, there is a unique measure \mathbb{P} on $\left(\prod_{i=1}^n \Omega_i, \overset{n}{\underset{i=1}{\times}} \mathcal{F}_i\right)$ such that for $E_i \in \mathcal{F}_i$ with $i \in \{1, \ldots, n\}$,

 $\mathbb{P}\left(\prod_{i=1}^n E_i\right) = \prod_{i=1}^n \mathbb{P}_i(E_i).$

It's hopefully natural here that one should aim an induction proof.

The theorem here then allows us to extend the matter to infinite products, although at this point we require that what we're dealing with are probability measures (to keep each term in [0,1] for convergence reasons), rather than just finite measures as could work with the previous statement of the theorem.

Independence

Definition 13 (Independence) With $(\Omega, \mathcal{F}, \mathbb{P})$ a probability space, $(\mathcal{G}_i)_{i=1}^n$ a collection of σ algebras, these σ -algebras are independent if for $E_i \in \mathcal{G}_i$ for $i \in \{1, \ldots, n\}$

Further, an arbitrary collection $(G_i)_{i\in I}$ of σ -algebras is independent if any finite subset of the collection is independent.

Note that this means $\{\emptyset, \Omega\}$ is independent of anything else.

Additionally, we say that a set $(X_i)_{i\in I}$ of random variables is independent iff $(\sigma(X_i))_{i\in I}$ is independent.

This definition requires a bit of work to deal with properly. One of the best general results we can attain quickly gives a fairly applicable result: **Theorem 11** With $(\Omega, \mathcal{F}, \mathbb{P})$ a probability space, $(\mathcal{A}_i)_{i \in I}$ an arbitrary collection of π -systems, then $(\sigma(A_i))_{i\in I}$ are independent iff for any finite $J\subseteq I$, $A_i\in A_i$ for $i\in J$: $\mathbb{P}\left(\bigcap_{i\in J}A_i\right) = \prod_{i\in J}\mathbb{P}(A_i)$

We also have the result that for any independent set of σ -algebras, any subset is also independent.

It takes a small bit of proving, but from the above results we get the lemma: **Lemma 12** With $(\Omega, \mathcal{F}, \mathbb{P})$, a family of independent random variables $X_i : \Omega \to E_i$, measurable functions $f_i: E_i \to \mathbb{R}$ for $i \in I$, then $(f(X_i))_{i \in I}$ are independent.

Tail events **Definition 14** For a sequence of random variables (X_n) , the tail σ -algebra is defined as

 $\mathcal{T} = \bigcap \sigma(\{X_k : k > n\})$

The intuition here is that all events in the tail σ -algebra contain sample information distinguishing the results of functions of infinitely many sequence elements.

Theorem 13 (Kolmogorov's 0-1 Law) Let (X_n) be a sequence of independent random vari-

ables. Then the tail σ -algebra of (X_n) contains only events with probability 0 or 1. To see this, we demonstrate that \mathcal{T} is independent of a σ -algebra containing it, and therefore that all of its events are independent with themselves.

Borel-Cantelli lemmas **Definition 15** With (A_n) a sequence of sets from \mathcal{F} :

 $\limsup A_n = \bigcap \bigcup A_m$ $= \{ \omega \in \Omega : \omega \in A_n \text{ for infinitely many } n \}$ $= \{A_n \text{ infinitely often}\}$ and $\liminf_{n\to\infty} A_n = \bigcup \bigcap A_m$ $= \{ \omega \in \Omega : \omega \in A_n \ eventually \}$ $= \{A_n \ eventually\}$

Lemma 14 (Fatou and Reverse Fatou for sets) With (A_n) a sequence of sets in \mathcal{F} , $\mathbb{P}(\liminf A_n) \leq \liminf \mathbb{P}(A_n)$ $\mathbb{P}(\limsup A_n) \ge \limsup \mathbb{P}(A_n).$

Lemma 15 (First Borel-Cantelli lemma) For (A_n) a sequence of events in \mathcal{F} , if

then $\mathbb{P}(A_n \ i.o.) = 0$.

Lemma 16 (Second Borel-Cantelli lemma) For (A_n) a sequence of independent events in

 $\sum \mathbb{P}(A_n) = \infty,$

 $\sum \mathbb{P}(A_n) < \infty,$

then $\mathbb{P}(A_n \ i.o.) = 1$.

By their nature, the BC lemmas are only informative in relation to almost sure events. While this may seem incredibly limited, by Kolmogorov's 0-1 Law, it turns out that many events of interest are in fact almost sure events.

As already covered in Part A Integration, we define integration as normal: **Definition 16 (Integral on simple functions)** For a measure space $(\Omega, \mathcal{F}, \mu)$, $f: \Omega \to [0, \infty]$ a non-negative simple function taking values $\{a_1, \ldots, a_n\} \subseteq \mathbb{R}$:

 $\int f d\mu = \sum_{i=1}^{n} a_i \mu \left(f^{-1} \left(\{ a_i \} \right) \right).$ **Definition 17 (Integral on non-negative functions)** For $f: \Omega \to [0, \infty]$ a non-negative mea-

 $\int f = \sup \left\{ \int g \, \mathrm{d}\mu : g \, simple \,, \, 0 \le g \le f \right\}$

 $f^{+} = \max(f, 0), f^{-} = -\min(f, 0), and$ $\int f \, \mathrm{d}\mu = \int f^+ \, \mathrm{d}\mu - \int f^- \, \mathrm{d}\mu$

Definition 18 (Integral) For $f: \Omega \to \mathbb{R}$ a measurable function and $\int |f| d\mu < \infty$, we write

Theorem 17 (Monotone convergence theorem) For (f_n) a sequence of non-negative functions measurable on $(\Omega, \mathcal{F}, \mu)$, such that $f_n \to f$ monotonically. Then

 $\int f_n d\mu \to \int f d\mu$. **Theorem 18 (Fatou's lemma)** For (f_n) a sequence of non-negative functions measurable on

 $\lim \inf_{n \to \infty} f_n \, \mathrm{d}\mu \le \lim \inf_{n \to \infty} \int f_n \, \mathrm{d}\mu$

Lemma 19 (Reverse Fatou's lemma) For (f_n) a sequence of non-negative functions measur-

able on $(\Omega, \mathcal{F}, \mu)$, assume that there is an integrable function g such that $f_n \leq g$ for $n \geq 1$.

 $\int \limsup f_n d\mu \ge \limsup \int f_n d\mu$ **Theorem 20 (Dominated convergence theorem)** For (f_n) a sequence of functions measurable on $(\Omega, \mathcal{F}, \mu)$ with $f_n \to f$ pointwise. Assume that there is an integrable function g such

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

Theorem 21 Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, $X : \Omega \to E$, $g : E \to \mathbb{R}$ measurable on their respective spaces. Then g is $(\mathbb{P} \circ X^{-1})$ -integrable iff $g \circ X$ is \mathbb{P} -integrable. Further, $\int_{-} g(x) d(\mathbb{P} \circ X^{-1})(x) = \int g(X(\omega)) d\mathbb{P}(\omega).$

Definition 19 (Expectation) For $X \in \mathcal{L}^1(\Omega, \mathcal{F}, \mathbb{P})$,

that $|f_n| \leq g$ for $n \geq 1$. Then

We say that X admits an nth moment if $X \in \mathcal{L}^n(\Omega, \mathcal{F}, \mathbb{P})$.

Definition 20 (Variance) For $X \in \mathcal{L}^2(\Omega, \mathcal{F}, \mathbb{P})$, $\operatorname{Var}(X) = \mathbb{E}\left[(X - \mathbb{E}[X])^2 \right].$

Definition 21 (Standardised moment) If $X \in \mathcal{L}^n(\Omega, \mathcal{F}, \mathbb{P})$, the nth standardised moment of $\mathbb{E}\left[\left(\frac{X - \mathbb{E}[X]}{\sqrt{\operatorname{Var}(X)}}\right)^n\right]$

Theorem 22 (Fubini-Tonelli) Let $(\Omega, \mathcal{F}, \mathbb{P})$ be the product of probability spaces $(\Omega_i, \mathcal{F}_i, \mathbb{P}_i)$ for $i \in \{1,2\}, \ and \ f: \Omega \to \mathbb{R} \ is \ a \ bounded \ measurable \ function. \ Then \ both$ $x \mapsto \int f(x,y) d\mathbb{P}_2(y)$

and $y \mapsto \int f(x,y) d\mathbb{P}_1(x)$

are measurable (respectively in \mathcal{F}_1 and \mathcal{F}_2).

If either $f \geq 0$ or f is \mathbb{P} -integrable over Ω , then $\int_{\mathbb{R}} f \, d\mathbb{P} = \int_{\mathbb{R}} \int_{\mathbb{R}} f(x, y) \, d\mathbb{P}_1(x) \, d\mathbb{P}_2(y) = \int_{\mathbb{R}} \int_{\mathbb{R}} f(x, y) \, d\mathbb{P}_2(y) \, d\mathbb{P}_1(x)$

Radon-Nikodym theorem

Integration as we've defined it gives a canonical method of defining a measure on a space: as the tegral of a fixed non-negative measurable function over the set being measured.

Concretely: with a measure space $(\Omega, \mathcal{F}, \mu)$, a measurable function $f: \Omega \to [0, \infty], A \in \mathcal{F}$,

is a measure on \mathcal{F} (via MCT).

We therefore want to characterise how often a measure can be characterised in this way (in particular, whether we can get this the case with respect to leb or the counting measure, both for which we have a wealth of tools).

Theorem 23 (Radon-Nikodym theorem) Let μ , ν be two probability measures on a σ algebra (Ω, \mathcal{F}) . Then $\nu \ll \mu$ if and only if there is a measurable function $f: \Omega \to [0, \infty]$ such that for $A \in \mathcal{F}$,

 $\nu(A) = \int f \, \mathrm{d}\mu.$

Further, $\nu \sim \mu$ if and only if $\mu(f^{-1}(\{0\})) = \nu(f^{-1}(\{0\})) = 0$.

We call f the radon-nikodym derivative of ν with respect to μ , $f = \frac{d\nu}{d\nu}$. If $\nu \sim \mu$, then $\frac{1}{f} = \frac{d\mu}{d\nu}$. This means that providing leb(A) = 0 implies that $\nu(A) = 0$, we can construct ν in this way. Using this theorem, we can define for $A, B \in \mathcal{F}$ the conditional distribution $\mathbb{P}(A \mid B)$, provided $\mathbb{P}(B) > 0$. If $\mathbb{P}(A) = 0$, then $\frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = 0$, so there is some $f_B : \Omega \to [0, \infty]$ measurable such

 $\frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \int f_B(\omega) \, \mathrm{d}\mathbb{P}(\omega).$

Convergence

We now consider modes of convergence of random variables using our results in integration.

Definition 22 (\mathcal{L}^p **spaces)** For $p \geq 0$,

 $\mathcal{L}^p(\Omega, \mathcal{F}, \mathbb{P}) = \{X : \Omega \to \mathbb{R} \ measurable \ s.t. \ \mathbb{E}[|X|^p] < \infty\}$ In particular, \mathcal{L}^0 is the space of all random variables, and \mathcal{L}^{∞} is the space of random variables which are bounded almost surely.

In these notes I'll refer almost exclusively to \mathcal{L}^p , although it's worth noting that there's a nuance here: the well-behaved space we generally want to refer to for useful results is not \mathcal{L}^p , but rather $L^p:=\mathcal{L}^p/\mathcal{N}$, where $\mathcal{N}:=\{X\in\mathcal{L}^0:\mathbb{E}[|X|]=0\}$. At the same time, in this course it's not particularly desirable to be working with $X + \mathcal{N}$ constantly (not least because there are some

instances where a property being 'just' almost sure is relevant), and therefore we principally use \mathcal{L}^p . An exception to the comments above comes for $0 \le p < 1$. L^p for these values is almost entirely useless, as it is not a normed space. \mathcal{L}^p does have some use however, in particular for p=0, in that it allows us to specify the set of measurable functions.

Definition 23 For a sequence (X_n) of random variables over $(\Omega, \mathcal{F}, \mathbb{P})$, we say that X_n converges to X: i. almost surely $(X_n \stackrel{\text{a.s.}}{\rightarrow} X \text{ or } X_n \rightarrow X \text{ a.s.})$ if

 $\mathbb{P}(X_n \to X \text{ as } n \to \infty) = 1.$ ii. in probability $(X_n \stackrel{\mathbb{P}}{\to} X)$ if for all $\varepsilon > 0$

 $\mathbb{P}(|X_n - X| > \varepsilon) \to 0$ $as n \to \infty$. iii. in \mathcal{L}^p $(X_n \stackrel{\mathcal{L}^p}{\to} X)$ if $X_n \in \mathcal{L}^p$ for $n \geq 1$ and

 $\mathbb{E}[|X_n - X|^p] \to 0$ as $n \to \infty$. iv. weakly in \mathcal{L}^1 if $X_n \in \mathcal{L}^1$ for $n \geq 1$ and for all $Y \in \mathcal{L}^{\infty}$ $\mathbb{E}[X_nY] \to \mathbb{E}[XY]$

as $n \to \infty$. v. in distribution $(X_n \xrightarrow{a} X)$ if for $x \in \mathbb{R}$ such that F_X is continuous, $F_{X_n}(x) \to F_X(x)$

as $n \to \infty$.

Of these, notion (v) of convergence in distribution is the odd one out, as it is independent of any particular instance of a random variable with the same distribution. This is a strictly weaker property than convergence in probability.

Convergence in \mathcal{L}^p is an identical notion to that in functional analysis, and (in a mathematical sense) leans very much into the measure-based notion of random variables. It is stronger both than weak convergence in \mathcal{L}^1 and convergence in probability.

Theorem 24 For a sequence (X_n) of random variables, i. If $X_n \stackrel{\text{a.s.}}{\to} X$ then $X_n \stackrel{\mathbb{P}}{\to} X$. ii. If $X_n \stackrel{\mathbb{P}}{\to} X$ then there is a subsequence (X_{n_k}) such that $X_{n_k} \stackrel{\text{a.s.}}{\to} X$.

Useful results **Lemma 25 (Markov's inequality)** Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, X a non-negative random variable. Then for $\lambda > 0$,

Corollary 1 (General Chebyshev's Inequality) For a measurable set $A \subseteq \mathbb{R}$, $X : \Omega \to A$ a random variable, $\varphi:A\to [0,\infty]$ an increasing measurable function. For $\lambda\in A$ with $\varphi(\lambda) < \infty$ we have

 $\mathbb{P}(X \ge \lambda) \le \frac{\mathbb{E}[X]}{\lambda}$

 $\mathbb{P}(X \ge \lambda) \le \frac{\mathbb{E}[\varphi(X)]}{\varphi(\lambda)}.$ This allows us to then demonstrate that for p > 0, $X_n \stackrel{\mathcal{L}^p}{\to} X$ implies $X_n \stackrel{\mathbb{P}}{\to} X$.

Further, we can also show the weak law of large numbers: Corollary 2 For (X_n) be a sequence of i.i.d. random variables with mean μ , variance σ^2 ,

 $\frac{1}{n}\sum X_k \to \mu$ as $n \to \infty$.

Theorem 26 (Jensen's inequality) Let $f: I \to \mathbb{R}$ be a convex function on an interval $I \subseteq \mathbb{R}$. If $X : \Omega \to I$ is an integrable random variable then

For considering the \mathcal{L}^p spaces, we define $\|\cdot\|_p := (\mathbb{E}[|X|^p])^{1/p}$. We can note immediately that for $0 \le p \le q, \, \mathcal{L}^q \subseteq \mathcal{L}^p.$

 $\mathbb{E}[f(X)] \ge f(\mathbb{E}[X]).$

Aside from this, note all the standard results regarding \mathcal{L}^p spaces, in particular Hölder's inequality. A particular application of use is to note that for $1 < p, q < \infty$ with 1/p + 1/q = 1, and for x > 0 $x\mathbb{P}(X \ge x) \le \mathbb{E}[Y\chi_{X>x}],$

In \mathcal{L}^2 , we are able to introduce some geometry to the space, defining an inner product $\langle X,Y\rangle:=$ $\mathbb{E}[XY]$. This allows us to consider \mathcal{L}^2 as a Hilbert space (modulus random variables equal to 0 almost everywhere). Consequently we get all of the results that we normally have for Hilbert spaces within \mathcal{L}^2 .

Uniform Integrability

i. $\{X_n : n \geq 1\}$ is uniformly integrable.

 $ii. X \in \mathcal{L}^1 \ and \ \mathbb{E}[|X_n - X|] \to 0 \ as \ n \to \infty.$

 $iii. X \in \mathcal{L}^1 \ and \ \mathbb{E}[|X_n|] \to \mathbb{E}[|X|] \ as \ n \to \infty.$

then $||X||_p \le q||Y||_p$.

We've now introduced the notions of convergence in \mathcal{L}^1 , and seen that it implies convergence in probability. We'd quite like to go the other direction however, and find a sufficient condition for which the two are the same.

Definition 24 (Uniform integrability) A collection C of random variables is called uniformly integrable (UI) if

 $\lim_{N \to \infty} \sup_{X \in \mathcal{C}} \mathbb{E}[|X|\chi_{|X|>N}] = 0$

Importantly, this is a property of collections, rather than individual random variables. The larger our collection is, the less likely it is to hold (and conversely, if \mathcal{C} is UI, then $\mathcal{D} \subseteq \mathcal{C}$ is also UI). We can see immediately therefore that this is a property that can only hold for collections of random variables in \mathcal{L}^1 , by considering the singleton sets which can be UI.

If we have a $Y \in \mathcal{L}^1$ such that for $X \in \mathcal{C}$, $|X| \leq Y$, then \mathcal{C} is uniformly integrable. This is one of the most common methods for demonstrating a collection is UI.

Another useful characterisation is found by reframing $|X|\chi_{|X|>N}$. From the below fact, we can see

that we can replace it with $(|X|-N)^+$ in the definition of uniform integrability without effect:

 $0 \le (|X| - N)^+ \le |X| \chi_{|X| > N} \le 2(|X| - N/2)^+$. A common formulation that allows us to avoid having to determine the sets on which |X| > N is given below: **Lemma 27** Let \mathcal{C} be a family of random variables. Then \mathcal{C} is UI iff $\sup \mathbb{E}[|X|] < \infty$ and

 $\sup \left\{ \mathbb{E} \left[|X|\chi_A \right] : X \in \mathcal{C}, A \in \mathcal{F}, \mathbb{P}(A) \leq \delta \right\} \to 0$

as $\delta \to 0$. **Theorem 28 (Vitali's convergence theorem)** Take (X_n) a sequence of integrable random variables which converge in probability to a random variable X. The following are equiv-

It's worth noting that (ii) is the definition of convergence in \mathcal{L}^1 , from which we get (iii) just via the reverse triangle inequality. The difficult part of this proof is demonstrating that (iii) entails (i) which requires that characterisation of UI in terms of $(|X|-N)^+$ in conjunction with convergence in probability.

By a non-examinable result (Dunford-Pettis), a collection is UI iff its closure is compact in $\sigma(L^1, L^{\infty})$, which is the weak topology on L^1 .

Filtrations

Definition 25 (Filtrations) For $(\Omega, \mathcal{F}, \mathbb{P})$ a probability space, a filtration (\mathcal{F}_n) is an sequence of σ -algebras ($\mathcal{F}_n \subseteq \mathcal{F}_{n+1}$ for $n \geq 1$). We then call $(\Omega, \mathcal{F}, (\mathcal{F}_n), \mathbb{P})$ a filtered probability space.

Definition 26 (Adapted stochastic process) Take (X_n) a sequence of random variables, (\mathcal{F}_n)

a filtration. If for $n \geq 1$, X_n is \mathcal{F}_n measurable, then (X_n) is said to be adapted to (\mathcal{F}_n) .

Note that for a stochastic process (X_n) , we get a natural filtration $\mathcal{F}_n := \sigma\left(\left\{X_k : k \le n\right\}\right).$

This is the smallest filtration possible.

It's worth noting that with $\mathcal{F}_{\infty} = \bigcup_{n>1} \mathcal{F}_n$, we don't necessarily have that $\mathcal{F}_{\infty} = \mathcal{F}$. Nonetheless, once we're working primarily within the filtered space, this shouldn't make a significant amount of difference, because the remainder of events will be inaccessible to any adapted process.

variable $\tau:\Omega\to\mathbb{N}\cup\{\infty\}$ is a stopping time with respect to (\mathcal{F}_n) if for $n\geq 1$, $\tau^{-1}(\{n\})\in\mathcal{F}_n$.

the maxima or minima of any two stopping times. Further, the first hitting time for an adapted

Definition 28 Let τ be a stopping time on $(\Omega, \mathcal{F}, (\mathcal{F}_n), \mathbb{P})$. The σ -algebra of information at

This definition isn't entirely straightforward to understand, but the key aspect is this: if we have the value of τ (say we know that $\omega \in {\tau = n}$), then we shouldn't be able to infer from $\omega \in A \in \mathcal{F}_{\tau}$ that $\omega \in B \in \mathcal{F}_{\infty} \setminus \mathcal{F}_n$. The above definition should then be clear as containing the events which

We can then get our intuition that if for stopping times τ and ρ , $\tau \leq \rho$, then $\mathcal{F}_{\tau} \subseteq \mathcal{F}_{\rho}$ (we get

The filtration $(\mathcal{F}_{\min(n,\tau)})$ comes up often, and can be conceived essentially as the information given

Martingales

We often consider notions of random walks in probability theory. The most interesting of these are random walks where the movement at each step has expectation 0, and this is a concept generalised by the notion of martingales. In this course we cover only the discrete case of martingales, although

Definition 29 (Martingales) Let $(\Omega, \mathcal{F}, (\mathcal{F}_n), \mathbb{P})$ be a filtered probability space. An integrable, (\mathcal{F}_n) -adapted stochastic process (X_n) is called i. a martingale if for $n \geq 0$, $\mathbb{E}[X_{n+1} | \mathcal{F}_n] = X_n$ a.s.;

aesthetic – negate one and you get the other. Additionally, we can see that a sequence of random variables is a martingale iff it is both a super- and submartingale. Consequently, statements about

Under this definition, the key property to observe for a submartingale (X_n) is this: for $m, n \geq 0$, $\mathbb{E}[X_n | \mathcal{F}_m] \geq X_{\min(m,n)},$ and the analogous relation follows for martingales and supermartingales. What this says is that a

martingale is a sequence of random variables for which we expect no change on average. Meanwhile,

a submartingale is a sequence which we expect to increase, and a supermartingale one we expect to decrease. To reflect this, we often refer to sequences (Y_n) satisfying $\mathbb{E}|Y_{n+1}|\mathcal{F}_n|=0$ as 'martingale difference sequences'

martingale in a smaller filtration. A very general example of a martingale can be given just as a sum of independent random variables (Y_n) each with mean 0, with respect to the natural filtration. Another cute example of a martingale involves taking an integrable random variable X and an arbitrary filtration, then

Lemma 29 Let (X_n) be a martingale with respect to (\mathcal{F}_n) , and $f: \mathbb{R} \to \mathbb{R}$ convex. Then provided $(f(X_n))$ is a sequence of integrable random variables, then it is a submartingale.

Definition 30 (Predictable process) A sequence (V_n) of random variables is predictable with respect to (\mathcal{F}_n) if $\sigma(V_n) \subseteq \mathcal{F}_{n-1}$ (V_n is \mathcal{F}_{n-1} -measurable) for $n \geq 1$.

 $X_n = \sum V_k(Y_k - Y_{k-1}),$

if each X_n is integrable then (X_n) is a martingale with respect to (\mathcal{F}_n) .

The definition of a submartingale prompts some questions about how exactly the sequence varies from being a martingale, and how exactly we might correct it back to a martingale.

 $X_n = X_0 + M_n + A_n$ where (M_n) is a martingale, (A_n) is predictable with respect to (\mathbb{F}_n) , and $M_0 = A_0 = 0$.

An important consequence of this is known as the angle bracket process. For (M_n) a martingale of

This is the conditional variance of $(M_{n+1} - M_n)$. **Theorem 32** Let $(\Omega, \mathcal{F}, (\mathcal{F}_n), \mathbb{P})$ be a filtered probability space, (X_n) a martingale, and τ a finite stopping time. Then $(X_{\min(n,\tau)})$ is a martingale with respect to (\mathcal{F}_n) and hence $(\mathcal{F}_{\min(n,\tau)})$.

 $\langle M \rangle_{n+1} - \langle M \rangle_n = \mathbb{E} \left[M_{n+1}^2 - M_n^2 \,|\, \mathcal{F}_n \right] = \mathbb{E} \left[(M_{n+1} - M_n)^2 \,|\, \mathcal{F}_n \right].$

of both this and (X_n) . **Theorem 33 (Doob's Optional Sampling theorem)** Let $(\Omega, \mathcal{F}, (\mathcal{F}_n), \mathbb{P})$ be a filtered probability space, (X_n) a martingale, and τ , ρ two finite stopping times, $\tau \leq \rho$. Then

and in particular, $\mathbb{E}[X_{\rho}] = \mathbb{E}[X_{\tau}] = \mathbb{E}[X_{0}]$.

Similarly, if (X_n) is a submartingale, then $\mathbb{E}[X_\rho | \mathcal{F}_\tau] \geq X_\tau$ a.s.. martingale $(X_{\min(n,\rho)} - X_{\min(n,\tau)})$ stopped at τ .

 $\mathbb{E}[X_{\tau}\chi_{\tau<\infty}] = \mathbb{E}[X_0].$ **Theorem 34 (Doob's maximal inequality)** Let (X_n) be a submartingale. Then, for $\lambda > 0$,

 $\lambda \mathbb{P}\left(\bigcup_{k=1}^{n} \{X_k \geq \lambda\}\right) \leq \mathbb{E}\left[X_n \chi_{\bigcup_{k=1}^{n} \{X_k \geq \lambda\}}\right] \leq \mathbb{E}\left[|X_n|\right].$

Essentially, we expect (X_n) to grow, but only so quickly. This should (I think?) follow from considering the hitting time of $[\lambda, \infty)$. We then get that if (X_n) is a martingale in \mathcal{L}^p ,

 $\mathbb{E}[X_n^p] \le \mathbb{E}\left[\max_{k \le n} X_k^p\right] \le \left(\frac{p}{p-1}\right)^p \mathbb{E}[X_n^p].$

Definition 27 (Stopping time) Take $(\Omega, \mathcal{F}, (\mathcal{F}_n), \mathbb{P})$ a filtered probability space. A random

stochastic process on any measurable set is a stopping time.

 $\mathcal{F}_{\tau} = \{ A \in \mathcal{F}_{\infty} : for \ n \geq 1, \ A \cap \{ \tau = n \} \in \mathcal{F}_n \}.$

are possible to have observed by time τ ('time' being defined in accordance to steps in the filtration).

more information if we have longer to discover it). Note now that for (X_n) an adapted process, τ a stopping time, $X^\tau := (X_{\min(n,\tau)})$ is a 'stopped process', and is adapted to both the filtration $(\mathcal{F}_{\min(n,\tau)})$ and hence (\mathcal{F}_n) .

that you stop observing things ('exit the room') at time τ , whenever that happens to be. This is what we mean by τ being a stopping time, in that it stops the observations continuing.

there is a rich theory concerning continuous martingales.

ii. a submartingale if for $n \geq 0$, $\mathbb{E}[X_{n+1} | \mathcal{F}_n] \geq X_n$ a.s.; and iii. a supermartingale if for $n \geq 0$, $\mathbb{E}[X_{n+1} | \mathcal{F}_n] \leq X_n$ a.s.. We can note immediately that the difference between a super- and submartingale is almost purely

submartingales are preferable for their generality to statements about martingales.

If a submartingale is adapted to a smaller filtration, then it is also a submartingale with respect to that filtration. This is because the only information a submartingale is using from the filtration is that in the natural filtration, and consequently the extra information is unnecessary. Note that this doesn't go the other direction though – it's straightforward to construct a large filtration (e.g. the constant filtration (\mathcal{F}) such that (X_n) is not a submartingale, despite potentially being a

defining $X_n := \mathbb{E}[X | \mathcal{F}_n]$.

See this as an application of Jensen's inequality, and note that this has some very wide-reaching

Theorem 30 For $(\Omega, \mathcal{F}, (\mathcal{F}_n), \mathbb{P})$ a filtered probability space, (Y_n) a martingale, (V_n) a predictable process, both with respect to (\mathcal{F}_n) , then defining for $n \geq 0$

Note that the sequence (X_n) is a martingale transform, sometimes denoted $((V \circ Y)_n)$.

Theorem 31 (Doob's Decomposition theorem) Let $(\Omega, \mathcal{F}, (\mathcal{F})_n, \mathbb{P})$ be a filtered probability space, (X_n) an integrable adapted process. i. (X_n) has a Doob decomposition

ii. For any (\widetilde{M}_n) a martingale, (\widetilde{A}_n) predictable such that $X_n = X_0 + M_n + \widetilde{A}_n$, $M_n = M_n$ and $A_n = A_n$ for all $n \geq 0$ almost surely. $iii. (X_n)$ is a sub(super)martingale iff (A_n) is non-decreasing (non-increasing) almost surely.

random variables in \mathcal{L}^2 , we have that (M_n^2) is a submartingale. We write the Doob decomposition

 $M_n^2 = M_0^2 + N_n + \langle M \rangle_n$ where (N_n) is a martingale, $(\langle M \rangle_n)$ an increasing predictable process. Using predictability, we see

This can be proven by noting that $(\chi_{n<\tau})$ is predictable, and thus we can write $(X_{\min(n,\tau)})$ in terms

 $\mathbb{E} ig[X_
ho \, | \, \mathcal{F}_ au ig] \stackrel{ ext{a.s.}}{=} X_ au,$

We can prove this in two steps. Firstly, prove the case for ρ constant, and then consider the

is a submartingale, and for $n \geq 1$,

 $Y_n^{\lambda} := (X_n - \lambda) \chi_{\bigcup_{k=1}^n \{X_k \ge \lambda\}}$

Theorem 35 (Doob's \mathcal{L}^p inequality) Let (X_n) be a non-negative submartingale in \mathcal{L}^p for $p \geq 1$ 1. Then $\max X_k \in \mathcal{L}^p$ and

This theorem is quite an important one, as we essentially just expand the definition of a martingale (and, via decomposition, super- and submartingales). The intuition should be that if (X_n) is a process stopping at ρ , but we only observe it up to τ , the properties of the martingale should allow us to calculate the expected value of X_o as X_τ . Importantly, we also generalise that a martingale's expectation is constant, even if stopping at a random time (as long as the time is bounded). We can also get a generalisation beyond just bounded stopping times: Corollary 3 Let (X_n) be a martingale, τ an a.s. finite stopping time. If either $\{X_n : n \geq 0\}$ is UI, or $\mathbb{E}[\tau] < \infty$ and the sequence $(\mathbb{E}[|X_{n+1} - X_n| | \mathcal{F}_n])$ is bounded, then We would like to keep characterising martingales further – in particular wishing to bound them.