

Chapter 1

Probability

1 Introduction

The following are complementary reading for the course.

- G. Grimmett and D. J. A. Welsh, Probability: An Introduction, 1986
- J. K. Blitzstein and J. Hwang, Introduction to Probability, 2019
- D. F. Anderson et al, Introduction to Probability, 2018
- S. M. Ross, Introduction to Probability Models, 2014
- G. Grimmett and D. Stirzaker, Probability and Random Processes, 2001
- G. Grimmett and D. Stirzaker, One Thousand Exercises in Probability, 2009

Notation. Common notation is all defined precisely in the aforementioned. The controversial and additional things are defined as such: $\mathbb{N} = \{1, 2, 3, \dots\}$, $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$, $\mathbb{R}^{>0} := (0, \infty)$.

1.1 Sample spaces and set theory

Definition 1.1.1. The **sample space** Ω is the set of all possible outcomes of an experiment. An element of the sample space $\omega \in \Omega$ is a **sample point**.

Examples 1.1.2. When flipping a coin $\Omega = \{H, T\}$. When rolling a standard die $\Omega = \{1, 2, 3, 4, 5, 6\}$.

Definition 1.1.3. Subsets of Ω are collections of sample points and called **events**.

Suppose events $A, B \subseteq \Omega$:

- $A \cup B$ is the event that A or B or both occur,
- $A \cap B$ is the event that A and B both occur,
- $A^c = \bar{A}$ is the event that occurs iff A does not occur.

Let \mathcal{I} be a general index set with $A_i \subseteq \Omega$, $\forall i \in \mathcal{I}$ and $B \subseteq \Omega$. The following identities hold.

$$\left(\bigcup_{i \in \mathcal{I}} A_i\right)^c = \bigcap_{i \in \mathcal{I}} A_i^c, \quad \left(\bigcap_{i \in \mathcal{I}} A_i\right)^c = \bigcup_{i \in \mathcal{I}} A_i^c, \quad B \cap \left(\bigcup_{i \in \mathcal{I}} A_i\right) = \bigcap_{i \in \mathcal{I}} (A_i \cup B), \quad B \cup \left(\bigcap_{i \in \mathcal{I}} A_i\right) = \bigcap_{i \in \mathcal{I}} (A_i \cap B).$$

These are **De Morgan's Laws** and **Distributivity** respectively.

Lecture 1
Monday
30/10/2023

Lecture 2
Tuesday
31/10/2023

1.2 Interpretation of probability

Definition 1.2.1. The **Cardinality** of a set, denoted $\text{card}(A)$ or $|A|$ is the number of elements in the set A .

Definition 1.2.2. Two sets have the same cardinality iff there exists a bijection between the them.

Definition 1.2.3. A is **finite** if it has as finite numbers of elements, A is **countably infinite** if there exists a bijection $f : A \rightarrow \mathbb{N}$, A is **countable** if it is finite or countable infinite, A is **uncountable** or **uncountable infinite** if it isn't countable.

Samples spaces can be countable or uncountable.

Definition 1.2.4 (Naive probability). Suppose $|A| < \infty$ and we want to assign a probability to $A \subseteq \Omega$.

$$P_{\text{Naive}}(A) := \frac{|A|}{|\Omega|} \implies P(A^c) = 1 - P(A).$$

This Naive example does not consider when $|A|$ is infinite but of finite area.

Example 1.2.5. Let $\Omega = \{(x, y) \in \mathbb{R}^2, x^2 + y^2 = 1\}$ and $A \subseteq \Omega$. Define:

$$P(A) := \frac{\text{area of } A}{\text{area of } \Omega}$$

In the case where $A = \{(x, y) \in \mathbb{R}^2, x^2 + y^2 = 0.5^2\}$ we have $P(A) = 0.25$

Remark 1.2.6. For classical / naive probability we require $|\Omega| < \infty$ or the “area” of Ω be finite.

Definition 1.2.7 (Limiting frequency). Consider n_{total} repetitions of an experiment and n_A the number of time A occurs.

$$P(A) := \lim_{n_{\text{total}} \rightarrow \infty} \frac{n_A}{n_{\text{total}}}$$

Unfortunately, $n_{\text{total}} \rightarrow \infty$ is often hard to conceive with finite representations not necessarily being representative.

Definition 1.2.8 (Subjective probability). For an event A assign the probability $P(A)$ based on our own personal beliefs. The subjective probability need not be the same for different individuals, and despite its appearance it remains a valid interpretation of probability.

Remark 1.2.9. All three interpretations of probability depend of assumptions about the experiment.

2 Counting

2.1 Multiplication principle

Computing naive probabilities often requires some combinatorics.

Definition 2.1.1 (Multiplication principle). If we perform an experiment A that has a possible outcomes and an experiment B with b possible outcomes (in any order) then the number of outcomes of the **compound experiment** will be ab .

Remark 2.1.2. When dealing with repetitions of the same experiment (with sample space Ω , the sample space is given by the Cartesian product of the individual samples spaces.

$$\Omega_1 \times \Omega_2 \times \cdots \times \Omega_n := \{(\omega_1, \omega_2, \dots, \omega_n) : \omega_i \in \Omega_i\}.$$

The cardinality of this samples space follows from the multiplication principle.

2.2 Power sets

Definition 2.2.1 (Power Set). Given a set A its **power set** is defined as:

$$\mathcal{P}(A) := \{X : X \subseteq A\}.$$

Theorem 2.2.2. If A is a finite set, $|\mathcal{P}(A)| = 2^{|A|}$.

Lecture 3
Friday
03/11/2023

2.3 Combinatorial coefficients

Definition 2.3.1 (Factorial). Let $n \in \mathbb{N}$ the **factorial** of n is defined as:

$$n! := \prod_{i=1}^n i.$$

Definition 2.3.2 (Descending factorial). Let $k, n \in \mathbb{N}$ with $k \leq n$ the **descending factorial** denoted $(n)_k$ is defined as:

$$(n)_k := n(n-1) \dots (n-k+1) = \prod_{i=0}^{k-1} (n-i) = \prod_{j=n-k+1}^n j = \frac{n!}{(n-k)!}.$$

Definition 2.3.3 (Binomial coefficient). Let $k, n \in \mathbb{N}_0$ the **binomial coefficient** is the number of subsets of size k of a set n :

$$\binom{n}{k} := \begin{cases} \frac{n(n-1) \dots (n-(k-1))}{k!} = \frac{(n)_k}{k!} = \frac{n!}{(n-k)!k!} & \text{if } k \leq n \\ 0 & \text{otherwise.} \end{cases}$$

Lecture 4
Monday
06/11/2023

2.4 Sampling with and without replacement

“Definitions” given in the context of drawing balls from an urn, $S = \{1, 2, \dots, n\}$.

Definition 2.4.1 (Ordered sampling with replacement). Take out a ball from S , note its number, put it back; repeat this k times. The sample space for this experiment is $\Omega = S^k$.

Definition 2.4.2 (Ordered sampling without replacement). Take out a ball from S , note its number but **do not** put it back; repeat $k < n$ times. There are $|\Omega| = (n)_k$ possible outcomes.

Definition 2.4.3 (Unordered sampling without replacement). We take k balls out of the urn, there are $\binom{n}{k}$ possibilities.

Definition 2.4.4 (Unordered sampling with replacement). We take k balls out of the urn, with the stars and bars argument: there must be k stars divided by $n-1$ bars giving us:

$$|\Omega| = \binom{n+k-1}{k} = \binom{n+k-1}{n-1}.$$

Lecture 5
Tuesday
07/11/2023

3 Axiomatic probability

3.1 Event space

We do not always want to consider all subsets of Ω so denote $\mathcal{F} \subseteq \mathcal{P}(\Omega)$ the **event space**, which contains the events we are allowed to consider. \mathcal{F} must always be a σ -algebra.

Definition 3.1.1 (Algebra). \mathcal{F} is an **algebra** (or a field) on Ω iff:

1. $\emptyset \in \mathcal{F}$,
2. $A \in \mathcal{F} \implies A^c \in \mathcal{F}$,
3. $A, B \in \mathcal{F} \implies A \cup B \in \mathcal{F}$.

Definition 3.1.2 (σ -algebra). \mathcal{F} is a **σ -algebra** (or a σ -field) on Ω iff:

1. $\emptyset \in \mathcal{F}$;
2. $A \in \mathcal{F} \implies A^c \in \mathcal{F}$,
3. For all i in some countable indexing set \mathcal{I} , $A_i \in \mathcal{F} \implies \bigcup_{i \in \mathcal{I}} A_i \in \mathcal{F}$.

Remark 3.1.3. 1. Any algebra is closed under finite unions and finite intersections,

2. Any σ -algebra is closed under countable intersections,

3. Any (σ -)algebra on Ω contains Ω .

Definition 3.1.4 (Trivial sigma algebra). The **trivial sigma algebra** on Ω is defined as $\mathcal{F}_{trivial} := \{\emptyset, \Omega\}$.

Example 3.1.5 (Smallest σ -algebra of an element). For some $A \subseteq \Omega$, the sigma algebra $\mathcal{F}_A := \{\emptyset, A, A^c, \Omega\}$ is the smallest σ -algebra on Ω (smallest cardinality) that contains A .

Lecture 6
Friday
10/11/2023

3.2 Probability measure

Definition 3.2.1 (Probability measure). A mapping $P : \mathcal{F} \rightarrow \mathbb{R}$ is a **probability measure** on (Ω, \mathcal{F}) iff:
 1. $P(A) \geq 0$ for all $A \in \mathcal{F}$; 2. $P(\Omega) = 1$; 3. for a countable, disjoint sequence of events $(A_i)_{i \in \mathcal{I}}$ on an indexing set \mathcal{I} :

$$P\left(\bigcup_{i \in \mathcal{I}} A_i\right) = \sum_{i \in \mathcal{I}} P(A_i).$$

3.3 Probability space

Definition 3.3.1 (Probability space). A **probability space** is a triple (Ω, \mathcal{F}, P) , with Ω a sample space, \mathcal{F} a σ -algebra on Ω , and P a probability measure on (Ω, \mathcal{F}) .

Corollary 3.3.2. For $A, B \in \mathcal{F}$:

1. $P(A^c) = 1 - P(A)$,
2. $A \subseteq B \implies P(A) \leq P(B)$,
3. $P(A \cup B) = P(A) + P(B) - P(A \cap B)$.

Lecture 7
Monday
13/11/2023

4 Conditional probability

Definition 4.0.1 (Conditional probability measure). Consider the probability space (Ω, \mathcal{F}, P) and some event $B \in \mathcal{F}$ with $P(B) > 0$, we construct the probability measure Q on (Ω, \mathcal{F}) by

$$Q(A) := \frac{P(A \cap B)}{P(B)}.$$

Denote the **conditional probability** of A given B by $P(A|B) = Q(A)$.

Lecture 8
Tuesday
14/11/2023

4.1 Bayes' rule and total probability

Theorem 4.1.1 (Bayes' rule). For $A, B \in \mathcal{F}$ with $P(A) > 0, P(B) > 0$ we have,

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}.$$

Definition 4.1.2 (Partition of a set). A partition of some set Ω is a collection $\{B_i, i \in \mathcal{I}\}$ for some countable index set \mathcal{I} with $B_i \cap B_j = \emptyset$ for all $i, j \in \mathcal{I}$ with $i \neq j$ and $\bigcup_{i \in \mathcal{I}} B_i = \Omega$.

Theorem 4.1.3 (Total probability). Given some partition $\{B_i, i \in \mathcal{I}\}$ of Ω with $P(B_i) > 0$ for all $i \in \mathcal{I}$ and some event $A \in \mathcal{F}$,

$$P(A) = \sum_{i \in \mathcal{I}} P(A \cap B_i) = \sum_{i \in \mathcal{I}} P(A|B_i)P(B_i).$$

These two theorems can then be combined to form the following.

Theorem 4.1.4 (Bayes' rule with extra conditioning). For events $A, B, E \in \mathcal{F}$ with $P(A \cap E) > 0, P(B \cap E) > 0$ we have

$$P(A|B \cap E) = \frac{P(B|A \cap E)P(A|E)}{P(B|E)}.$$

Theorem 4.1.5 (Total probability with extra conditioning). Given events $A, E \in \mathcal{F}$ with $P(E) > 0$ and some partition $\{B_i, i \in \mathcal{I}\}$ of Ω with $P(B_i \cap E) > 0$ for all $i \in \mathcal{I}$,

$$P(A|E) = \sum_{i \in \mathcal{I}} \frac{P(A \cap B_i \cap E)}{P(E)} = \sum_{i \in \mathcal{I}} P(A|B_i \cap E)P(B_i|E).$$

Lecture 9
Friday
17/11/2023

5 Independence

5.1 Event independence

Two events $A, B \in \mathcal{F}$ will be independent iff the occurrence of one does not effect the probability the other occurs, i.e $P(A|B) = P(A)$ and vice versa.

Definition 5.1.1 (Independent events). Two events $A, B \in \mathcal{F}$ are said to be **independent** iff

$$P(A \cap B) = P(A)P(B),$$

and **dependent** otherwise.

Corollary 5.1.2. If A and B are independent then so are all pairs of their complements.

Definition 5.1.3 (General independence). A finite collection of events $\{A_1, A_2, \dots, A_n\}$ is independent iff

$$P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1)P(A_2) \dots P(A_n),$$

and similarly a countably or uncountably infinite collection of events is independent iff each finite subcollection is independent.

5.2 Conditional independence

Definition 5.2.1 (Conditional independence). Given the events $A, B, C \in \mathcal{F}$ with $P(C) > 0$ we say A and B are **conditional independent** given C iff,

$$P(A \cap B|C) = P(A|C)P(B|C).$$

5.3 Product rule for general independence

The upcoming subsection may seem disparate, they are however necessary parts to the omitted proof of the product rule for general independence and therefore deemed relevant.

Definition 5.3.1 (Set difference). Given two set $A, B \in \Omega$ the **set difference** of A and B is defined as, $A \setminus B := A \cap B^c$.

Lemma 5.3.2. Any countable union of sets can be written as a countable union of disjoint sets.

Definition 5.3.3 (Increasing and decreasing sets). A sequence of sets $(A_i)_{i=1}^{\infty}$ is said to **increase** to A (written $A_i \uparrow A$) iff $A_1 \subseteq A_2 \subseteq \dots$ and $\bigcup_{i=1}^{\infty} A_i = A$. The definition for a sequence of sets $(B_i)_{i=1}^{\infty}$ to **decrease** to a set B ($B_i \downarrow B$) is defined similarly.

Theorem 5.3.4 (Continuity property of probability measures). If $A_1, A_2, \dots \in \mathcal{F}$ with $A_i \uparrow A$ or $A_i \downarrow A$ for some $A \in \mathcal{F}$,

$$\lim_{i \rightarrow \infty} P(A_i) = P(\lim_{i \rightarrow \infty} A_i) = P(A).$$

Theorem 5.3.5 (Product rule for general independence). Given a countably infinite set of independent events $A_1, A_2, \dots \in \mathcal{F}$,

$$P\left(\bigcap_{i=1}^{\infty} A_i\right) = \prod_{i=1}^{\infty} P(A_i).$$

6 Discrete random variables

6.1 Images and their properties

throughout this subsection we will be considering the function $f: \mathcal{X} \rightarrow \mathcal{Y}$.

Definition 6.1.1 (Image). For some subset $A \subseteq \mathcal{X}$ we define the **image** of A under f by,

$$f(A) := \{y \in \mathcal{Y} : \exists x \in A, y = f(x)\} = \{f(x) : x \in A\}.$$

When $A = \mathcal{X}$, $f(\mathcal{X}) = \text{im } f$.

Lecture 10
Monday
20/11/2023

Lecture 11
Tuesday
21/11/2023

Definition 6.1.2 (Pre-image). For some subset $B \subseteq \mathcal{Y}$ we now define the **pre-image** of B under f by,

$$f^{-1}(B) := \{x \in \mathcal{X} : f(x) \in B\}.$$

Despite the similar notation to the inverse function of f they are not the same thing. Notably, the pre-image under f always exists while the inverse function need not exist.

Lemma 6.1.3. For a collection of subsets $B_i \in \mathcal{F}$ for i in some indexing set \mathcal{I} we have,

$$f^{-1}\left(\bigcup_{i \in \mathcal{I}} B_i\right) = \bigcup_{i \in \mathcal{I}} f^{-1}(B_i).$$

6.2 DRVs and their distributions

Definition 6.2.1 (Discrete random variable). A **discrete random variable** (DRV) on the probability space (Ω, \mathcal{F}, P) is a function $X : \Omega \rightarrow \mathbb{R}$ that satisfies the following properties:

- $\text{im } X = \{X(\omega) : \omega \in \Omega\}$ must be a countable subset of \mathbb{R} ,
- $X^{-1}(x) \in \mathcal{F}$ for all $x \in \mathbb{R}$.

Remark 6.2.2. The nomenclature of X being discrete stems from the fact that its image is a countable subset of \mathbb{R} and so can be mapped to \mathbb{N} which we see as being discrete.

Definition 6.2.3 (Probability mass function). The **probability mass function** (pmf) of a DRV X is defined as a function $p_X : \mathbb{R} \rightarrow [0, 1]$ such that,

$$p_X(x) := P(X^{-1}(x)).$$

This is commonly denoted by $p_X(x) = P(X = x)$.

Remark 6.2.4. Some useful properties of the pmf extending from the definition are:

- $x \notin \text{im } X \implies p_X(x) = 0$,
- For $x_1, x_2 \in \text{im } X$ with $x_1 \neq x_2$, $X^{-1}(x_1) \cap X^{-1}(x_2) = \emptyset$,
- $\sum_{x \in \text{im } X} p_X(x) = \sum_{x \in \mathbb{R}} p_X(x) = 1$.

Theorem 6.2.5. Suppose \mathcal{I} is some indexing set and $S = \{s_i \in \mathbb{R} : i \in \mathcal{I}\}$ is countable and $\{\pi_i : i \in \mathcal{I}\}$ is a collection such that $\pi_i \geq 0$ for all $i \in \mathcal{I}$ and $\sum_{i \in \mathcal{I}} \pi_i = 1$. Then there exists some probability space (Ω, \mathcal{F}, P) and a DRV X on said probability space such that $p_X(s_i) = \pi_i$ for all $i \in \mathcal{I}$ and $p_X(s) = 0$ otherwise.

7 Common DRVs

All DRVs within this section will be considered over the probability space (Ω, \mathcal{F}, P) .

7.1 Bernoulli distribution

Definition 7.1.1 (Bernoulli distribution). A DRV X is said to have **Bernoulli distribution** with parameter $p \in (0, 1)$ if $\text{im } X = \{0, 1\}$ with $p_X(1) = p$. This is denoted by $X \sim \text{Bern}(p)$.

Definition 7.1.2 (Indicator variable). Given some event $A \in \mathcal{F}$ the **indicator variable** of the event A is given by,

$$\mathbb{I}_A(\omega) := \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases}.$$

Remark 7.1.3. $\mathbb{I}_A \sim \text{Bern}(P(A))$.

7.2 Binomial distribution

Definition 7.2.1 (Binomial distribution). Consider a sequence of $n \in \mathbb{N}$ iid Bernoulli trials with parameter p , count the number of successes and denote this by the random variable X then $\text{im } X = [0, n]$ and,

$$p_X(x) = \binom{n}{x} p^x (1-p)^{n-x} \quad \text{for } x \in [0, n].$$

We say X follows a **binomial distribution** and this is denoted by $X \sim \text{Bin}(n, p)$.

7.3 Hypergeometric distribution

As we have done previously, consider of urn of $N \in \mathbb{N}$ balls with $K \in \mathbb{N}$ of these being white and the remainder being black from which we will draw $n \in \mathbb{N}$ balls and want to consider the DRV (X) for the number of white balls drawn. When drawing with replacement we have $X \sim \text{Bin}(n, K/N)$. However, when drawing without replacement X follows the hypergeometric distribution.

Definition 7.3.1 (Hypergeometric distribution). A DRV X follows the **hypergeometric distribution** with three parameters $N \in \mathbb{N}_0, K \in \mathbb{N}, n \in [0, N]$ if $\text{im } X = [0, \min(n, K)]$ and,

$$p_X(x) = \frac{\binom{K}{x} \binom{N-K}{n-x}}{\binom{N}{n}} \quad \text{for } x \in [0, K].$$

Lemma 7.3.2 (Vandemonde's identity). **Vandemonde's identity** is an important tool in the derivation of the pmf for the hypergeometric distribution and so is included here. The identity is as follows, for $k, m, n \in \mathbb{N}$ with $k \leq m + n$, we have:

$$\binom{m+n}{k} = \sum_{i=0}^k \binom{m}{i} \binom{n}{k-i}.$$

7.4 Discrete uniform distribution

Definition 7.4.1 (Discrete uniform distribution). A DRV X follows the **discrete uniform distribution** over a nonempty set of numbers C , denoted $X \sim \text{DUnif}(C)$, if $\text{im } X = C$ and,

$$p_X(x) = \begin{cases} \frac{1}{\text{card}(C)} & \text{for } x \in C \\ 0 & \text{otherwise} \end{cases}.$$

7.5 Poisson distribution

The poisson distribution is commonly used for modelling the number of events occurring in a certain time period. Its pdf is derived by taking the $\lim_{n \rightarrow \infty} p_X(x)$ where $X \sim \text{Bin}(n, \frac{\lambda}{n})$ for some $\lambda \in \mathbb{R}$.

Definition 7.5.1 (Poisson distribution). A DRV X follows the **poisson distribution** with parameter $\lambda \in \mathbb{R}^{>0}$, denoted $X \sim \text{Poi}(\lambda)$, if $\text{im } X = \mathbb{N}_0$ and,

$$p_X(x) = \frac{\lambda^x}{x!} e^{-\lambda} \quad \text{for } x \in \mathbb{N}_0.$$

7.6 Geometric distribution

Definition 7.6.1 (Geometric distribution). A DRV X follows the **geometric distribution** with parameter $p \in (0, 1)$, denoted $X \sim \text{Geom}(p)$, if $\text{im } X = \mathbb{N}$ and,

$$p_X(x) = (1-p)^{x-1} p \quad \text{for } x \in \mathbb{N}.$$

This can be seen as counting the number of Bernoulli trials with parameter p that occur before a failure.

7.7 Negative binomial distribution

Definition 7.7.1 (Generalised binomial coefficient). Let $\alpha \in \mathbb{C}$ and $k \in \mathbb{N}$ and define the **generalised binomial coefficient** by,

$$\binom{\alpha}{k} := \frac{\alpha(\alpha-1)\dots(\alpha-k+1)}{k!}.$$

Lemma 7.7.2. For $x \in \mathbb{N}_0$ and $r \in \mathbb{N}$ the following identity holds,

$$\binom{x+r-1}{r-1} = (-1)^x \binom{-r}{x}.$$

The generalised binomial coefficient as well as this lemma are necessary to have a well defined and valid pdf for the negative binomial distribution.

Definition 7.7.3 (Negative binomial distribution). A DRV X follows the **negative binomial distribution** with parameters $r \in \mathbb{N}$ and $p \in (0, 1)$, denoted $X \sim \text{NBin}(r, p)$, if $\text{im } X = \mathbb{N}_0$ and,

$$p_X(x) = \binom{x+r-1}{r-1} p^r (1-p)^x \quad \text{for } x \in \mathbb{N}_0.$$

This is the distribution of the number of failed ii Bernoulli trials with parameter p before r successes have occurred.

8 Continuous random variables

8.1 General random variables and their distributions

Definition 8.1.1 (Random variable). A **random variable (RV)** on the probability space (Ω, \mathcal{F}, P) is a mapping $X : \Omega \rightarrow \mathbb{R}$ such that $X^{-1}((-\infty, x]) = \{\omega \in \Omega : X(\omega) \leq x\} \in \mathcal{F}$ for all $x \in \mathbb{R}$. By taking the countable union of pre-images of all $\omega \leq x$ in \mathcal{F} , it can be seen that a DRV satisfies this condition.

Definition 8.1.2 (Cumulative distribution function). For some RV X on the probability space (Ω, \mathcal{F}, P) , the **cumulative distribution function (CDF)** of X is defined as the mapping $F_X : \mathbb{R} \rightarrow [0, 1]$ given by,

$$F_X(x) = P(X^{-1}((-\infty, x])),$$

often denoted $F_X(x) = P(X \leq x)$.

Theorem 8.1.3 (cdf properties). For some RV X on the probability space (Ω, \mathcal{F}, P) the following properties hold:

1. F_X is monotonically non-decreasing,
2. F_X is right-continuous ($(x_n) \downarrow x \implies F_X(x_n) \rightarrow F_X(x)$ as $n \rightarrow \infty$),
3. $\lim_{x \rightarrow -\infty} F_X(x) = 0$ and $\lim_{x \rightarrow \infty} F_X(x) = 1$.

Theorem 8.1.4. For $a, b \in \mathbb{R}$ if $a < b$, then $P(a < X \leq b) = F_X(b) - F_X(a)$.

Remark 8.1.5. In general the cdf of an RV is not left continuous.

8.2 CRVs and pdfs

Definition 8.2.1 (Continuous random variable). A random variable X on the probability space (Ω, \mathcal{F}, P) is called a **continuous random variable (CRV)** iff its cdf can be written as:

$$F_X(x) = \int_{-\infty}^x f_X(u) du \quad \text{for all } x \in \mathbb{R},$$

where $f_X : \mathbb{R} \rightarrow \mathbb{R}$ satisfies: $f_X(u) \geq 0$ for all $u \in \mathbb{R}$ and $\int_{-\infty}^{\infty} f_X(u) du = 1$. We call f_X the **probability density function (pdf)** of X .

Theorem 8.2.2. If X is a CRV on the probability space (Ω, \mathcal{F}, P) with pdf f_X , $P(X = x) = 0$ for all $x \in \mathbb{R}$.

Theorem 8.2.3. With the same conditions, $P(a \leq X \leq b) = \int_a^b f_X(u) du$ for all $a, b \in \mathbb{R}$ with $a \leq b$.

Remark 8.2.4. Combining the results from this section leads to the conclusion that the cdf of a CRV is continuous.

9 Common CRVs

All CRVs X within this section will be considered over the probability space (Ω, \mathcal{F}, P) with the natural notation for their pdf and cdf. These distribution will be uniquely identified by their pdfs.

9.1 Uniform distribution

Definition 9.1.1 (Uniform distribution). A CRV X follows the **uniform distribution** on the interval (a, b) for $a, b \in \mathbb{R}$ with $a < b$, denoted $X \sim U(a, b)$ if it satisfies:

$$f_X(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}, \quad F_X(x) = \begin{cases} 0 & \text{if } x \leq a \\ \frac{1}{b-a} & \text{if } a < x < b \\ 1 & \text{if } x \geq b \end{cases}.$$

9.2 Exponential distribution

Definition 9.2.1 (Exponential distribution). A CRV X follows the **exponential distribution** with parameter $\lambda \in \mathbb{R}^{>0}$, denoted $X \sim \text{Exp}(\lambda)$ if it satisfies:

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}, \quad F_X(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 - e^{-\lambda x} & \text{if } x > 0 \end{cases}.$$

9.3 Gamma distribution

Definition 9.3.1 (Gamma function). For $t \in \mathbb{R}$ with $t > 0$ we define the **gamma function** by,

$$\Gamma(t) := \int_0^\infty x^{t-1} e^{-x} dx.$$

One of the gamma function's many interesting properties is that $\Gamma(t) = (t-1)\Gamma(t-1)$ for $t > 1$.

Definition 9.3.2 (Gamma distribution). A CRV X follows the **gamma distribution** with shape and rate parameter $\alpha, \beta \in \mathbb{R}^{>0}$ respectively, denoted $X \sim \text{Gamma}(\alpha, \beta)$ if it satisfies:

$$f_X(x) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}.$$

Its cdf cannot be written in a closed form so must be left as an integral of the pdf or approximated.

9.4 Chi-squared distribution

Definition 9.4.1 (Chi-squared distribution). A CRV X follows the **chi-squared distribution** with $n \in \mathbb{N}$ degrees of freedom, denoted $X \sim \chi^2(n)$ if it satisfies:

$$f_X(x) = \begin{cases} \frac{1}{2\Gamma(n/2)} \left(\frac{x}{2}\right)^{n/2-1} e^{-x/2} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}.$$

Its cdf can also not be written in a closed form. The $\chi^2(n)$ distribution is the same as the $\text{Gamma}(\frac{n}{2}, \frac{1}{2})$ distribution.

9.5 F-distribution

These pdfs are getting tough.

Definition 9.5.1 (F-distribution). A CRV X follows the **f-distribution** with $d_1, d_2 \in \mathbb{R}^{>0}$ degrees of freedom, denoted $X \sim F(d_1, d_2)$ if it satisfies:

$$f_X(x) = \begin{cases} \frac{\Gamma\left(\frac{d_1+d_2}{2}\right) \left(\frac{d_1}{d_2}\right)^{d_1/2} x^{d_1/2-1}}{\Gamma\left(\frac{d_1}{2}\right) \Gamma\left(\frac{d_2}{2}\right) \left(1 + \frac{d_1}{d_2}x\right)^{(d_1+d_2)/2}} & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}.$$

Its cdf can also not be written in a closed form. It is important to note that d_1, d_2 are not restricted to integer values, and that $X = \frac{X_1/d_1}{X_2/d_2}$ where $X_1 \sim \chi^2(d_1)$ and $X_2 \sim \chi^2(d_2)$.

9.6 Beta distribution

Definition 9.6.1 (Beta function). For $\alpha, \beta \in \mathbb{R}^{>0}$ we define the **beta function** by,

$$B(\alpha, \beta) := \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}.$$

Definition 9.6.2 (Beta distribution). A CRV X follows the **beta distribution** with parameters $\alpha, \beta \in \mathbb{R}^{>0}$, denoted $X \sim \text{Beta}(\alpha, \beta)$ if it satisfies:

$$f_X(x) = \begin{cases} \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} & \text{if } 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}.$$

Its cdf can also not be written in a closed form.

9.7 Normal distribution

Definition 9.7.1 (Standard normal distribution). A CRV X follows the **standard normal distribution** or **Gaussian distribution**, denoted $X \sim N(0, 1)$ if it satisfies,

$$f_X(x) = \phi(x) := \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \quad \text{for } x \in \mathbb{R}, \quad F_X(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt \quad \text{for } x \in \mathbb{R}.$$

Where, once again, there is no explicit formula for the cdf.

Definition 9.7.2 (Normal distribution). A CRV X follows the **normal distribution** with mean $\mu \in \mathbb{R}$ and variance σ^2 for $\sigma \in \mathbb{R}^{>0}$ denoted $X \sim N(\mu, \sigma^2)$ if it satisfies,

$$f_X(x) = \phi(x) := \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \text{for } x \in \mathbb{R}.$$

9.8 Cauchy distribution

Definition 9.8.1 (Cauchy distribution). A CRV X follows the **Cauchy distribution** if it satisfies,

$$f_X(x) = \frac{1}{\pi(1+x^2)} \quad \text{for } x \in \mathbb{R}, \quad F_X(x) = \frac{1}{\pi} \arctan(x) + \frac{1}{2} \quad \text{for } x \in \mathbb{R}.$$

If $X, Y \sim N(0, 1)$, then $Z = X/Y$ follows the Cauchy distribution.

9.9 Student t-distribution

Definition 9.9.1 ((Student's) t-distribution). A CRV X follows the **Student t-distribution** with $\nu \in \mathbb{R}^{>0}$ degrees of freedom if it satisfies,

$$f_X(x) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}} \quad \text{for } x \in \mathbb{R}.$$

Its cdf cannot be written in a closed form.

Remark 9.9.2. Not all RVs are either discrete or continuous.

10 Transformations of random variables

10.1 DRVs

Theorem 10.1.1. Let X be a DRV on (Ω, \mathcal{F}, P) and let $g : \mathbb{R} \rightarrow \mathbb{R}$ be a deterministic function, then $Y = g(X)$ is a DRV with pmf:

$$p_Y(y) = \sum_{\{x \in \text{im } X : g(x)=y\}} p_X(x) \quad \text{if } y \in \text{im } Y \text{ and } 0 \text{ otherwise.}$$

10.2 CRVs

Theorem 10.2.1. Let X be a CRV on (Ω, \mathcal{F}, P) and let $g : \mathbb{R} \rightarrow \mathbb{R}$ be a strictly monotonic and differentiable function with inverse $g^{-1} : \mathbb{R} \rightarrow \mathbb{R}$, then $Y = g(X)$ is a CRV with pdf:

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} [g^{-1}(y)] \right| \quad \text{for all } y \in \mathbb{R}.$$

Remark 10.2.2. The term $\left| \frac{d}{dy} [g^{-1}(y)] \right|$ is often called the **Jacobian** of the transformation.

11 Expectation of random variables

Throughout this section, unless otherwise specified, all infinite sums will be assumed to converge absolutely and all integrals will be assumed to be $< \infty$.

11.1 Definition

Definition 11.1.1 (Expectation of a DRV). Let X be a DRV on (Ω, \mathcal{F}, P) then the **expectation** of X is defined by,

$$E(X) := \sum_{x \in \text{im } X} xp_X(x).$$

Definition 11.1.2 (Expectation of a CRV). Let X be a CRV on (Ω, \mathcal{F}, P) then the **expectation** of X is defined by,

$$E(X) := \int_{-\infty}^{\infty} xf_X(x)dx.$$

11.2 LOTUS

Theorem 11.2.1 (Discrete LOTUS). Let X be a DRV on (Ω, \mathcal{F}, P) and $g : \mathbb{R} \rightarrow \mathbb{R}$, we have,

$$E(g(X)) = \sum_{x \in \text{im } X} g(x)p_X(x).$$

Theorem 11.2.2 (Continuous LOTUS). Let X be a CRV on (Ω, \mathcal{F}, P) and $g : \mathbb{R} \rightarrow \mathbb{R}$, we have,

$$E(g(X)) = \int_{-\infty}^{\infty} g(x)f_X(x)dx.$$

Note that this is one of the few theorems throughout the course given without proof.

Theorem 11.2.3. If X is non-negative then $E(X) \geq 0$.

11.3 Variance

Definition 11.3.1 (Variance). Let X be a discrete/continuous random variable, then the **variance** of X is defined by,

$$\text{Var}(X) := E[X - E(X)]^2.$$

Theorem 11.3.2. For a discrete/continuous random variable with finite variance,

$$\text{Var}(X) = E(X^2) - [E(X)]^2.$$

12 Multivariate random variables

12.1 Multivariate distributions

Definition 12.1.1 (Joint distribution function). Consider the sequence of random variables X_1, X_2, \dots, X_n all on (Ω, \mathcal{F}, P) and write $\mathbf{X} = (X_1, X_2, \dots, X_n)$ and $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$. Then the **joint distribution function** of \mathbf{X} is $F_{\mathbf{X}} : \mathbb{R}^n \rightarrow [0, 1]$ defined by:

$$F_{\mathbf{X}}(\mathbf{x}) := P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n) \quad \text{for all } \mathbf{x} \in \mathbb{R}^n.$$

12.2 Independence

Definition 12.2.1 (Pairwise independence for n random variables). We call the sequence of RVs, X_1, X_2, \dots, X_n , **pairwise independent** iff,

$$F_{X_i, X_j}(x_i, x_j) = F_{X_i}(x_i)F_{X_j}(x_j) \quad \text{for all } x_i, x_j \in \mathbb{R} \text{ with } i \neq j.$$

Definition 12.2.2 (Independence of a family of random variables). Given some indexing set $\mathcal{I} \subset \mathbb{R}$, a family of random variables $\{X_i : i \in \mathcal{I}\}$ is **independent** iff for all finite $\mathcal{J} \subseteq \mathcal{I}$:

$$P\left(\bigcap_{j \in \mathcal{J}} \{X_j \leq x_j\}\right) = \prod_{j \in \mathcal{J}} P(\{X_j \leq x_j\}) \quad \text{for all } (x_j)_{j \in \mathcal{J}} \text{ with } x_j \in \mathbb{R}.$$

(All finite subfamilies of the family of random variables is independent by the natural definition)

12.3 Multivariate DRVs

Definition 12.3.1 (Joint probability mass functions). Let X_1, X_2, \dots, X_n all be DRVs on (Ω, \mathcal{F}, P) that form the random vector \mathbf{X} , their **joint probability mass function**, $p_{\mathbf{X}} : \mathbb{R}^n \rightarrow [0, 1]$ is defined as,

$$p_{\mathbf{X}}(x_1, x_2, \dots, x_n) := P(\{\omega \in \Omega : X_1(\omega) = x_1, X_2(\omega) = x_2, \dots, X_n(\omega) = x_n\}) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n).$$

The **marginal probability mass function** of $X_i \in \mathbf{X}$ is given by,

$$p_{X_i}(k) = \sum_{(x_1, x_2, \dots, x_n) \in \mathbb{R}^n} \cdots \sum p_{\mathbf{X}}(x_1, x_2, \dots, x_{i-1}, k, x_{i+1}, \dots, x_n).$$

It can be obtained that for any sufficiently "nice" set $A \in \mathbb{R}^n$,

$$P(\mathbf{X} \in A) = \sum_{(x_1, x_2, \dots, x_n) \in A} \cdots \sum P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n).$$

Definition 12.3.2 (Independence of DRVs). Given some indexing set $\mathcal{I} \in \mathbb{R}$ a family of DRVs, $\{X_i : i \in \mathcal{I}\}$ with joint pmf $p_{\mathbf{X}}$, is **independent** iff for all finite $\mathcal{J} \in \mathcal{I}$:

$$p_{\mathbf{X}}(\mathbf{x}) = \prod_{i=1}^n p_{X_i}(x_i) \quad \text{for all } \mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n.$$

12.4 Multivariate CRVs

Definition 12.4.1 (Continuous random vector). The random vector $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a **continuous random vector** iff,

$$F_{\mathbf{X}}(\mathbf{x}) = \int \cdots \int_B f_{\mathbf{X}}(\mathbf{x}) dx_1 dx_2 \cdots dx_n \quad \text{with } B = (\infty, x_1] \times (\infty, x_2] \times \cdots \times (\infty, x_n], \quad \text{for all } \mathbf{x} \in \mathbb{R}^n;$$

for some $f_{\mathbf{X}} : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying: $f_{\mathbf{X}}(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$ and $\int \cdots \int_{\mathbb{R}^n} f_{\mathbf{X}}(\mathbf{x}) dx_1 dx_2 \cdots dx_n = 1$.

Note that $f_{\mathbf{X}}(\mathbf{x}) = \frac{\partial^n}{\partial x_1 \partial x_2 \cdots \partial x_n} F_{\mathbf{X}}(\mathbf{x})$ and $P(\mathbf{X} \in A) = \int \cdots \int_A f_{\mathbf{X}}(\mathbf{x}) d^n \mathbf{x}$.

Definition 12.4.2 (Independence of CRVs). Given some indexing set $\mathcal{I} \in \mathbb{R}$ a family of CRVs, $\{X_i : i \in \mathcal{I}\}$ with joint pdf $f_{\mathbf{X}}$, is **independent** iff for all finite $\mathcal{J} \in \mathcal{I}$:

$$f_{\mathbf{X}}(\mathbf{x}) = \prod_{i=1}^n f_{X_i}(x_i) \quad \text{for all } \mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n.$$

12.5 Transformations of random vector

Definition 12.5.1 (Transformation). We are going to **transform** the random vector \mathbf{X} with joint pdf $f_{\mathbf{X}}$ to $\mathbf{U} = (u_1(\mathbf{X}), u_2(\mathbf{X}), \dots, u_n(\mathbf{X}))$ with $u_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$ for all $i \in [1, n]$. Firstly, define $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ by $T(\mathbf{x}) = (u_1(\mathbf{x}), u_2(\mathbf{x}), \dots, u_n(\mathbf{x}))$ and assume T is bijective on $D = \{\mathbf{x} \in \mathbb{R}^n : f_{\mathbf{X}}(\mathbf{x}) > 0\}$ with range $S \subseteq \mathbb{R}^n$. Secondly, have the Jacobian determinant of $T^{-1} : S \rightarrow D$, $J(\mathbf{u}) = \det([a_{ij}]_{m \times n})$ with $a_{ij} = \frac{\partial x_i}{\partial u_j}$. Finally, define:

$$f_{\mathbf{U}}(\mathbf{u}) := \begin{cases} f_{\mathbf{X}}(T^{-1}(\mathbf{u}))|J(\mathbf{u})| & \text{if } \mathbf{u} \in S \\ 0 & \text{otherwise} \end{cases}.$$

12.6 Multivariate LOTUS

Theorem 12.6.1 (Discrete multivariate LOTUS). If X_1, X_2, \dots, X_n are DRVs on (Ω, \mathcal{F}, P) and form the random vector \mathbf{X} with $g : \mathbb{R}^n \rightarrow \mathbb{R}$, then $Y = g(\mathbf{X})$ is a DRV on (Ω, \mathcal{F}, P) with expectation,

$$E(g(\mathbf{X})) = \sum_{x_i \in \text{im } X_i} \cdots \sum g(\mathbf{x}) P(\mathbf{X} = \mathbf{x}) \quad \text{for all } \mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n.$$

Theorem 12.6.2 (Continuous multivariate LOTUS). If X_1, X_2, \dots, X_n are DRVs on (Ω, \mathcal{F}, P) and form the random vector \mathbf{X} with $h : \mathbb{R}^n \rightarrow \mathbb{R}$ we have,

$$E(h(\mathbf{X})) = \int \cdots \int_{\mathbb{R}^n} g(\mathbf{x}) f_{\mathbf{X}}(\mathbf{x}) dx_1 dx_2 \cdots dx_n \quad \text{for all } \mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n.$$

12.7 Covariance

Definition 12.7.1 (Covariance). Given two random variable X and Y on the same probability space with expectations μ_X and μ_Y respectively. The **covariance** of X and Y is defined as,

$$\text{Cov}(X, Y) := E[(X - \mu_X)(Y - \mu_Y)] \quad \text{assuming both expectation take finite values.}$$

Definition 12.7.2 (Correlation). Given the same X and Y the **correlation** of X and Y is defined as,

$$\text{Cor}(X, Y) := \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}.$$

Theorem 12.7.3. For jointly discrete/continuous RVs with finite expectations the following hold:

1. when $X = Y$, $\text{Cov}(X, Y) = E[(X - \mu_X)^2] = \text{Var}(X)$,
2. $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$,
3. when X and Y are independent, $E(XY) = E(X)E(Y)$,
4. when variances are also finite, $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$.

13 Generating functions

13.1 Probability generating functions

Definition 13.1.1 (Probability generating functions). Given a DRV X with $\text{im}(X) \subseteq \mathbb{N}_0$, denote,

$$\mathcal{S}_X := \left\{ s \in \mathbb{R} : \sum_{x=0}^{\infty} |s|^x P(X = x) < \infty \right\},$$

and define the **probability generating function (pgf)** of X as the function $G_X : \mathcal{S}_X \rightarrow \mathbb{R}$ given by,

$$G_X(s) := E(s^X) = \sum_{x=0}^{\infty} s^x P(X = x),$$

noting that the pgf is well defined for $|s| < 1$ and $G_X(0) = P(X = 0)$ and $G_X(1) = 1$.

Theorem 13.1.2 (Uniqueness of pgfs). Given two DRVs X and Y with pgfs G_X and G_Y respectively,

$$G_X(s) = G_Y(s) \quad \text{for all } s \in \mathcal{S}_X \cap \mathcal{S}_Y \iff p_X(x) = p_Y(x) \quad \text{for all } x \in \mathbb{N}_0.$$

Theorem 13.1.3. Let X, Y be independent DRVs with $\text{im } X, \text{im } Y \in \mathbb{N}_0$, then

$$G_{X+Y}(s) = G_X(s)G_Y(s) \quad \text{for all } s \in \mathcal{S}_X \cap \mathcal{S}_Y.$$

Theorem 13.1.4 (Pgfs of sum of independent DRVs). Given a collection of n independent DRVs X_1, X_2, \dots, X_n ,

$$G_{\sum_{i=1}^n X_i}(s) = \prod_{i=1}^n G_{X_i}(s) \quad \text{for all } s \in \bigcap_{i=1}^n \mathcal{S}_{X_i}.$$

Theorem 13.1.5 (Moments). Given a DRV X with $\text{im } X \subseteq \mathbb{N}_0$, the k th derivative of G_X , for $k \in \mathbb{N}$ is given by,

$$\left. \frac{d^k}{ds^k} G_X(s) \right|_{s=1} = G_X^{(k)}(1) = E[X(X-1)\dots(X-k+1)].$$

13.2 Common pgfs

Example 13.2.1 (Bernoulli distribution). Let $X \sim \text{Bern}(p)$, then $G_X(s) = 1 - p + sp$ for all $s \in \mathbb{R}$.

Example 13.2.2 (Binomial distribution). Let $X \sim \text{Bin}(n, p)$, then $G_X(s) = (1 - p + sp)^n$ for all $s \in \mathbb{R}$.

Example 13.2.3 (Poisson distribution). Let $X \sim \text{Poi}(\lambda)$, then $G_X(s) = \exp(\lambda(s - 1))$ for all $s \in \mathbb{R}$.

13.3 Moment generating functions

Definition 13.3.1 (Moment generating function). Let X be a RV, its **moment generating function** (mgf) is defined as,

$$M_X(t) = E(e^{tX}),$$

assuming the expectation exists in some neighbourhood of zero.

Remark 13.3.2. If X is a RV with a mgf, $M_X(t) = E(e^{tX}) = G_X(e^t)$

Theorem 13.3.3. If X is a RV with a mgf, the k th moment of X is $E(X^k) = M_X^{(k)}(0)$.

Theorem 13.3.4. If X_1, X_2, \dots, X_n are a family of independent RVs on the same probability space with mgfs $M_{X_1}, M_{X_2}, \dots, M_{X_n}$ respectively, we have,

$$M_{\sum_{i=1}^n X_i}(t) = \prod_{i=1}^n M_{X_i}(t).$$

Theorem 13.3.5 (Characterisation by mgf). If the RVs X, Y have existent mgfs M_X, M_Y respectively such that $M_X(t) = M_Y(t)$ for all t in some neighbourhood of 0, we have,

$$F_X(u) = F_Y(u) \quad \text{for all } u.$$

14 Conditional distribution and expectation

14.1 Discrete: Conditional expectation and total expectation

Definition 14.1.1 (Condition distribution and expectation of a DRV). Given a DRV Y on (Ω, \mathcal{F}, P) and some event $B \in \mathcal{F}$ with $P(B) > 0$, the **conditional distribution** of Y given B is defined as,

$$P(Y = y|B) := \frac{P(\{Y = y\} \cap B)}{P(B)} \quad \text{for } y \in \mathbb{R};$$

with the **conditional expectation** of Y given B defined as,

$$E(Y|B) := \sum_{i \in \text{im } Y} e P(Y = y|B).$$

Theorem 14.1.2 (Discrete law of total expectation). Given a DRV Y on (Ω, \mathcal{F}, P) and some partition $\{B_i : i \in \mathcal{I}\}$ of Ω with $P(B_i) > 0$ for all $i \in \mathcal{I}$ we have,

$$E(Y) = \sum_{i \in \mathcal{I}} E(Y|B_i)P(B_i).$$

14.2 Conditioning on a DRV

Definition 14.2.1 (Conditional probability mass function). Given two joint DRVs (X, Y) , the **conditional probability mass function** of Y given $X = x$ is given by,

$$p_{Y|X}(y|x) := \frac{p_{X,Y}(x, y)}{p_X(x)} \quad \text{for } y \in \mathbb{R}.$$

Theorem 14.2.2 (Conditional expectation). Given two joint DRVs (X, Y) , the **conditional expectation** of Y given $X = x$ is given by,

$$E(Y|X = x) = \sum_{y \in \text{im } Y} y p_{Y|X}(y|x).$$

Independence, LOTUS and a Bayes' rule formulation all follow naturally from this as they do for the non-distribution settings.

14.3 Continuous: Conditional density, distribution and expectation

Definition 14.3.1 (Conditional distribution and conditional density). For two joint CRVs (X, Y) the **conditional density** of Y given $X = x$ is defined as,

$$f_{Y|X}(y|x) := \frac{f_{X,Y}(x, y)}{f_X(x)} \quad \text{for all } y, x \in \mathbb{R} \text{ with } f_X(x) > 0;$$

with the corresponding **conditional distribution** of Y given $X = x$ defined as,

$$F_{Y|X=x}(y|x) := \frac{1}{f_X(x)} \int_{-\infty}^y f_{X,Y}(x, t) dt \quad \text{for all } y, x \in \mathbb{R} \text{ with } f_X(x) > 0.$$

Where, once again, familiar formulations for independence and Bayes' rule can be easily derived.

Definition 14.3.2 (Conditional expectation). Given two joint CRVs (X, Y) , the **conditional expectation** of Y given $X = x$ is defined as,

$$E(Y|X = x) := \int_{-\infty}^{\infty} y f_{Y|X}(y|x) dy \quad \text{provided } f_X(x) > 0.$$

Theorem 14.3.3 (Continuous law of total expectation). Given two joint CRVs (X, Y) with $E(|Y|) < \infty$, we have,

$$E(Y) = \int_{\{x: f_X(x) > 0\}} E(Y|X = x) f_X(x) dx.$$