

Measure Theory

Let X be a set. Let $\mathcal{A} \subseteq \mathcal{P}(X)$ be a subset of the power set.

We call \mathcal{A} a σ -algebra if the following is true:

and we call an element of \mathcal{A} an \mathcal{A} -measurable set.

- $\emptyset, X \in \mathcal{A}$;
- $A \in \mathcal{A} \Rightarrow A^c := X \setminus A \in \mathcal{A}$;
- $A_i \in \mathcal{A}, i \in \mathbb{N} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{A}$.
countably many



Examples: (1) $\mathcal{A} = \{\emptyset, X\}$ (smallest σ -algebra).

(2) $\mathcal{A} = \mathcal{P}(X)$ (largest σ -algebra).

arbitrary index set

It is "easy to show" the following. Let \mathcal{A}_i be σ -algebras on X , where $i \in I$. Then $\bigcap_{i \in I} \mathcal{A}_i$ is a σ -algebra.

For $M \subseteq \mathcal{P}(X)$, there is a smallest σ -algebra containing M : $\bigcap_{M \subseteq \mathcal{A}} \mathcal{A}$ where \mathcal{A} are σ -algebras.

We define $\sigma(M) := \bigcap_{\substack{M \subseteq \mathcal{A} \\ \mathcal{A} \text{ a } \sigma\text{-algebra}}} \mathcal{A}$ often called the σ -algebra generated by M .

Let $X = \{a, b, c, d\}$ and $M = \{\{a\}, \{b\}\}$. Then $\sigma(M) = \{\emptyset, X, \{a\}, \{b\}, \{a, b\}, \{b, c, d\}, \{a, c, d\}, \{c, d\}\}$.

Let X be a topological space. We define the Borel σ -algebra $\mathcal{B}(X)$ to be the σ -algebra generated by open sets, i.e. $\mathcal{B}(X) := \sigma(\mathcal{T})$ on some space (X, \mathcal{T}) .

A set X and σ -algebra \mathcal{A} form a measurable space (X, \mathcal{A}) .

A map $\mu: \mathcal{A} \rightarrow [0, \infty]$ is called a measure if the following is fulfilled:

$[0, \infty] \cup \{\infty\}$
volume

- $\mu(\emptyset) = 0$ (empty set has no volume);
- $\mu\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n \mu(A_i)$ with $A_i \cap A_j = \emptyset, \forall A_i \in \mathcal{A}$,
(additivity of μ).



but we also want the notion of volume approximations. Consider the picture

We construct the sequence $\{A_1, A_2, \dots\}$ and can still take unions.

Therefore $\mu\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i)$ with $A_i \cap A_j = \emptyset, \forall A_i \in \mathcal{A}$,
(σ -additivity).

Then we call (X, \mathcal{A}, μ) a measure space.

Let X be a set and \mathcal{A} be a given σ -algebra on X .

(a) We can define the counting measure for $A \in \mathcal{A}$ by $\mu(A) := \begin{cases} |A| & \text{if } A \text{ has finitely many elements;} \\ \infty & \text{otherwise.} \end{cases}$

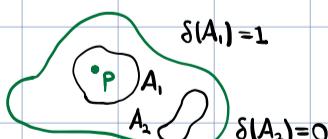
We then have the rules

$$x + \infty := \infty$$

$$x \cdot \infty := \infty$$

$$0 \cdot \infty := 0 \quad \text{only applicable in measure theory.}$$

(b) We can define the Dirac measure for $p \in X$ by $\delta_p(A) := \begin{cases} 1, & \text{if } p \in A; \\ 0, & \text{if } p \notin A. \end{cases}$



(c) We want a measure on $X = \mathbb{R}^n$ to have the following properties:

- the lebesgue measure
- $\mu([0, 1]^n) = 1$ (the unit cube has "volume" 1);
 - $\mu(x + A) = \mu(A)$ for all $x \in \mathbb{R}^n$ ("volume" measurement is translation invariant).

We shall define the concept of a probability model formally.

A probability model is a set Ω called the sample space together with a collection of subsets \mathcal{A} of Ω called events, and a real-valued function $P: \mathcal{A} \rightarrow [0, 1] \subset \mathbb{R}$ called a probability measure.

Probability model: $\begin{cases} \Omega \text{ sample space} \\ \mathcal{A} \text{ events (collection of subsets of } \Omega) \\ P \text{ probability measure} \end{cases}$

The collection \mathcal{A} of events satisfies the following:

- \mathcal{A} |
 - $\Omega \in \mathcal{A}$;
 - If $A \in \mathcal{A}$, then $\Omega \setminus A \in \mathcal{A}$;
 - If $\{A_i\}_{i \in \mathbb{N}} \subseteq \mathcal{A}$ is a (possibly infinite) sequence of events in \mathcal{A} , then $\bigcup_{i \in \mathbb{N}} A_i \in \mathcal{A}$.

The probability measure P satisfies the following:

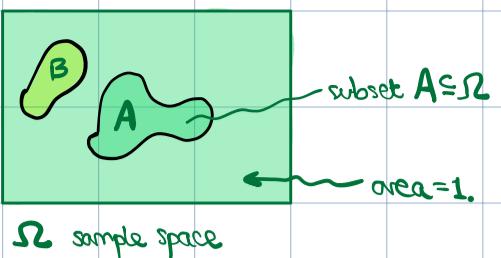
- P |
 - $P(A) \geq 0$ for all $A \in \mathcal{A}$;
 - $P(\Omega) = 1$;
 - The probability measure is countably additive, i.e.

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i) \text{ for any infinite sequence of pairwise disjoint } A_i.$$

Exercise: Check whether the following \mathcal{A} satisfy the axioms above.

- $\mathcal{A} = P(\Omega)$ the power set. Yes, trivially.
 - $\emptyset \in \mathcal{A} \Rightarrow \Omega \in \mathcal{A}$. ✓
 - $A \in \mathcal{A} \Rightarrow \Omega \setminus A \in \mathcal{A}$ by definition (either A finite $\Rightarrow \Omega \setminus A \in \mathcal{A}$, or A infinite $\Rightarrow A = \Omega \setminus B$ for some $B \in \mathcal{A} \Rightarrow \Omega \setminus A = B \in \mathcal{A}$). ✓
 - Not satisfied. Take $A_i = \{i\}$ for $i \in \mathbb{N}$ so $\{A_i\}_{i \in \mathbb{N}} \subseteq \mathcal{A}$, but $\bigcup_{i \in \mathbb{N}} A_i \notin \mathcal{A}$.
- $\Omega = \mathbb{R}$ and $\mathcal{A} = (\text{all finite } S \subset \mathbb{R} \text{ and complements})$.
 - 1, 2 See above. ✓
 - 3 Should be true since countable unions of countable sets are countable. ✓
- $\mathcal{A} = \{\emptyset, \Omega\}$ trivially.
- $\Omega = \{1, 2, 3, \dots\}$ and \mathcal{A} is all the sets $A \subset \Omega$ such that $\lim_{n \rightarrow \infty} \frac{|A \cap \{1, 2, \dots, n\}|}{n}$ exists.

From measure theory, probability measures are measures with total mass = 1.



$P: \mathcal{A} \rightarrow [0, 1]$

collection of subsets of Ω

$$\bullet P(\Omega) = 1;$$

$$\bullet P(A) \in [0, 1] \text{ for } A \in \mathcal{A};$$

$$\bullet P(A \cup B) = P(A) + P(B) \text{ whenever } A \cap B = \emptyset;$$

(disjoint A and B)

$$\bullet P\left(\bigcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} P(A_j).$$

Let Ω be a set. A collection of subsets $\mathcal{A} \subseteq P(\Omega)$ is called a σ -algebra if:

- $\emptyset, \Omega \in \mathcal{A}$;
- If $A \in \mathcal{A}$ then $\Omega \setminus A \in \mathcal{A}$;
- If $A_1, A_2, \dots \in \mathcal{A}$ then $\bigcup_{j=1}^{\infty} A_j \in \mathcal{A}$.

} the elements of a σ -algebra are called events in probability theory.

Example: Take a die, roll it once. Then $\Omega = \{1, 2, 3, 4, 5, 6\}$ (sample space).

Say $\mathcal{A} = P(\Omega)$. Since every side has the same probability, we have $P: \mathcal{A} \rightarrow [0, 1]$ given by $P(A) = \frac{\#A}{\#\Omega}$.

(e.g. $P(\{2\}) = \frac{1}{6}$, or $P(\{2, 4, 6\}) = \frac{3}{6} = \frac{1}{2}$).

Exercise: Prove that $P(\Omega \setminus A) = 1 - P(A)$ for $A \in \mathcal{A}$.

σ -algebra

Probability problems can be categorized into discrete, (absolutely) continuous, and mixed.

Discrete	Continuous
sample space Ω finite or countable	sample space $\Omega \subseteq \mathbb{R}^n$ uncountable with $\Omega \in \mathcal{B}(\mathbb{R}^n)$ a Borel set
σ -algebra: could take $\mathcal{A} = P(\Omega)$ (standard)	σ -algebra: could take $\mathcal{A} = \mathcal{B}(\Omega)$ the Borel σ -algebra of Ω .
$P: \mathcal{A} \rightarrow [0, 1]$ determined by $P(\{\omega\})$ for all $\omega \in \Omega$ (singletons)	$P: \mathcal{A} \rightarrow [0, 1]$ determined by a probability density function $f: \Omega \rightarrow \mathbb{R}$ with $f(x) \geq 0$ measurable $\int_{\Omega} f(x) dx = 1$

\Rightarrow define $P(A) = \int_A f(x) dx$.

Examples: Roll a biased die, so $\Omega = \{1, 2, 3, 4, 5, 6\}$

and set $P_1 = P_2 = P_3 = P_4 = P_5 = \frac{1}{10}$ and $P_6 = \frac{1}{2}$.

$$\text{Then } P(\{1, 2, 3, 4, 5\}) = \sum_{i=1}^5 P_i = 5 \cdot \frac{1}{10} = \frac{1}{2}.$$

Consider $\Omega = [0, 1]$ and $f: \Omega \rightarrow \mathbb{R}$ given by $f(x) = \frac{1}{2}$.

$$\text{Then } \int_{\Omega} f(x) dx = \int_{[0, 1]} \frac{1}{2} dx = 1.$$

Here in general, we have

$$\int_{\Omega} f(x) dx = \frac{1}{2} \int_{\Omega} 1 dx = \frac{1}{2} \cdot (\text{Lebesgue measure } (\Omega)).$$

Consider tossing a coin with possible outcomes H, T. Then $P(H) = \frac{a}{a+b} \in \mathbb{Q} \cap [0, 1]$ for some $a, b \in \mathbb{N}$.

(a fair coin has $a=b$).

Now consider drawing a ball from a box, where balls are either labelled H or T.

$$\text{Then } P(H) = \frac{a}{a+b} \in \mathbb{Q} \cap [0, 1] \text{ again.}$$

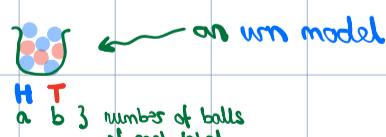
assume

In both cases above, we have $\Omega = \{H, T\}$ and $P(\{H\}) = \frac{a}{a+b}$ and $P(\{T\}) = \frac{b}{a+b}$. $\left. \begin{array}{l} P_H \text{ and } P_T \text{ define } P \\ \text{in this discrete case} \end{array} \right\}$

Both represent binomial distributions:

- n tosses of the same coin, counting the heads;
- draw n balls with replacement, counting the heads;
- generally sample size n , unordered, with replacement counting problem.

→ without order



an urn model

H T

a b } number of balls

of each label

probability mass functions

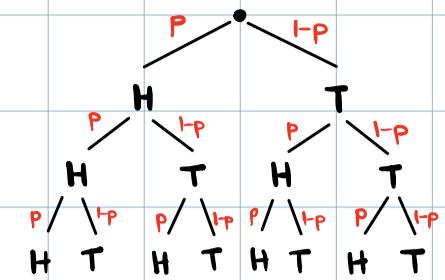
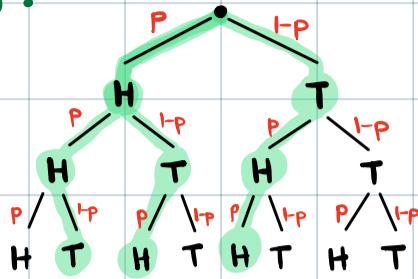
P

Then we would have $\Omega = \{0, 1, 2, \dots, n\}$ and $P(\{k\}) = \binom{n}{k} p^k (1-p)^{n-k}$ with parameters (n, p) .

Consider arbitrary $p \in [0, 1]$ and $n=3$ (shown on the decision tree on the right).

Then the probability $P(\{k\})$ of choosing k heads H is given by

$$P(\{k\}) = \binom{n}{k} p^k (1-p)^{n-k} = \binom{3}{k} p^k (1-p)^{3-k}.$$



Recall that a **probability space** is given by (Ω, \mathcal{A}, P) .

$\mathcal{A} \subseteq P(\Omega)$
σ-algebra
sample space

"events"
probability measure

$P: \mathcal{A} \rightarrow [0, 1]$

and why should we?

Why? Consider the procedure: 1 Throw a die; 2 Throw a point into $[1, 1] \subset \mathbb{R}$. A possible outcome is $(3, \frac{1}{4})$. What is its probability?

First probability space $P_1: (\Omega_1, \mathcal{A}_1, P_1)$ (discrete)
 $\Omega_1 = \{1, 2, 3, 4, 5, 6\}$
 $P_1(\Omega) = \sum_{k=1}^6 \frac{1}{6}$

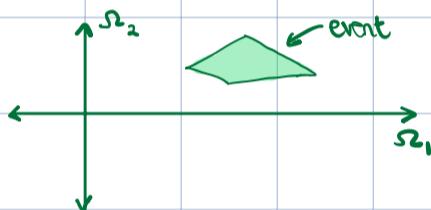
We have a new product probability space:

$(\Omega_1 \times \Omega_2, \sigma(\mathcal{A}_1 \times \mathcal{A}_2), P)$

\uparrow Cartesian product
product σ-algebra

\downarrow product measure

Second probability space $P_2: (\Omega_2, \mathcal{A}_2, P_2)$ (continuous)
 $\Omega_2 = [1, 1]$
 $P_2(\Omega) = \int \frac{1}{2} dx$



Consider a collection $(\Omega_n, \mathcal{A}_n, P_n)$ of probability spaces with $n \in \mathbb{N}_{>0}$. Take the case $n=2$.

The **product (probability) space** is given by $(\Omega_1 \times \Omega_2, \sigma(\mathcal{A}_1 \times \mathcal{A}_2), P)$ where P is a **product measure**

satisfying for $A_1 \in \mathcal{A}_1$ and $A_2 \in \mathcal{A}_2$ we have $P(A_1 \times A_2) = P_1(A_1) \cdot P_2(A_2)$.

\uparrow lives in product σ-algebra

Consider the event of $((\text{rolling a } 2 \text{ or } 3), (\text{hitting } [1, 0]))$.

$$\text{Then } P(\{2, 3\} \times [1, 0]) = P_1(\{2, 3\}) \cdot P_2([1, 0]) = \sum_{k \in \{2, 3\}} \frac{1}{6} \cdot \int_{[1, 0]} \frac{1}{2} dx = \frac{2}{6} \cdot \frac{1}{2} = \frac{1}{6}.$$

Let us provide a general definition. Consider a collection $(\Omega_n, \mathcal{A}_n, P_n)$ of probability spaces with $n \in \mathbb{N}_{>0}$.

The **product (probability) space** P given by (Ω, \mathcal{A}, P) is defined by:

- $\Omega = \Omega_1 \times \Omega_2 \times \dots = \prod_{j \in \mathbb{N}_{>0}} \Omega_j$ having elements sequences $(\omega_1, \omega_2, \omega_3, \dots)$
 $\omega_1, \omega_2, \omega_3 \in \Omega_j$

- $\mathcal{A} = \sigma(\text{"cylinder sets"})$ the product σ-algebra.

\uparrow take the smallest σ-algebra
 $A_1 \times \Omega_2 \times \Omega_3 \times \dots$
 $\Omega_1 \times A_2 \times \Omega_3 \times \dots$ containing all these subsets

- P given by $P(A_1 \times \dots \times A_m \times \Omega_{m+1} \times \Omega_{m+2} \times \dots) = \prod_{j=1}^m P_j(A_j)$ the **product measure**.

Example: Throw a dice infinitely many times. Since each probability space is the same, denote each as $(\Omega_0, \mathcal{A}_0, P_0)$.

Then the produce space is (Ω, \mathcal{A}, P) product measure
 $\Omega = \Omega_0 \times \Omega_1 \times \dots$ product σ-algebra

by the definition above.

$$\Omega_0 = \{1, 2, 3, 4, 5, 6\}$$

$$P_0(A) = \sum_{k \in A} \frac{1}{6}$$

Consider the event $A \in \mathcal{A}$ given by A : "at the 100th throw, we roll the first 6."

Thus $A = \underbrace{\{6\}^c \times \dots \times \{6\}^c}_{99 \text{ times}} \times \{6\} \times \underbrace{\Omega_0 \times \Omega_0 \times \dots}_{\substack{\text{we don't care} \\ \text{what happens after}}} \quad \text{where } \{6\}^c = \Omega_0 \setminus \{6\} \text{ the complement.}$

$$\Rightarrow P(A) = P_0(\{6\}^c) \dots P_0(\{6\}^c) \cdot P_0(\{6\}) = P_0(\{6\}^c)^{99} \cdot P_0(\{6\}) = \left(\frac{5}{6}\right)^{99} \cdot \frac{1}{6}.$$

The Hypergeometric Distribution

Keep 3 things in mind: size $n \in \mathbb{N}$, sample unordered, no replacement. The urn model provides a nice intuition again.

Consider drawing n balls from a jar containing ball colours C (some finite set). 

One possible outcome for $n=3$ may be  so we essentially have some function $f: C \rightarrow \mathbb{N}_{\geq 0}$ 

Then $\Omega = \text{set of all functions of that form} = \{(k_c)_{c \in C} \in \mathbb{N}_{\geq 0}^C \mid \sum_{c \in C} k_c = n\}$. 

(Indexing colours can allow us to write outcomes as tuples)
adding to $n=3$

As all discrete cases, we can take our σ -algebra $\mathcal{A} = \mathcal{P}(\Omega)$.

For our example, we have $\Omega = \{(k_0, k_1, k_2) \in \mathbb{N}_{\geq 0}^3 \mid k_0 + k_1 + k_2 = n\}$.

To know the probability of an outcome, we must know the initial counts of all colours. We introduce $N_c = \text{number of balls of colour } c \text{ in the urn}$

with $N = \sum_{c \in C} N_c$ total count.

Then we have $P((k_0, k_1, k_2)) = \frac{\binom{N_0}{k_0} \cdot \binom{N_1}{k_1} \cdot \binom{N_2}{k_2}}{\binom{N}{n}}$. 

We can use probability mass functions in this case.
(discrete)

The (multivariate) hypergeometric distribution is described by the probability mass function

$$P(\{(k_c)_{c \in C}\}) = \frac{\prod_{c \in C} \binom{N_c}{k_c}}{\binom{N}{n}} \quad \text{where } n \text{ are chosen.}$$

Hypergeometric distribution for 2 colours:

$$C = \{0, 1\} \quad \text{and} \quad N_0 + N_1 = N.$$

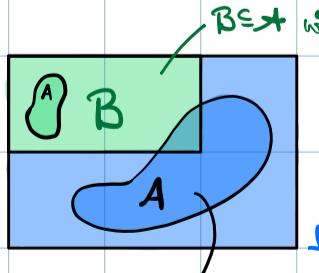
Since only two colours exist, we can simply count for one colour to get the other.

Count the 1s: $\Omega = \{0, 1, \dots, n\}$ since when n balls are picked, these are our possibilities.

$$\text{Then we can define } P: \mathcal{P}(\Omega) \rightarrow [0, 1] \text{ by } P(\{k\}) = \frac{\binom{N_1}{k} \binom{N_0}{n-k}}{\binom{N}{n}}.$$

Conditional probability is a special type of probability measure defined relative to another probability measure.

Let (Ω, \mathcal{A}, P) define a probability space, and let $B \subseteq \mathcal{A}$ be a subset of our σ -algebra (of events) with $P(B) \neq 0$.



when A not contained in B

We can define new probability space(s)

$$(B, \tilde{\mathcal{A}}, \tilde{P}) \quad \tilde{P}(A) = \frac{P(A)}{P(B)} \quad \text{normalized but what if we don't want to change our sample space } \Omega \text{ and } \sigma\text{-algebra } \mathcal{A}?$$

$$\rightsquigarrow (\Omega, \mathcal{A}, P_B) \quad \text{where } P_B(A) = \frac{P(A \cap B)}{P(B)}.$$

$A \in \mathcal{A}$ may not be contained in B !

Now we define conditional probability.

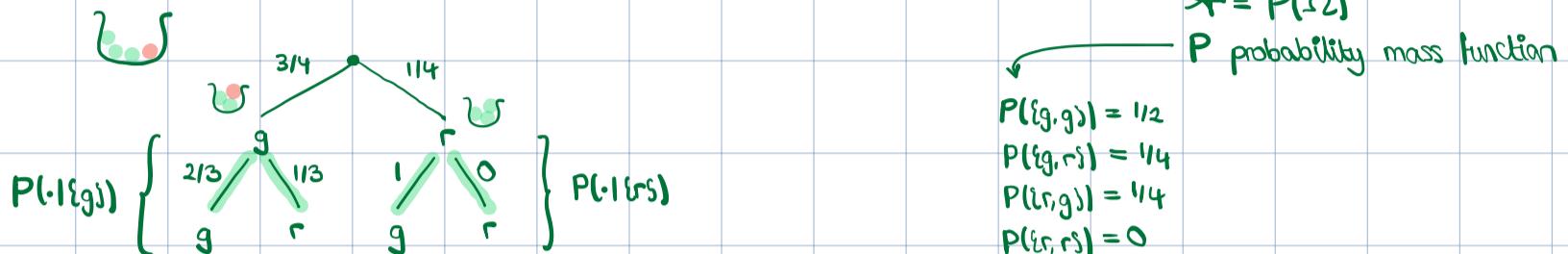
Let (Ω, \mathcal{A}, P) define a probability space, and $B \in \mathcal{A}$ an element with $P(B) \neq 0$.

We define $P_B(A) = P(A|B) := \frac{P(A \cap B)}{P(B)}$ the conditional probability of A under B .

The conditional probability (measure) given B is given by $P(\cdot|B) : \mathcal{A} \rightarrow [0, 1]$. Notice $P(B|B) = 1$ by definition.

If $P(B) = 0$ we define $P(A|B) = 0$. Here, we don't have a well-defined probability measure.

Example: Urn model: ordered, without replacement. Define colour set $C := \{r, g\}$ so $\Omega = C \times C$ when 2 are picked.



Suppose $B = \text{"first ball green"}$. Then $B = \{(g, g), (g, r)\}$ and if $A = \{(g, r)\} \in \mathcal{A}$ then

$$P(A|B) = P(\{(g, r)\}|B) = \frac{P(\{(g, r)\} \cap B)}{P(B)} = \frac{P(\{(g, r)\})}{P(B)} = \frac{\frac{1}{4}}{\frac{3}{4}} = \frac{1}{3}.$$

Counting Principles (Combinatorics)

Consider a bag with n elements, where we select k of these elements at random. We want to figure out:

# Possibilities	We replace	We don't replace
order matters	$\frac{n(n-1)\dots(n-(k-1))}{(n-k)!} = \frac{n!}{(n-k)!}$	n^k
order doesn't matter	$\frac{n!}{k!(n-k)!} = \binom{n}{k}$	$\frac{(k+n-1)!}{k!(n-1)!}$

non-trivial case: we have n things, pick k , order doesn't matter, we replace.

$\boxed{1 \quad 2 \dots n}$ with k tennis balls.
n buckets

Q: How many ways to put k tennis balls in n buckets?

= k balls and $n-1$ barriers: $\underbrace{\circ \circ | \circ \circ \dots \circ \circ}_{k \text{ balls}}$ = $aabbaba\dots aba$

$$= \frac{(k+n-1)!}{k!(n-1)!}$$

permute 'a's in $k!$ ways
permute 'b's in $(n-1)!$ ways
arrange $(k+n-1)$ objects,
overcounting by $k!(n-1)!$

The Birthday Problem: We have k people, and we want to know the probability of a birthday match.

Assume 365 days and a uniform distribution.

$\boxed{27 \quad 310 \quad 89 \dots 64}$ so $|\Omega| = 365^k$

$\Rightarrow \Omega = \{(a_1, a_2, \dots, a_k) \mid a_i \in \mathbb{Z}/365\mathbb{Z}\}$ and $\mathcal{A} = P(\Omega)$ because discrete problem.

Say $A \in \mathcal{A}$ represents events with a match. Then $P(A) = 1 - P(\Omega \setminus A) = 1 - P(A^c)$.

$$A^c: \underbrace{\bullet \quad \bullet \quad \dots \quad \bullet}_{365 \quad 364 \quad 365-(k-1) \quad \text{choices}} \xrightarrow{\text{entry}} |A^c| = \frac{365!}{(365-k)!} \text{ so } P(A^c) = \frac{|A^c|}{|\Omega|} = \frac{365!}{365^k} \frac{365^{-k}}{(365-k)!}$$

$$\Rightarrow P(A) = 1 - \frac{365!}{(365-k)!} 365^{-k}$$

Theorem (Bayes): Let (Ω, \mathcal{A}, P) define a probability space with $A, B \in \mathcal{A}$ of nonzero probability.

Then $P(A|B) P(B) = P(B|A) P(A)$.

Proof: Simple. We have the conditional probability measure defined above, so

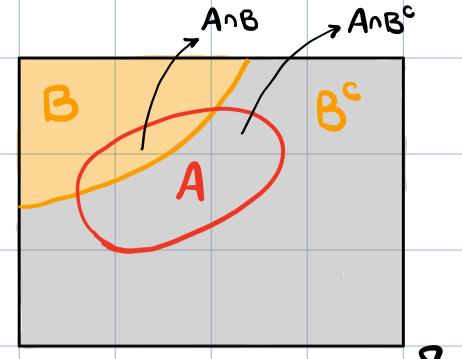
$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{and} \quad P(B|A) = \frac{P(B \cap A)}{P(A)} = \frac{P(A \cap B)}{P(A)} \Rightarrow P(B|A) P(A) = P(A|B) P(B). \quad \square$$

" $P(A \cap B)$ "

Law of Total Probability

Consider a probability space (Ω, \mathcal{A}, P) . Say $A \in \mathcal{A}$ and $B \in \mathcal{A}$.

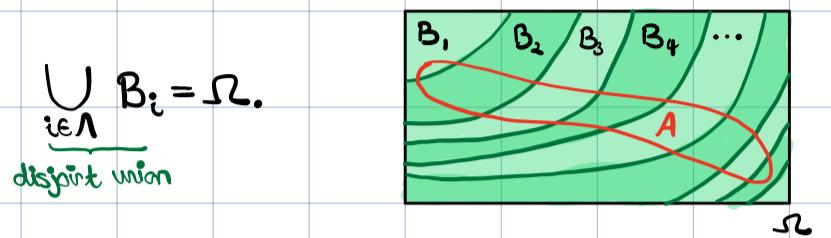
$$\begin{aligned} \text{We have } P(A) &= P((A \cap B) \cup (A \cap B^c)) \quad \leftarrow \text{disjoint because } \\ &= P(A \cap B) + P(A \cap B^c) \\ &= P(A|B)P(B) + P(A|B^c)P(B^c). \end{aligned}$$



$$\underline{B \cup B^c = \Omega}$$

disjoint union

Consider the case with countably many $B_i \in \mathcal{A}$ for $i \in \Lambda \subseteq \mathbb{N}$ with $\bigcup_{i \in \Lambda} B_i = \Omega$.

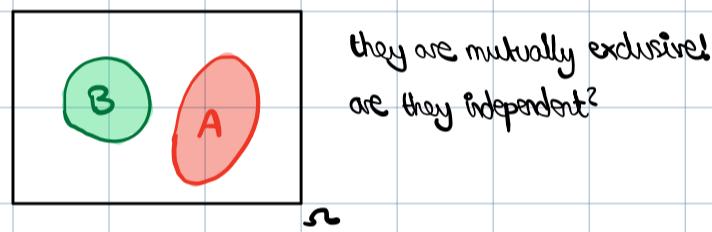
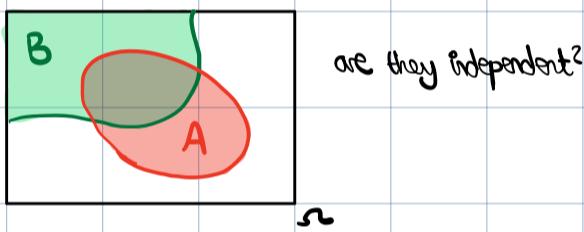


$$\text{Then } P(A) = P\left(\bigcup_{i \in \Lambda} (A \cap B_i)\right) = \sum_{i \in \Lambda} P(A \cap B_i) = \sum_{i \in \Lambda} P(A|B_i) P(B_i)$$

disjoint union

Example: Monty Hall Problem.

Consider a space (Ω, \mathcal{A}, P) and events $A, B \in \mathcal{A}$. We want to deduce something about **Independence**.



$$\text{We want } P(A|B) = P(A) \text{ and } P(B|A) = P(B).$$

$$\text{Notice then } P(A \cap B) = P(A)P(B).$$

$$\left. \begin{array}{l} P(A) = \frac{1}{2} \\ P(B) = \frac{1}{2} \end{array} \right\} \Rightarrow P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{\left(\frac{1}{4}\right)}{\left(\frac{1}{2}\right)} = \frac{1}{2},$$

and $P(B|A) = \frac{1}{2}$.

⇒ independent!

Let (Ω, \mathcal{A}, P) be a probability space. Two events $A, B \in \mathcal{A}$ are **independent** if $P(A \cap B) = P(A)P(B)$.

A family $(A_j)_{j \in \Lambda}$ with $A_j \in \mathcal{A}$ is called **independent** if $P(\bigcap_{j \in J} A_j) = \prod_{j \in J} P(A_j)$ for all finite $J \subseteq \Lambda$.

independence and disjoint
are NOT the same thing!

Example: 2 dice rolls with order: (Ω, \mathcal{A}, P) and $A = \text{"first throw gives 6"}, B = \text{"sum of throws = 7"}$

$\Omega = \{(1,1), (1,2), \dots, (6,6)\} \cong \mathbb{Z}^2$ "uniform distribution" $P(\{(w_1, w_2)\}) = \frac{1}{36}$

$$\text{Then } P(A) = \frac{1}{6} \text{ and } P(B) = P(\{(1,6), (2,5), (3,4), (4,3), (5,2), (6,1)\}) = \frac{6}{36} = \frac{1}{6}.$$

$$\Rightarrow P(A \cap B) = \frac{1}{6} P(\{(1,6)\}) = \frac{1}{36} = P(A)P(B) \Rightarrow A \text{ and } B \text{ are independent.}$$

Example: throw a point into $[0,1]$. Then (Ω, \mathcal{A}, P)

$\Omega = [0,1] \cong \mathbb{R}$ "uniform density distribution" $f: \Omega \rightarrow \mathbb{R}$ the indicator function $\mathbf{1}_{\Omega}(x) = \begin{cases} 1 & \text{if } x \in \Omega; \\ 0 & \text{otherwise.} \end{cases}$

For independent $A, B \in \mathcal{A}$ we have:

$$\int_{A \cap B} \mathbf{1}_{\Omega}(x) dx = P(A \cap B) = P(A)P(B) = \int_A \mathbf{1}_{\Omega}(x) dx \int_B \mathbf{1}_{\Omega}(x) dx$$

|| ||

$$\int_{[0,1]} \mathbb{1}_{A \cap B}(x) dx = \int_{[0,1]} \mathbb{1}_A(x) dx \int_{[0,1]} \mathbb{1}_B(x) dx$$

Theorem (Bayes' Formula): Let (Ω, \mathcal{A}, P) define a probability space. Let $B_1, \dots, B_k \in \mathcal{A}$ be events that partition Ω , so $B_i \cap B_j = \emptyset$ and $\bigcup_{i=1}^k B_i = \Omega$. Say $P(B_i) > 0$ for all i . Let $A \in \mathcal{A}$.

$$\text{Then } P(B_i | A) = \frac{P(A|B_i)P(B_i)}{\sum_{j=1}^k P(A|B_j)P(B_j)}$$

by definition of the
conditional probability measure

Proof: For each i , we have $P(B_i | A) = \frac{P(A \cap B_i)}{P(A)} = \frac{P(A|B_i)P(B_i)}{P(A)}$

$$= \frac{P(A|B_i)P(B_i)}{\sum_{j=1}^k P(A|B_j)P(B_j)}.$$

Done. \square

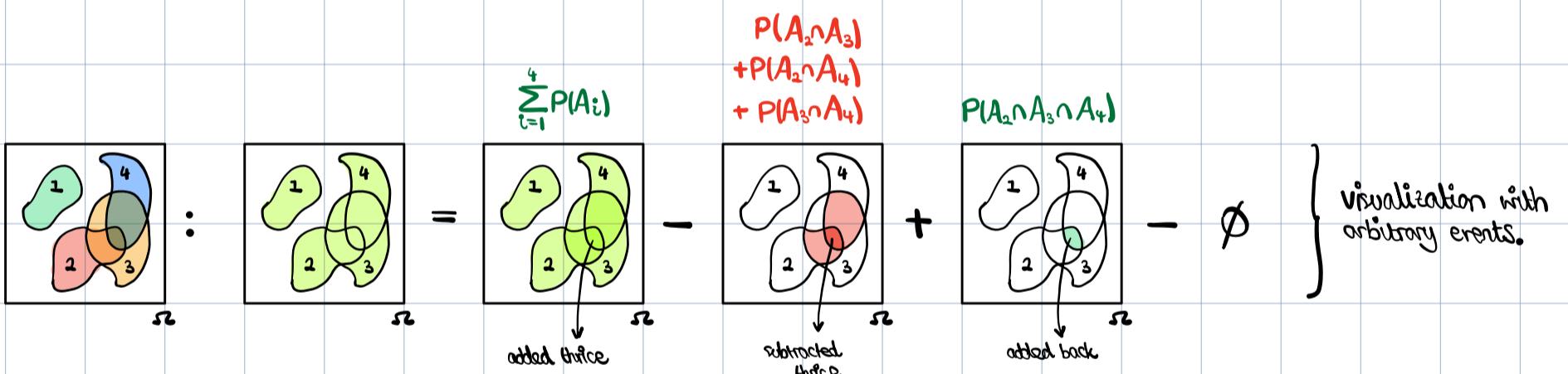
by Law of Total
Probability

a family $(A_i)_{i \in \mathbb{Z}}$ with $A_i \in \mathcal{A}$ is independent if
 $P(\bigcap_{j \in J} A_j) = \prod_{j \in J} P(A_j)$ for any finite $J \subseteq \mathbb{Z}$.

We say $A_1, \dots, A_n \in \mathcal{A}$ are **mutually** (or simply) **independent** if the equation $P(A_{i_1} \cap \dots \cap A_{i_k}) = P(A_{i_1}) \cdots P(A_{i_k})$ holds for all combinations of indices i_1, \dots, i_k distinct.

Proposition (Inclusion-Exclusion): Let (Ω, \mathcal{A}, P) define a space. For arbitrary $A_1, \dots, A_n \in \mathcal{A}$, we have

$$P(\bigcup_{i=1}^n A_i) = \sum_{i=1}^n P(A_i) - \sum_{i < j} P(A_i \cap A_j) + \sum_{i < j < k} P(A_i \cap A_j \cap A_k) - \dots + (-1)^{n-1} P(\bigcap_{i=1}^n A_i).$$



Exercise: There are n letters and n envelopes. Letters are randomly assigned to envelope. What is the probability that no letter is assigned to the envelope it belongs to?

Let $A_i := \text{"envelope } i \text{ is assigned to letter } i"$ so we want $P(\bigcup_{i=1}^n A_i)$.

Notice that $P(A_i) = \frac{1}{n}$ for all i .

$P(A_1 \cup A_2 \cup \dots \cup A_n)$
 $\approx A_1 \text{ or } A_2 \text{ or } \dots \text{ or } A_n$

there are $n-1$ choices left
given A_i was successful

$$P(A_i \cap A_j) = P(A_j | A_i) P(A_i) = \frac{1}{n-1} \cdot \frac{1}{n} \text{ for } i < j.$$

$$P(A_i \cap A_j \cap A_k) = P(A_j \cap A_k | A_i) P(A_i)$$

$$\vdots = \frac{1}{(n-2)(n-1)} \cdot \frac{1}{n} \text{ for } i < j < k.$$

$$P(A_1 \cap \dots \cap A_n) = \frac{1}{n!}.$$

$$\Rightarrow P(\bigcup_{i=1}^n A_i) = \sum_{i=1}^n \frac{1}{n!} - \sum_{i < j} \frac{(n-2)!}{n!} + \sum_{i < j < k} \frac{(n-3)!}{n!} - \dots + (-1)^{n-1} \frac{1}{n!} \quad \text{and now we count:}$$

$$= n \cdot \frac{1}{n} - \binom{n}{2} \frac{(n-2)!}{n!} + \binom{n}{3} \frac{(n-3)!}{n!} - \dots + (-1)^{n-1} \frac{1}{n!}$$

$$= 1 - \sum_{r=2}^n (-1)^{r-1} \binom{n}{r} \frac{(n-r)!}{n!} = 1 - \sum_{r=2}^n (-1)^{r-1} \frac{n!}{r!(n-r)!} \frac{(n-r)!}{n!} = 1 - \sum_{r=2}^n (-1)^{r-1} \frac{1}{r!}.$$

Example: Randomly take 13 cards from a deck of 52. Define $B := \text{"get exactly 6 spades, or exactly one ace."}$. Find $P(B)$.

Let $A_1 = \text{"getting exactly 6 spades"}$ and $A_2 = \text{"getting exactly 1 ace"}$.

Then $P(B) = P(A_1 \cup A_2) = P(A_1) + P(A_2) - P(A_1 \cap A_2)$.

$$\begin{array}{c} \left(\begin{array}{c} 13 \\ 6 \end{array} \right) \left(\begin{array}{c} 39 \\ 7 \end{array} \right) \\ \frac{1}{52!} \left(\begin{array}{c} 13 \\ 13 \end{array} \right) \end{array}$$

pick 6 cards from 13 spades

pick 7 cards from 39 non-spades

$\left(\begin{array}{c} 4 \\ 1 \end{array} \right) \left(\begin{array}{c} 48 \\ 12 \end{array} \right)$

$\frac{1}{52!} \left(\begin{array}{c} 52 \\ 13 \end{array} \right)$

$\left(\begin{array}{c} 13 \\ 13 \end{array} \right)$ pick 13 cards from 52 cards

$A_1 \cap A_2 := \text{getting 6 spades and 1 ace.}$

$\left(\begin{array}{c} 12 \\ 6 \end{array} \right) \left(\begin{array}{c} 3 \\ 1 \end{array} \right) \left(\begin{array}{c} 36 \\ 5 \end{array} \right) + \left(\begin{array}{c} 12 \\ 5 \end{array} \right) \left(\begin{array}{c} 1 \\ 1 \end{array} \right) \left(\begin{array}{c} 36 \\ 7 \end{array} \right)$

6 non-ace spades

1 non-spade ace

getting a non-spade ace

getting a spade ace

$\left(\begin{array}{c} 12 \\ 6 \end{array} \right) \left(\begin{array}{c} 3 \\ 1 \end{array} \right) \left(\begin{array}{c} 36 \\ 5 \end{array} \right) + \left(\begin{array}{c} 12 \\ 5 \end{array} \right) \left(\begin{array}{c} 1 \\ 1 \end{array} \right) \left(\begin{array}{c} 36 \\ 7 \end{array} \right)$

$\frac{1}{52!} \left(\begin{array}{c} 13 \\ 13 \end{array} \right) + \frac{1}{52!} \left(\begin{array}{c} 13 \\ 13 \end{array} \right)$

5 non-ace spades

1 spade ace

7 non-ace non-spade (-3)

7 non-spade (-13)

Let (Ω, \mathcal{A}, P) be a probability space. We say $A, B \in \mathcal{A}$ are conditionally independent given $C \in \mathcal{A}$

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

Exercise: say $A, B \in \mathcal{A}$ with $P(A) = 0.4$ and $P(B) = 0.7$, and $P(A \cup B) = 0.8$.

a) Find $P(A|B)$. We have $P(A|B) = \frac{P(A \cap B)}{P(B)}$ so we need $P(A \cap B)$. Notice $P(A \cap B) = P(A) + P(B) - P(A \cup B)$.
so we are done (use calculator).

b) Find $P(A|B^c)$. We have $P(A|B^c) = \frac{P(A \cap B^c)}{P(B^c)} \rightsquigarrow P(A \cap B^c) = P(A) - \underbrace{P(A \cap B)}$. Done.

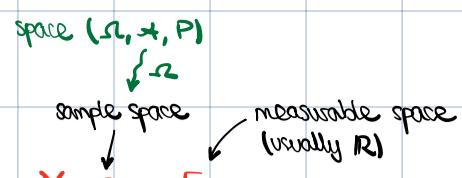
$$\frac{P(B^c)}{1 - P(B)}$$

from (a)

c) Find $P(A^c|B^c)$. We have $P(A^c|B^c) = \frac{P(A^c \cap B^c)}{P(B^c)} \rightsquigarrow P(A^c \cap B^c) = P((A \cup B)^c) = 1 - P(A \cup B)$. Done.

$$\frac{P(B^c)}{1 - P(B)}$$

A random variable X is a measurable function $X: \Omega \rightarrow E$.



The probability that X takes on a value in a measurable $S \subseteq E$ is denoted $P(X \in S) = P(\{w \in \Omega | X(w) \in S\})$.

probability measure $E \in \mathcal{A}$

When $|\text{Image}(X)| \leq \aleph_0$, we call X a discrete random variable and its distribution discrete.

⇒ defined by probability mass functions

When $\text{Image}(X) \subseteq E$ is uncountable we call X a continuous random variable.

In the special case that X is absolutely continuous, its distribution can be defined by a probability density function.
 (in particular, each point must have probability zero)

Here's a more intuitive (yet still accurate) explanation.

Suppose we roll 2 dice, so we have a space (Ω, \mathcal{A}, P) . We may be interested in probabilities of their sums.

Define a function $X: \Omega \rightarrow \mathbb{R}$ by

$$X((x_1, x_2)) = x_1 + x_2$$

$\Omega = \{(1, 1), (1, 2), \dots, (6, 6)\}$

$P(\Omega)$

uniform

discrete

and this is a random variable.

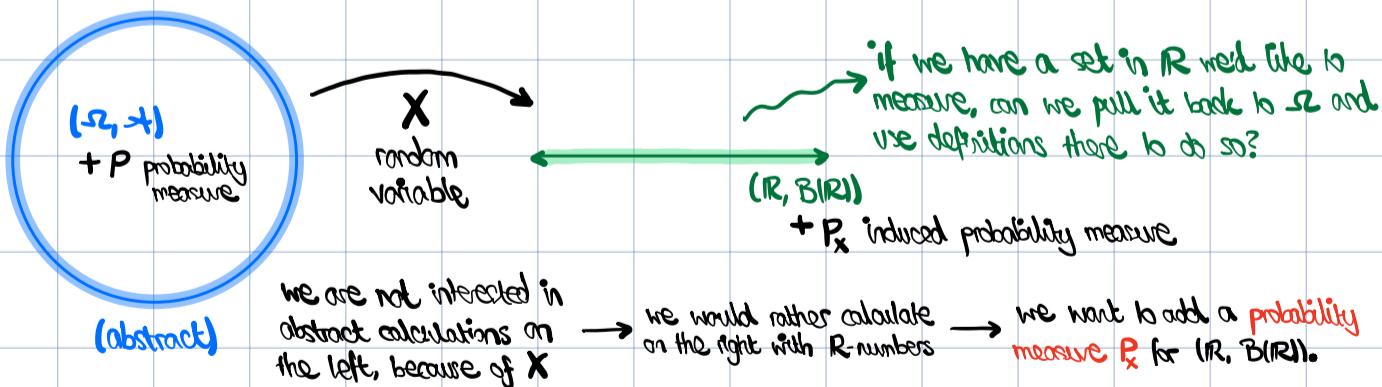
Let (Ω, \mathcal{A}) and $(\tilde{\Omega}, \tilde{\mathcal{A}})$ be measurable spaces (event spaces). A (measurable) map $X: \Omega \rightarrow \tilde{\Omega}$ is called a random variable if $X^{-1}(\tilde{A}) \in \mathcal{A}$ for all $\tilde{A} \in \tilde{\mathcal{A}}$.

Define (Ω, \mathcal{A}) and $(\tilde{\Omega}, \tilde{\mathcal{A}})$ for rolling two dice. Define $X: \mathcal{A} \rightarrow \tilde{\mathcal{A}} = \mathbb{R}$.

$$(w_1, w_2) \mapsto w_1 + w_2$$

Then $X^{-1}(\tilde{A}) \in \mathcal{A} = P(\Omega)$ is trivially fulfilled. ↗ why? I don't fully understand here.

Let (Ω, \mathcal{A}) and $(\tilde{\Omega}, \tilde{\mathcal{A}})$ be measurable spaces (event spaces). On one hand, we have an abstract event space, and on the other, a very concrete one.



Let (Ω, \mathcal{A}, P) be a probability space, and $X: \Omega \rightarrow \mathbb{R}$ a random variable.

We define $P_X: B(\mathbb{R}) \rightarrow [0, 1] \subseteq \mathbb{R}$ with $P_X(B) = P(X^{-1}(B)) = P(\{w \in \Omega | X(w) \in B\})$ to be the probability distribution of X .

\uparrow
 $\in B(\mathbb{R})$
 pre-image of $B \in B(\mathbb{R})$ under random var X

another notation

Proposition (P): Let (Ω, \mathcal{A}, P) be a probability space and $X: \Omega \rightarrow \mathbb{R}$ a random variable.

Then P_X defines a probability measure on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$.

Proof: Notice $X^{-1}(\Omega) = \Omega$, and $X^{-1}(\emptyset) = \emptyset$. We proceed by definition.

$$1 \quad P_X(\Omega) = P(X \in \Omega) = P(X^{-1}(\Omega)) = \underline{P(\Omega)} = 1. \quad \Delta$$

$$2 \quad P_X(\emptyset) = P(X \in \emptyset) = P(X^{-1}(\emptyset)) = \underline{P(\emptyset)} = 0. \quad \Delta$$

3 **T-additivity:** Say $\{A_i\}_{i \in \Lambda}$ is a pairwise disjoint countable collection with $A_i \in \mathcal{B}(\mathbb{R})$.

Thus by the stability of pre-images under \cap , we have $i \neq j \Rightarrow X^{-1}(A_i) \cap X^{-1}(A_j) = X^{-1}(A_i \cap A_j) = \emptyset$.

$\Rightarrow \{X^{-1}(A_i)\}_{i \in \Lambda}$ pairwise disjoint. Δ

We have $P_X(\bigcup_{i \in \Lambda} A_i) = P(X \in \bigcup_{i \in \Lambda} A_i) = P(X^{-1}(\bigcup_{i \in \Lambda} A_i))$ but pre-images are stable under unions, so

$$= P\left(\bigcup_{i \in \Lambda} X^{-1}(A_i)\right) = \sum_{i \in \Lambda} P(X^{-1}(A_i)) = \sum_{i \in \Lambda} P_X(A_i). \quad \Delta$$

↑-additivity of P since $X^{-1}(A_i) \in \mathcal{A}$
and $\{X^{-1}(A_i)\}_{i \in \Lambda}$ pairwise disjoint

4 Trivially $P_X(A) \geq 0$ for $A \in \mathcal{B}(\mathbb{R})$ since $P_X(A) = P(X^{-1}(A)) \geq 0$. $\Delta \square$

If \tilde{P} is a probability measure on \mathbb{R} , and $P_X = \tilde{P}$, then we write $X \sim \tilde{P}$ and say X is distributed as \tilde{P} .

Example: Consider n tosses of a coin with probability p_n of tossing a heads. Define $X: \Omega \rightarrow \mathbb{R}$ by $X(\omega) = \#\text{1s in } \omega$. $\rightarrow X \sim \text{Bin}(n, p)$

Then (Ω, \mathcal{A}, P)
 $\begin{array}{ll} \{\text{H}, \text{T}\} & \text{Beroulli} \\ P(\omega) & \xrightarrow{\text{number of 1s in } \omega} \\ \text{TH} & P(\{\omega\}) = p_n^{\#1} (1-p_n)^{\#0} \end{array}$

X
random variable
 $X \sim \text{Bin}(n, p)$

Example: Roll two 4-sided dice such that (Ω, \mathcal{A}, P) is defined as standard. Define $X: \Omega \rightarrow \mathbb{R}$ the sum of the rolls.

$$\text{Then } X^{-1}\{4\} = \{\omega \in \Omega \mid X(\omega) = 4\} = \{(1, 3), (2, 2), (3, 1)\} \subseteq \mathcal{A}.$$

$$\text{Then } P_X(\{4\}) = P(X \in \{4\}) = P(X^{-1}\{4\}) = P(\{(1, 3), (2, 2), (3, 1)\}).$$

Remark: An alternate common notation for $P_X(\{a_1, \dots, a_n\})$ is $P(X \in \{a_1, \dots, a_n\})$.

$$P_X([a, b]) \text{ is } P(\underbrace{a \leq X \leq b}_{X \in [a, b]}).$$

We define a cumulative distribution function (cdf) and will show that every random variable X has a cdf.

(Ω, \mathcal{A}, P)
abstract

X
random variable
 $(\mathbb{R}, \mathcal{B}(\mathbb{R}), P_X)$

We define the cdf by $F_X: \mathbb{R} \rightarrow [0, 1]$ with $F_X(x) = P_X(-\infty, x]$ i.e.

$$= P(X \in (-\infty, x]) \text{ i.e.}$$

$$= P(X \leq x).$$

- nondecreasing (monotonically increasing)
- $\lim_{x \rightarrow -\infty} F_X = 0$
- $\lim_{x \rightarrow \infty} F_X = 1$
- F_X is right-continuous ($\lim_{x \rightarrow x_0} F_X(x) = F(x_0)$)

reason: x included in interval
 \Rightarrow singleton gives nonempty set

Let (Ω, \mathcal{A}, P) be a discrete probability space, and $X: \Omega \rightarrow \mathbb{R}$ a random variable.

We define the expected value (or expectation) of X as the weighted sum $E[X] = \sum_{x \in \Omega} P_X(x) \cdot x$.

$$E[X] = \sum_{x \in \Omega} P_X(x) \cdot x$$

ER
P_X(x) > 0
ER, so
countable sum in the
discrete case.

Proposition (Linearity of Expectation): Let X be a discrete random variable. Then for $a, b \in \mathbb{R}$, we have $E[aX+b] = aE[X]+b$ (linearity).

Proof: Computation: $E[aX+b] = \sum_x (ax+b) P_X(x) = a \sum_x x P_X(x) + b \sum_x P_X(x) = aE[X]+b$ by definition. \square

see Theorem below!

Theorem (Law of Unconscious Statistician): Let (Ω, \mathcal{A}, P) be a discrete probability space. Let $X: \Omega \rightarrow \mathbb{R}$ be a random variable, and

say $g: \mathbb{R} \rightarrow \mathbb{R}$ is measurable. Then $E[g(X)] = \sum_{x \in \Omega} g(x) P_X(x)$.

Proof: By definition, we have $X: \Omega \rightarrow \mathbb{R}$ a measurable function. Main idea: Define a composition random variable $\tilde{X} = g(X) = X \circ g$.

We have $\tilde{X}: \Omega \xrightarrow{X} \mathbb{R} \xrightarrow{g} \mathbb{R}$ and \tilde{X} measurable since X and g measurable.

pullback of A
under \tilde{X}

Consider the induced probability measure $P_g(x) = P_{\tilde{X}}$ defined by $P_{\tilde{X}}(A) = P(\tilde{X} \in A) = P(\tilde{X}^{-1}(A))$.

For $A \in \Omega$ we have $\underset{\Omega}{A} \mapsto \underset{\mathbb{R}}{X(A)} \mapsto \underset{\mathbb{R}}{g(X(A))}$, so $\tilde{X}^{-1}(A) = (X \circ g)^{-1}(A) = (g^{-1} \circ X^{-1})(A) = X^{-1}(g^{-1}(A))$.

We have $E[g(X)] = E[\tilde{X}] = \sum_x x \cdot P_{\tilde{X}}(x)$
 $= \sum_x x \cdot P(X^{-1}(g^{-1}(x)))$ so consider $y = g^{-1}(x) \Rightarrow x = g(y)$:
 $= \sum_y g(y) \cdot P(X^{-1}(y)) = \sum_y g(y) P_X(y)$. We are done. \square

Elementary Distributions For this section, we do computation. Implicitly assume (Ω, \mathcal{A}, P) defines our abstract space.

Proposition (E[Binomial]): Let $X \sim \text{Bin}(n, p)$, so $P_X(x) = \binom{n}{x} p^x (1-p)^{n-x}$ where p is the probability of event success.

Then $E[X] = np$.

Proof: By definition, we have $E[X] = \sum_{x \in \Omega} P_X(x) \cdot x$
 $P_X(x) > 0 \quad \leftarrow n \text{ trials so at most } x=n \text{ successes}$
 $\text{at least } x=0 \text{ successful}$

(recall we have $P(\Omega) = P_X(\mathbb{R}) = \sum_{x=0}^n \binom{n}{x} p^x (1-p)^{n-x} = 1$)

by definition of
a probability measure

$$\begin{aligned} &= \sum_{x=0}^n x \cdot \binom{n}{x} p^x (1-p)^{n-x} \\ &= \sum_{x=0}^n x \cdot \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} \\ &= \sum_{x=1}^n x \cdot \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} \\ &= np \sum_{x=1}^n \frac{(n-1)!}{(x-1)!(n-x)!} p^{x-1} (1-p)^{n-x} \quad \text{and now we re-index:} \\ &= np \sum_{x=1}^n \frac{(n-1)!}{(x-1)!(n-1-(x-1))!} p^{x-1} (1-p)^{n-1-(x-1)} \\ &= np \sum_{x=0}^m \frac{m!}{x!(m-x)!} p^x (1-p)^{m-x} = np \quad \text{and we are done.} \quad \square \end{aligned}$$

1

We define the Poisson distribution to be $\frac{e^{-\mu} \mu^x}{x!}$ where μ is an intensity parameter. We revisit this when looking at Stochastic processes.

Proposition (E[Risson]): Let $X \sim \text{Poisson}(\lambda)$. Then $E[X] = \lambda$.

Proof: By definitions. We have $E[X] = \sum_{x=0}^{\infty} x \cdot P_X(x) = \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} = \sum_{x=1}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} = \sum_{x=1}^{\infty} \frac{e^{-\lambda} \lambda^x}{(x-1)!} = \lambda \sum_{x=0}^{\infty} \frac{e^{-\lambda} \lambda^{x+1}}{(x+1)!} = \lambda$ and we are done. \square

Let (Ω, \mathcal{A}, P) be a probability space, and $X: \Omega \rightarrow \mathbb{R}$ a random variable.

We define the variance by $\text{Var}(X) = E[(X - \mu)^2]$ where $\mu = E[X]$ (so $\text{Var}(X) = E[(X - E[X])^2]$).
 "measure of spread" measuring how much X deviates from mean

Proposition (Properties of Variance): Let X be a random variable. Then:

$$1. \text{Var}(aX+b) = a^2 \text{Var}(X) \text{ for } a, b \in \mathbb{R};$$

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

Proof (1): Definitions. $\text{Var}(X) = E[(X - E[X])^2]$

$$\begin{aligned} \Rightarrow \text{Var}(aX+b) &= E[(aX+b - E[aX+b])^2] = E[(aX+b - aE[X]-b)^2] \\ &= E[a^2(X-E[X])^2] = a^2 E[(X-E[X])^2] = a^2 \text{Var}(X). \quad \Delta \end{aligned}$$

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

$$\begin{aligned} &= E[X^2] - 2\mu E[X] + \mu^2 = E[X^2] - \mu^2. \quad \Delta \square \end{aligned}$$

Proposition (Poisson Approximation for Binomial Distribution): Let $X \sim \text{Binomial}(n, p)$ and $Y \sim \text{Poisson}(\lambda)$ where $\lambda = np$.

Then $X \approx Y$ for large n , and $X \xrightarrow{n \rightarrow \infty} Y$.

Proof (Rough): Notice $E[X] = np = E[Y]$. This justifies our choice $\lambda = np$. Thus $p = \frac{\lambda}{n}$.

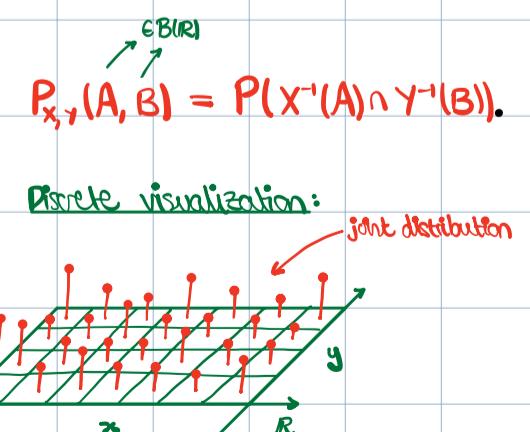
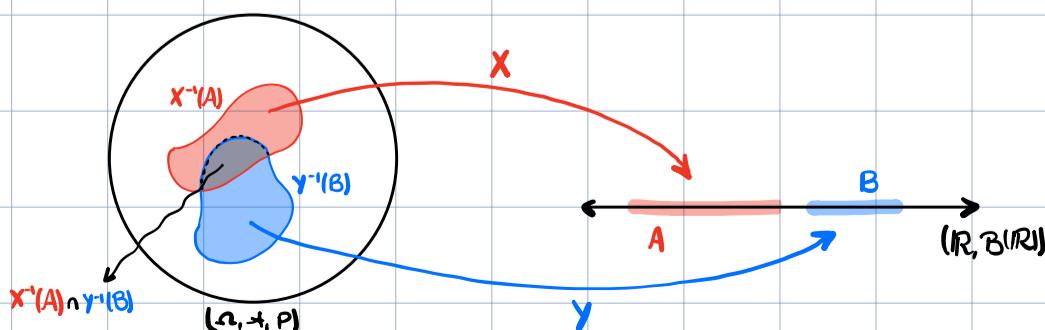
$$\begin{aligned} \text{We have } P_X(x) &= \binom{n}{x} p^x (1-p)^{n-x} \\ &= \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} \\ &= \frac{1}{x!} \underbrace{n(n-1)\dots(n-x+1)}_{\substack{n \dots n \\ x \text{ times}}} p^x (1-p)^{n-x} \xrightarrow{n \rightarrow \infty} \frac{1}{x!} n^x \left(\frac{\lambda}{n}\right)^x (1-\frac{\lambda}{n})^n = \frac{1}{x!} \lambda^x e^{-\lambda} = P_Y(x). \quad \square \end{aligned}$$

Multivariate Distributions

Let (Ω, \mathcal{A}, P) be a probability space, and let $X: \Omega \rightarrow \mathbb{R}$ be random variables on Ω .
 $Y: \Omega \rightarrow \mathbb{R}$

Q: How do X and Y relate/interact?

We can define a bivariate probability function $P_{X,Y}(x,y)$ induced by X and Y , by $P_{X,Y}(A,B) = P(X^{-1}(A) \cap Y^{-1}(B))$.
 $P_{X,Y}: \mathcal{B}(\mathbb{R})^2 \rightarrow [0,1], A, B \subseteq \mathbb{R}$



Let (Ω, \mathcal{A}, P) be a probability space and $X, Y: \Omega \rightarrow \mathbb{R}$ be random variables. Then X and Y are independent random variables

if $X^{-1}((-\infty, x])$ and $Y^{-1}((-\infty, y])$ are independent events for any $x, y \in \mathbb{R}$.

$$\Leftrightarrow P(X^{-1}((-\infty, x]) \cap Y^{-1}((-\infty, y])) = P(X^{-1}((-\infty, x])) P(Y^{-1}((-\infty, y))) \Leftrightarrow \underbrace{P(X \leq x, Y \leq y)}_{F_{(X,Y)}(x,y)} = F_X(x) F_Y(y).$$

We define $F_{(X,Y)}(x,y) = P(X \leq x, Y \leq y) = P(X^{-1}((-\infty, x]) \cap Y^{-1}((-\infty, y]))$ the joint cumulative distribution function

$$F_{(X,Y)}: \Omega \rightarrow \mathbb{R}^2. \text{ Thus random variables } X \text{ and } Y \text{ are independent} \Leftrightarrow F_{(X,Y)}(x,y) = F_X(x) F_Y(y).$$

Example: Consider Ω a product space $\Omega = \Omega_1 \times \Omega_2$. Then any random variables of the form

$$X: \Omega \rightarrow \mathbb{R} \text{ and } Y: \Omega \rightarrow \mathbb{R} \text{ are independent.}$$

$(\omega_1, \omega_2) \mapsto f(\omega_1)$ $(\omega_1, \omega_2) \mapsto g(\omega_2)$

A family $(X_i)_{i \in I}$ of random variables is independent if $P((X_j \leq x_j)_{j \in J}) = \prod_{j \in J} P(X_j \leq x_j)$ for all finite $J \subseteq I$, for all $x_j \in \mathbb{R}$.

Let (Ω, \mathcal{A}, P) be a probability space, and $X: \Omega \rightarrow \mathbb{R}$ a random variable. We define the (general) expectation

$$E[X] = \int_{\Omega} X dP.$$

measure P
Lebesgue integral

Say we have a measurable function $g: \mathbb{R} \rightarrow \mathbb{R}$. We can define a new random variable $g(X): \Omega \rightarrow \mathbb{R}$.

How does $E[X]$ relate to $E[g(X)]$?

$$\text{Consider } \int_A g(X) dP \text{ for } A \in \mathcal{A}. \text{ We have } \int_A g(X) dP = \int_A g(X(w)) dP(w) = \int_{X(A)} g(x) d(P \circ X^{-1})(x) = \int_{X(A)} g(x) dP_X(x).$$

Integration in some abstract space
change of variables
 $\text{so } X^{-1}(x) = \omega$
 P_X new measure
Integration in \mathbb{R}

$$\text{Discrete case: } \sum_{x \in X(A)} g(x) P_X(\{x\}).$$

$$\text{Continuous case: } \int_{X(A)} g(x) f_X(x) dx.$$

Proposition (Monotonicity of Expectation): Let (Ω, \mathcal{A}, P) be a probability space, and $X, Y: \Omega \rightarrow \mathbb{R}$ random variables.

If $X \leq Y$ almost surely (so $P(\{\omega \in \Omega | X(\omega) \leq Y(\omega)\}) = 1$), then $E[X] \leq E[Y]$.

We define the standard deviation $\sigma(X)$ of a random variable X by $\sigma(X) = \sqrt{\text{Var}(X)}$.

Proposition: Let X and Y be independent random variables. Then $E[XY] = E[X]E[Y]$ and $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$.

Proof: By definitions and algebra. Not very insightful. \square

Proposition: Let X and Y be independent random variables, and $g, h: \mathbb{R} \rightarrow \mathbb{R}$ be measurable. Then $E[g(X)h(Y)] = E[g(X)]E[h(Y)]$.

Proof: Denote our probability space (Ω, \mathcal{A}, P) . We have $E[g(X)] = \int_{\Omega} g(X) dP = \int_{\mathbb{R}} g(x) dP_X(x)$.

$$E[h(Y)] = \int_{\Omega} h(Y) dP = \int_{\mathbb{R}} h(y) dP_Y(y).$$

We prove the following Theorem instead:

Theorem (Functions Preserve Independence): Let X and Y be independent random variables $X, Y: \Omega \rightarrow \mathbb{R}$.

Let $g, h: \mathbb{R} \rightarrow \mathbb{R}$ be functions. Then $g(X)$ and $h(Y)$ are independent random variables.

Proof: Let $A, B \subseteq \mathbb{R}$. Notice $g'(A), h'(B) \subseteq \mathbb{R}$. Then we have

$$P(g(X) \in A, h(Y) \in B) = P(X \in g^{-1}(A), Y \in h^{-1}(B)) \xrightarrow{\text{independence}} P(X \in g^{-1}(A)) P(Y \in h^{-1}(B)) = P(g(X) \in A) P(h(Y) \in B).$$

Done by definitions. \square

We showed $g(X)$ and $h(Y)$ are independent, so we are done. \square

Let (Ω, \mathcal{A}, P) be a probability space, and $X, Y: \Omega \rightarrow \mathbb{R}$ random variables. We define the covariance of X and Y by

$$\text{Cov}(X, Y) = E[\underbrace{(X - E[X])}_{X - \mu_X} \underbrace{(Y - E[Y])}_{Y - \mu_Y}] \quad (\text{comparing how } (X - \mu_X) \text{ and } (Y - \mu_Y) \text{ relate}).$$

Corollary: We have $\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$.

Corollary: Let X and Y be independent random variables. Then $\text{Cov}(X, Y) = 0$.

Proposition (Cauchy-Schwarz): Let X and Y be random variables. Then $\text{Cov}(X, Y)^2 \leq \text{Var}(X)\text{Var}(Y)$.

Let X and Y be random variables. We define the correlation coefficient ρ by $\rho = \frac{\text{Cov}(X, Y)}{\sigma(X)\sigma(Y)} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$.

Proposition (Correlation): Let X, Y be random variables. Then $|\rho| \leq 1$ and $|\rho|=1 \iff Y = aX+b$ for some $a, b \in \mathbb{R}$.

Proof: We have $|\rho| \leq 1$ immediate by Cauchy-Schwarz. \square

Example: Suppose 15 televisions are to be purchased by a local bar from a large production run of 500.

Say 450 will last at least 5 years without needing repair.

a) Calculate the probability that ≥ 12 last ≥ 5 years.

We have a Binomial distribution with $p_{\text{success}} = 0.9$. Then say $X \sim \text{Bin}(15, 0.9)$

$$\begin{aligned} \text{so we want } P(X \geq 12) &= P(X=12) + P(X=13) + P(X=14) + P(X=15) \\ &= \sum_{x=12}^5 \binom{15}{x} p^x (1-p)^{15-x} = \binom{15}{12} (0.9)^{12} (0.1)^3 + \binom{15}{13} (0.9)^{13} (0.1)^2 \approx 0.944. \\ &\quad + \binom{15}{14} (0.9)^{14} (0.1)^1 + \binom{15}{15} (0.9)^{15} (0.1)^0 \end{aligned}$$

We can also approximate this using a Poisson process, where on average 0.02 per television dies annually, or

$X \sim \text{Poisson}(0.02)$. With 15 televisions and 5 years, we'd take $X \sim \text{Poisson}(0.10)$ and ≈ 0.10 per television

$P_X(x) = \frac{e^{-\lambda} \lambda^x}{x!} = \frac{e^{-0.10} 0.10^x}{x!}$ describes the probability of seeing x such events over 5 years.

$$\text{We want } \sum_{x=0}^3 P_X(x) = (e^{-0.10}) \left(\frac{0.10^0}{0!} + \frac{0.10^1}{1!} + \frac{0.10^2}{2!} + \frac{0.10^3}{3!} \right) \approx 0.934.$$

Consider wanting to model the distribution of events over time, with the following assumptions:

1 **Independence.** The probability of events occurring in non-overlapping intervals are independent.

$$\lim_{\Delta t \rightarrow 0} \frac{o(\Delta t)}{\Delta t} = 0$$

2 **Individuality.** Probability of simultaneous events occurring is 0, so $P[\geq 2 \text{ events in } [t, t+\Delta t]] = o(\Delta t)$

3 **Homogeneity.** Regardless of time interval $[t, t+\Delta t]$, we have $P(X \text{ occurs in } [t, t+\Delta t]) = \lambda \Delta t + o(\Delta t)$.

Consider a unit time interval, and slice it uniformly into n intervals:

and define $\Delta t = \frac{t}{n}$.

Notice $P[1 \text{ event in interval of size } \Delta t] \xrightarrow{\Delta t \rightarrow 0} \lambda \Delta t = \lambda \frac{t}{n}$. We can model this binomially:

$$P[1 \text{ event in } \Delta t] = \lambda \frac{t}{n} \Rightarrow P[k \text{ events in } t] = \binom{n}{k} \left(\lambda \frac{t}{n}\right)^k \left(1 - \lambda \frac{t}{n}\right)^{n-k} = \frac{n!}{k!(n-k)!} \lambda^k \left(\frac{t}{n}\right)^k \left(1 - \lambda \frac{t}{n}\right)^{n-k}$$

n trials of $\lambda \frac{t}{n}$

$$\lim_{n \rightarrow \infty} \frac{1}{k!} \frac{n \cdot (n-1) \cdot \dots \cdot (n-(k+1))}{n^k} \left(\lambda t\right)^k \left(1 - \frac{\lambda t}{n}\right)^{n-k} \left(1 - \frac{\lambda t}{n}\right)^{-k}$$

$$= \frac{(\lambda t)^k e^{-\lambda t}}{k!}.$$

error term going to 0 as $\Delta t \rightarrow 0$

We call the process above a **Poisson process**, and for $X \sim \text{Poisson}(\lambda)$, define $P_X(x) = \frac{e^{-\lambda} \lambda^x}{x!}$ the **Poisson distribution**.

Let X be a random variable with $\mu_X = E[X]$. Let $k \in N = \{1, 2, \dots, 3\}$ be positive.

We define the k^{th} moment of X to be $E[X^k]$. When $k=1$ this is expectation.

k^{th} central moment of X to be $E[(X-\mu_X)^k]$. When $k=2$ this is variance.

Since $\text{Var}(X) = E[X^2] - E[X]^2$ we see that calculating the 2nd moment of X can be useful. In general, how can we compute these?

Let X be a random variable with a distribution function F_X . The moment generating function of F_X is defined by

$$M(t) = E[e^{tX}] \text{ provided this exists for } t \in B(0, \varepsilon) \text{ for some } \varepsilon \in \mathbb{R}_{>0}.$$

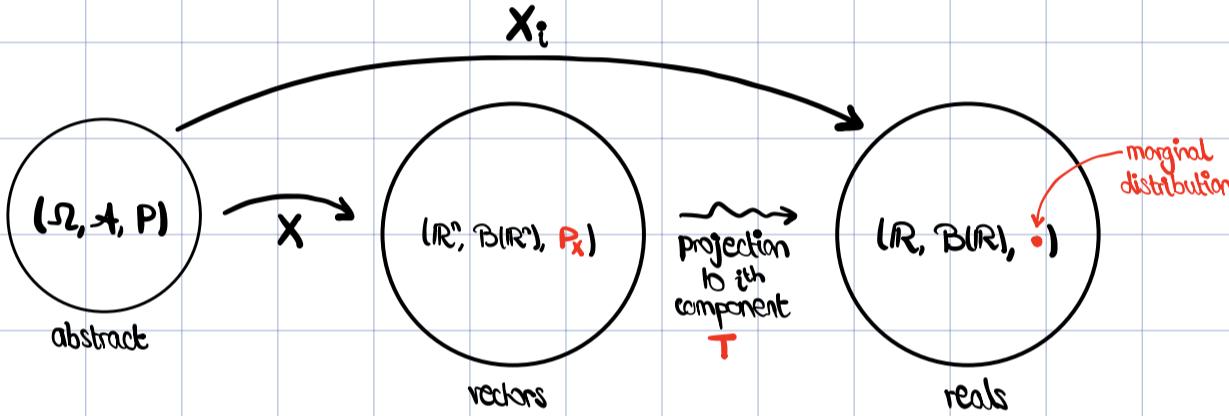
Theorem (Moment Generating Functions): Say X is a random variable with moment generating function $M(t)$.

Then $E[X^n] = M^{(n)}(0)$ (the n^{th} moment of X is the n^{th} derivative of $M(t)$ at $t=0$).

Let (Ω, \mathcal{A}, P) be a probability space, and $X: \Omega \rightarrow \mathbb{R}^n$ be a random variable (or a random vector).

We can write X in terms of component random variables, so $X(\omega) = \begin{pmatrix} X_1(\omega) \\ \vdots \\ X_n(\omega) \end{pmatrix}$ where $X_i: \Omega \rightarrow \mathbb{R}$.

If we know the distribution of X , what can we find out about the distribution of X_i ?

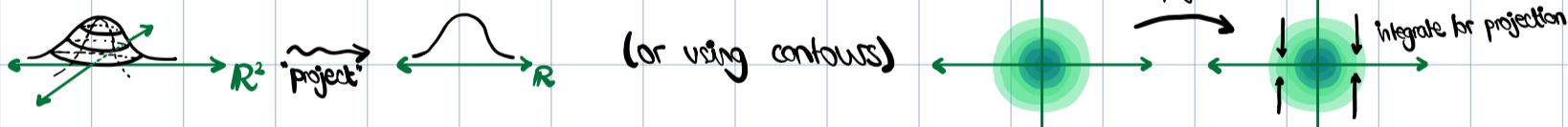


The function $P_{X_i} = (P_X)_T$ is called the marginal distribution of X with respect to component i .

$$\begin{aligned} \text{We can calculate the marginal cumulative distribution function } F_{X_i}(t) &= P_{X_i}((-\infty, t]) = P_X(\mathbb{R} \times \dots \times (-\infty, t] \times \mathbb{R} \times \dots \times \mathbb{R}) \\ &= P_X(X_i \leq t, X_j \in \mathbb{R}), \end{aligned}$$

↑ *ith component*
↓ *j ≠ i, j = 1, 2, …, n*

Continuous Case: P_X has a probability density function $f_X: \mathbb{R}^n \rightarrow \mathbb{R}_{>0}$.



$$\text{Then the distribution } f_{X_i}(t) \text{ can be defined } f_{X_i}(t) = \int_{\mathbb{R}^{n-1}} f_X(x_1, \dots, t, \dots, x_n) d(x_1, \dots, \underset{i^{\text{th}} \text{ component}}{x_i}, \dots, x_{i+1}, \dots, x_n).$$

n-1 integrations

This is the marginal probability density function with respect to the i^{th} variable.

Discrete Case: P_X has a probability mass function $(p_X)_{x \in \mathbb{R}^n}$.

$$\text{Then the marginal probability mass function } (p_i)_{t \in \mathbb{R}} \text{ is } p_i = \sum_{\substack{x_1, \dots \\ \in \mathbb{R}}} p_{x_1, \dots, x_i=t, \dots, x_n}.$$

Recall that, for a space (Ω, \mathcal{A}, P) and $B \in \mathcal{A}$, we have an induced conditional space $(\Omega, \mathcal{A}, P(\cdot|B))$.

Notice a random variable $X: \Omega \rightarrow \mathbb{R}$ naturally lives on both spaces. However, its behaviour may change. given events

For X , we define $E[X] = \int_{\Omega} X dP \rightsquigarrow E[X|B] = \int_{\Omega} X dP(\cdot|B)$ the conditional expectation of X given B .

$$\text{If we write } P(A|B) = P(A \cap B) = \frac{1}{P(B)} \int_A \mathbf{1}_B dP \text{ then } E[X|B] = \frac{1}{P(B)} \int_{\Omega} X \mathbf{1}_B dP = \frac{1}{P(B)} E[\mathbf{1}_B X].$$

indicator function random variable
↑ indicator
↑ scaling factor

Example: Let $X \sim N(0, 1^2)$ so $f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$. Say $B = \mathbb{R}_{>0} = \{x \in \mathbb{R} \mid x > 0\}$.

$$\begin{aligned} \text{Then } E[X|B] &= \frac{1}{P(B)} \int_{\Omega} X(\omega) \mathbf{1}_B(\omega) dP(\omega) = \frac{1}{P(B)} \int_{\Omega} x \underbrace{\mathbf{1}_B(X^{-1}(x))}_{\begin{cases} 1 & \text{for } x > 0, \\ 0 & \text{for } x \leq 0. \end{cases}} f_X(x) dx \\ &= \frac{1}{P(B)} \int_0^{\infty} x f_X(x) dx = 2 \cdot \frac{1}{\sqrt{2\pi}} \int_0^{\infty} x e^{-\frac{1}{2}x^2} dx = 1. \end{aligned}$$

Notice $E[\mathbf{1}_A|B] = \int_{\Omega} \mathbf{1}_A dP(\cdot|B) = \int_A dP(\cdot|B) = P(A|B)$ so we can write conditional probability as conditional expectation.

of indicator functions ↑

Example: Consider a container of N balls, M of which are red, and n are picked without replacement.

Say $X: \Omega \rightarrow \mathbb{R}$ is the random variable counting the number of red balls picked. What is $E[X]$?

(Use indicator functions, define $Y: \Omega \rightarrow \mathbb{R}^n$ a random variable and consider marginals).

Define $Y: \Omega \rightarrow \mathbb{R}^n$ by $Y((w_1, \dots, w_n)) = \begin{pmatrix} Y_1(w_1) \\ \vdots \\ Y_n(w_n) \end{pmatrix}$ where $Y_i(w_i) = \begin{cases} 1 & \text{if } w_i \text{ red;} \\ 0 & \text{otherwise.} \end{cases}$ Notice $(Y_i)_{i=1}^n$ is independent.

Therefore, by $X = \sum_{i=1}^n Y_i$, we have $E[X] = E[\sum_{i=1}^n Y_i] = \sum_{i=1}^n E[Y_i]$.

Consider the marginal $P_i(w_i) = \sum_{x=1}^N P(x_1, \dots, x_i, \dots, x_n)$. Notice $P_i(w_i) = P_j(w_j)$ for any i, j by this.

$$\Rightarrow E[Y_i] = E[Y_j] \text{ for all pairs } i, j. \text{ In particular, } \sum_{i=1}^n E[Y_i] = n E[Y_i] = n(P(Y_i=1) \cdot 1 + 0) = n \left(\frac{M}{N}\right).$$

What's $\text{Var}(X)$? Well, $\text{Var}(X) = E[X^2] - E[X]^2 = E[X^2] - \left(n^2 \frac{M^2}{N^2}\right)$.

We have $E[X^2] = E[(\sum_{i=1}^n Y_i)^2]$. By independence, $E[Y_i Y_j] = E[Y_i] E[Y_j]$ for any $i \neq j$.

$$\begin{aligned} &= E[\sum_{i=1}^n Y_i^2] + \sum_{i < j} E[Y_i] E[Y_j] = E[\sum_{i=1}^n Y_i^2] + \sum_{i < j} E[Y_i]^2 \text{ because } E[Y_i] = E[Y_j]. \\ &= \sum_{i=1}^n E[Y_i^2] + \binom{n}{2} E[Y_i]^2. \end{aligned}$$

$$\text{Notice } E[Y_i^2] = \sum_{x=0}^1 x^2 P(Y_i=x) = 1^2 P(Y_i=1) = \frac{M}{N}.$$

Say $X_1 \sim \text{Poisson}(\mu_1)$ and $X_2 \sim \text{Poisson}(\mu_2)$ are independent random variables. Find the conditional probability mass function of X_1 given $X_1 + X_2 = t$ for some $t \in \mathbb{N}$.

$$\text{We have } P[X_1=x_1 | X_1 + X_2 = t] = \frac{P[X_1=x_1, X_1 + X_2 = t]}{P[X_1 + X_2 = t]} = \frac{P[X_1=x_1, X_2=t-x_1]}{P[X_1 + X_2 = t]} = \frac{P[X_1=x_1] P[X_2=t-x_1]}{P[X_1 + X_2 = t]}$$

definition independence

$$\begin{aligned} \Rightarrow P[X_1=x_1 | X_1 + X_2 = t] &= \left(\frac{e^{-\mu_1} \mu_1^{x_1}}{x_1!} \right) \left(\frac{e^{-\mu_2} \mu_2^{t-x_1}}{(t-x_1)!} \right) (P[X_1 + X_2 = t]). \\ &= \left(\frac{e^{-\mu_1} \mu_1^{x_1}}{x_1!} \right) \left(\frac{e^{-\mu_2} \mu_2^{t-x_1}}{(t-x_1)!} \right) \text{ Poisson}(\mu_1 + \mu_2) \\ &\stackrel{\text{Independence}}{=} \sum_{x_1=0}^t P[X_1=x_1, X_2=t-x_1] \\ &= \sum_{x_1=0}^t P[X_1=x_1] P[X_2=t-x_1], \\ &= \sum_{x_1=0}^t \frac{e^{-\mu_1} \mu_1^{x_1}}{x_1!} \frac{e^{-\mu_2} \mu_2^{t-x_1}}{(t-x_1)!} \\ &= e^{-(\mu_1 + \mu_2)} \mu_2^t \sum_{x_1=0}^t \frac{1}{x_1! (t-x_1)!} \left(\frac{\mu_1}{\mu_2}\right)^{x_1} \end{aligned}$$

$$\begin{aligned}
&= \binom{t}{x_1} \left(\frac{\mu_1}{\mu_1 + \mu_2} \right)^{x_1} \left(1 - \frac{\mu_1}{\mu_1 + \mu_2} \right)^{t-x_1} \\
&= \text{Binomial}(t, \frac{\mu_1}{\mu_1 + \mu_2}). \\
\text{Therefore } P[X_1=x_1 | X_1+X_2=t] &= \text{Binomial}\left(t, \frac{\mu_1}{\mu_1 + \mu_2}\right). \\
&= \frac{e^{-(\mu_1+\mu_2)t}}{t!} \mu_1^t \sum_{x_1=0}^t \frac{t!}{x_1!(t-x_1)!} \left(\frac{\mu_1}{\mu_2}\right)^{x_1} 1^{t-x_1} \\
&\quad \text{binomial} \\
&= \frac{e^{-(\mu_1+\mu_2)t}}{t!} \mu_2^t \left(1 + \frac{\mu_1}{\mu_2}\right)^t \\
&= \text{Poisson}(\mu_1 + \mu_2).
\end{aligned}$$

People arrive at a house from 5:30pm-9:00pm according to a Poisson process, averaging 12 people per hour.

What number of people is most likely to arrive?

We want $\underset{x \in \mathbb{N}}{\operatorname{argmax}} P(M=x)$ where M is the total number of people who arrive. Notice $M \sim \text{Poisson}(\lambda)$

$= \underset{x \in \mathbb{N}}{\operatorname{argmax}} \frac{e^{-\lambda} \lambda^x}{x!}$ but we cannot compute this. How can we proceed?

→ Consider $\frac{P(M=x)}{P(M=x+1)}$ and see how this behaves.

$\lambda = 42$ intensity parameter

Suppose $X \sim \text{Hypergeometric}(N, M, n)$. Find $E[X]$.

$$\text{We have } P_X(x) = \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}} \text{ so } E[X] = \sum_{x=0}^n x P_X(x) = \sum_{x=0}^n x \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}}.$$

$$\begin{aligned}
\sum_{x=0}^n x \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}} &= \sum_{x=1}^n x \frac{M!}{x!(M-x)!} \frac{(N-M)!}{(n-x)!(N-n)!} \cdot \frac{n!(N-n)!}{N!} \\
&= \frac{Mn}{N} \sum_{x=1}^n \frac{(M-1)!}{(x-1)!(M-1-(x-1))!} \frac{(N-M)!}{(n-x)!(N-M-(n-x))!} \frac{(n-1)!(N-1-(n-1))!}{(N-1)!} \\
&= \frac{Mn}{N} \sum_{x=0}^n \frac{M^x}{x^x (M^x - x^x)!} \frac{(N-M)!}{(n-x)!(N-M-(n-x))!} \frac{n^n (N^n - n^n)!}{N^n!} \\
&= \frac{Mn}{N} \mathbf{1} = \frac{Mn}{N}.
\end{aligned}$$

Consider n Bernoulli trials, and let $X_i: \Omega \rightarrow \mathbb{R}$ be the indicator that the i^{th} trial is a success.

In other words, define $X: \Omega \rightarrow \mathbb{R}^n$ by $X((w_1, \dots, w_n)) = \begin{pmatrix} X_1(w_1) \\ \vdots \\ X_n(w_n) \end{pmatrix}$ and notice all the X_i are independent.

Let $Y \sim \text{Bernoulli}(n, p)$ counting the number of successes, so $Y = \sum_{i=1}^n X_i$.

$$\begin{aligned}
\text{Then } \text{Var}(Y) &= \text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{Var}(X_i) + 2 \sum_{i < j} \text{Cov}(X_i, X_j) \text{ but since } X_i \perp X_j, \text{ we have } \text{Cov}(X_i, X_j) = 0. \\
&= \sum_{i=1}^n \text{Var}(X_i) \text{ but each } X_i \text{ is "identical" so } \text{Var}(X_i) = \text{Var}(X_j) \\
&= n \text{Var}(X_1) = n(E[X_1^2] - E[X_1]^2) \\
&= n \left(\sum_{x=0}^1 x^2 P_X(x) - p^2 \right) = n(p^2 - p^2) = np(1-p).
\end{aligned}$$

Say $X \sim \text{Hypergeometric}(N, M, n)$. What's $\text{Var}(X)$?

Well, we have X counting the number of successes from n picks out of N samples.

Consider $Y: \Omega \rightarrow \mathbb{R}^n$ given by $Y(w) = \begin{pmatrix} Y_1(w_1) \\ \vdots \\ Y_n(w_n) \end{pmatrix}$ where $Y_i: \Omega \rightarrow \mathbb{R}$ are indicator random variables.

Notice $X = \sum Y_i$ so $\text{Var}(X) = \text{Var}\left(\sum Y_i\right) = \sum \text{Var}(Y_i) + 2 \sum \text{Cov}(Y_i, Y_j)$.

We have $E[X_i] = \sum_{x=0}^1 x P_{X_i}(x) = P(X_i=1) = \frac{M}{N}$. Notice this is invariant of $i=1, \dots, n$.

$$\text{Therefore } \text{Var}(X_i) = E[X_i^2] - E[X_i]^2 = \sum_{x=0}^1 x^2 P_{X_i}(x) - \left(\frac{M}{N}\right)^2 = 1^2 \frac{M}{N} - \frac{M^2}{N^2}$$

We also have $\text{Cov}(X_i, X_j) = E[X_i X_j] - E[X_i] E[X_j]$

$$= \underbrace{\sum_{x_1=0}^1 \sum_{x_2=0}^1 x_1 x_2 P(X_i=x_1, X_j=x_2)}_{\hookrightarrow 1 \cdot 1 P(X_i=1, X_j=1)} - \left(\frac{M}{N}\right)^2$$

$$\Rightarrow \text{Cov}(X_i, X_j) = \frac{M-1}{N-1} \left(\frac{M}{N}\right) - \left(\frac{M}{N}\right)^2$$

$$\Rightarrow \text{Var}(X) = \sum_{i=1}^n \text{Var}(X_i) + 2 \sum_{i < j} \text{Cov}(X_i, X_j) = n \left(\frac{M}{N} - \frac{M^2}{N^2}\right) + 2 \binom{n}{2} \left(\frac{M-1}{N-1} \left(\frac{M}{N}\right) - \left(\frac{M}{N}\right)^2\right).$$

Theorem (Multivariate LOTUS): Say X_1 and X_2 are random variables $X_1, X_2: \Omega \rightarrow \mathbb{R}$ such that $(X_1, X_2) \sim P(x_1, x_2)$.

Let $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ be measurable. Then $E[g(X_1, X_2)] = \sum_{x_1} \sum_{x_2} g(x_1, x_2) P(X_1=x_1, X_2=x_2)$.

Proof: Recall $P(X_1 \in A, X_2 \in B)$ is defined by $P(X_1^{-1}(A) \cap X_2^{-1}(B)) = P_{X_1, X_2}(A, B)$ the joint (bivariate) distribution.

Notice $(X_1, X_2): \Omega \rightarrow \mathbb{R}^2$ defines a random variable pair. Define a new random variable by

$$\tilde{Z} = g((X_1, X_2)): \Omega \xrightarrow{(X_1, X_2)} \mathbb{R}^2 \xrightarrow{g} \mathbb{R}. \text{ Then notice } \tilde{Z}^{-1}(A) = (X_1, X_2)^{-1}(g^{-1}(A)) \text{ for } A \in \mathcal{B}(\mathbb{R}).$$

$$\text{We have } E[\tilde{Z}] = \sum_z z P_{\tilde{Z}}(z)$$

$$= \sum_z z P(\tilde{Z}^{-1}(\{z\})) = \sum_z z P((X_1, X_2)^{-1}(g^{-1}(\{z\}))) \quad \text{and let } S \subseteq \mathbb{R}^2 \text{ be defined}$$

$$= \sum_S g(S) P((X_1, X_2)^{-1}(S))$$

$$= \sum_{x_1} \sum_{x_2} g(x_1, x_2) P((X_1, X_2)^{-1}((x_1, x_2)))$$

$\overbrace{X_1^{-1}(\{x_1\}) \cap X_2^{-1}(\{x_2\})}$

$$= \sum_{x_1} \sum_{x_2} g(x_1, x_2) P(X_1=x_1, X_2=x_2) \text{ as required. } \square$$

Some of the examples above showed the usefulness of indicator functions in computing expectation and variance.

Let's look at some more.

probability p probability $1-p$

Consider a chain of N beads composed of two colours: A and B. Let $X: \Omega \rightarrow \mathbb{R}$ count the number of "alternates".

Alternates are defined by a sequence "AB" or "BA", so "ABAB" has 3 alternates. Find $E[X]$ and $\text{Var}(X)$.

Define a random variable $Y: \Omega \rightarrow \mathbb{R}^N$ by $Y((w_1, \dots, w_N)) = (Y_1(w_1, w_2), \dots, Y_{N-1}(w_{N-1}, w_N), 0)$.

Each Y_i is an indicator random variable $Y_i(w_i, w_{i+1}) = \begin{cases} 1 & \text{if } w_i \neq w_{i+1} \text{ (alternates).} \\ 0 & \text{otherwise.} \end{cases}$

$$\text{Notice } X = \sum_{i=1}^N Y_i \text{ so } E[X] = E\left[\sum_{i=1}^N Y_i\right] = \sum_{i=1}^N E[Y_i] = \sum_{i=1}^{N-1} E[Y_i] + E[Y_N].$$

$$\text{We have } E[Y_i] = \sum_x x P(Y_i=x) = \sum_{x=0}^1 x P(Y_i=x) = P[Y_i=1]$$

$$\text{where } P[Y_i=j] = P[Y_j] \text{ for all } i, j, \text{ and } P[Y_i=1] = P[w_{i+1}=A | w_i=B] P[w_i=B] = p(1-p) + (1-p)p.$$

$$+ P[w_{i+1}=B | w_i=A] P[w_i=A]$$

$$\Rightarrow E[X] = \sum_{i=1}^{N-1} E[Y_i] = (N-1)(2p(1-p))$$

$$\text{We have } \text{Var}(X) = \text{Var}\left(\sum_{i=1}^N Y_i\right) = \sum_{i=1}^N \text{Var}(Y_i) + 2 \sum_{i < j} \text{Cov}(Y_i, Y_j).$$

$$\text{Here, we have } \text{Var}(Y_i) = E[Y_i^2] - E[Y_i]^2 = \sum_{x=0}^1 x^2 P[Y_i=x] - (2p(1-p))^2 = P[Y_i=1] - (2p(1-p))^2$$

$$= 2p(1-p) - (2p(1-p))^2.$$

$$\text{and } \text{Cov}(Y_i, Y_j) = E[Y_i Y_j] - E[Y_i]E[Y_j]. \text{ Notice that } E[Y_i Y_j] = \sum_x \sum_y x y P[Y_i=x, Y_j=y]$$

$$= P(Y_i=1, Y_j=1)$$

$$= P(w_i \neq w_{i+1}, w_j \neq w_{j+1})$$

↑ independent if $i+1 \neq j$
otherwise dependent

Some of the above examples rely on the following notion.

Consider discrete random variables $X, Y: \Omega \rightarrow \mathbb{R}$. Recall for an event $B \in \mathcal{A}$, we had $E(X|B) = \frac{1}{P(B)} E[\mathbb{1}_B(X)]$.

$$\text{Similarly, we have } E[X|Y=y] = \sum_x x P(X=x|Y=y)$$

↓ some function of y

$$= \sum_x x \frac{P(X=x, Y=y)}{P(Y=y)}$$

↑ joint probability mass function

Notice $E[X|Y=y] = f(y): \mathbb{R} \rightarrow \mathbb{R}$ and $Y: \Omega \rightarrow \mathbb{R}$ so $f(Y): \Omega \xrightarrow{Y} \mathbb{R} \xrightarrow{f} \mathbb{R}$ defines a new random variable.

Then $f(Y)$ is called the conditional expectation of X given Y denoted $E[X|Y]$.

Remark: Notice $E[X|Y]$ is a random variable, not a value.

Example: Consider a die throw, and a random variable $X: \Omega \rightarrow \mathbb{R}$ checking if the roll is even, so $X(\omega) = \begin{cases} 1 & \text{if } \omega \in \{2, 4, 6\}, \\ 0 & \text{otherwise.} \end{cases}$

$Y: \Omega \rightarrow \mathbb{R}$ checking if we rolled a six, so $Y(\omega) = \begin{cases} 1 & \text{if } \omega=6, \\ 0 & \text{otherwise.} \end{cases}$

Then $E[X|Y]: \Omega \rightarrow \mathbb{R}$ is a random variable and $E[X|Y](\omega) = \begin{cases} E[X|Y=0], & \omega \in \{1, \dots, 5\}; \\ E[X|Y=1], & \omega \in \{6\}. \end{cases}$

Midterm Practice:

$$1a) \quad P[X|Y=1] = \frac{P[X, Y=1]}{P[Y=1]} \quad \text{where } P[Y=1] = P[X=1, Y=1] + P[X=2, Y=1].$$

$$= \frac{2}{8} = \frac{1}{4}$$

$$\text{For } x=1: \quad P[X=1, Y=1] = 4 \cdot \frac{1}{8} \quad \text{and} \quad \text{for } x=2: \quad P[X=2, Y=1] = 4 \cdot \frac{1}{8}.$$

$$f(1|1) \quad P[Y=1] \qquad \qquad \qquad f(2|1) \quad P[Y=1]$$

$$P[X|Y=2] = \frac{P[X, Y=2]}{P[Y=2]} \quad \text{where} \quad P[Y=2] = \frac{1}{2} + \frac{1}{4} = \frac{3}{4}$$

$$\text{For } x=1: \quad P[X=1, Y=2] = \frac{4}{3} \cdot \frac{1}{4} = \frac{1}{3} \quad \text{and} \quad \text{for } x=2: \quad P[X=2, Y=2] = \frac{4}{3} \cdot \frac{1}{2} = \frac{2}{3}.$$

$$f(1|2) \quad P[Y=2] \qquad \qquad \qquad f(2|2) \quad P[Y=2]$$

b) X and Y are not independent because $f(2|1) = \frac{1}{2} \neq f(2) = \frac{5}{8}$.

c) $P[X+Y \leq 3]$. Define $S=XY$ so $S \in \{0, 1, 2\}$. Notice $P[S \leq 3] = 1 - P[S=4]$

$$= 1 - P[X=2, Y=2] = 1 - \frac{1}{8} = \frac{7}{8}$$

$P[X+Y > 1]$. Define $S=XY$ so $S \in \{0, 1, 2\}$. We have $P[S > 1] = P[S=2] = P[X=2, Y=1] = \frac{1}{8}$.

$$P[X+Y > 2] = 1 - P[X+Y \leq 2] = 1 - P[X=1, Y=1] = 1 - \frac{1}{8} = \frac{7}{8}.$$

3 a) Roll a fair die 100 times independently. Denote $X_i =$ the number of i 's, for $i \in \{1, \dots, 6\}$.

Find the p.m.f. of $T_0 = X_1 + X_3 + X_5$.

We have T_0 the count of odd die rolls, so each roll has $P_{\text{odd}} = \frac{1}{2}$.

Thus we have $T_0 \sim \text{Binomial}(n, p_{\text{odd}}) = \text{Binomial}(100, \frac{1}{2})$ so $P_{T_0}(x) = \binom{100}{x} \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{100-x}$

b) Find the p.m.f. of X_5 given $T_0 = t_0$.

We have to find $P[X_5 | X_1 + X_3 + X_5 = t_0]$. Consider $S = X_1 + X_3$ so $X_5 + S = t_0$.

\Rightarrow find $P[X_5 | X_5 + S = t_0]$.

Continuous Distributions

We have developed a theory of probability above. We will now introduce several useful distributions.

1 Exponential: $X \sim \text{Exp}(\lambda)$ with $f(x) = \lambda e^{-\lambda x}$.

a) Show that X follows cdf $F(x) = \begin{cases} 1 - e^{-\lambda x} & \text{for } x > 0; \\ 0 & \text{for } x \leq 0. \end{cases}$

• models lifetime and/or decay.

• models "wait time" by considering λ to be a Poisson parameter (so X represents waiting time between events).

$$\text{so } F_X(x) = P[X \leq x] = 1 - P[X > x] = 1 - P[\text{no events in } [0, x]] = 1 - \frac{e^{-\lambda x} (\lambda x)^0}{0!} = 1 - e^{-\lambda x} \quad \text{for all } x > 0.$$

• "memoryless" property: $P[X > t+s | X > s] = P[X > t]$ (probability of a "leap" only depends on "leap" size, not its location).

*Exercise

2 Gamma function: $\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$, and $X \sim \text{Gamma}(\alpha, \beta)$ with $f(x) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta}$.

$\Gamma(n+1) = n!$ for $n \in \mathbb{N}$,
 $n \Gamma(n) = \Gamma(n+1)$.

*Exercise (Integration by parts)

↑ normalization constant

• extension of $\text{Exp}(\lambda)$. When $\alpha=1$, we have $\text{Gamma}(1, \beta) = \text{Exp}(\beta)$ $\Rightarrow \text{Gamma}(1, \frac{1}{\lambda}) = \text{Exp}(\lambda)$.

3 Uniform distribution.

4 Beta-distribution: $\beta(\alpha_1, \alpha_2) = \int_0^1 x^{\alpha_1-1} (1-x)^{\alpha_2-1} dx = \frac{\Gamma(\alpha_1) \Gamma(\alpha_2)}{\Gamma(\alpha_1 + \alpha_2)}$. If $X \sim \text{Beta}(\alpha_1, \alpha_2)$ then

$$f(x) = \frac{1}{\beta(\alpha_1, \alpha_2)} x^{\alpha_1-1} (1-x)^{\alpha_2-1} \quad \text{for } x \in [0, 1].$$

5 Normal distribution: $X \sim N(0, 1)$ with $f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ (standard normal).

• if $X \sim N(0, 1)$ then $F_X(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$ (no closed form).

6 (General) Normal distribution $Z \sim N(\mu, \sigma^2)$ with $f(z) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(z-\mu)^2}{2\sigma^2}}$ (general normal).

• if $X \sim N(\mu, \sigma^2)$ then $F_X(x) = \int_{-\infty}^x \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$ (no closed form).

Proposition (Uniformity): Let X be a random variable such that $F_X(x)$ is strictly increasing (so F_X' exists).

Let $Y \sim \text{Unif}(0, 1)$. Then $F_X^{-1}(Y) \stackrel{\text{cdf}}{\sim} F_X(x)$.

Proof: (rough) We have $F_X(x) = P[X \leq x]$. Define $Y \sim \text{Unif}(0, 1)$ so $X = F_X^{-1}(Y)$.

$$\text{Then } P[X \leq x] = P[F_X^{-1}(Y) \leq x] = P[Y \leq F_X(x)] \stackrel{\substack{Y \in [0, 1] \\ \text{since c.d.f.}}}{=} F_X(x). \quad \square$$

Proposition ($N(0, 1)$): Let $X \sim N(0, 1)$. Then $E[X] = 0$ and $\text{Var}(X) = 1$.

Proof: By definitions. We have $E[X] = \int_{-\infty}^{\infty} x f(x) dx = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \frac{1}{\sqrt{2\pi}} \left(\int_{-\infty}^0 x e^{-\frac{x^2}{2}} dx + \int_0^{\infty} x e^{-\frac{x^2}{2}} dx \right) \xrightarrow{\text{by symmetry}} 0$.

$$\text{and } \text{Var}(X) = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2 = \int_{-\infty}^{\infty} x^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \xrightarrow{\text{compute}} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^2 e^{-\frac{x^2}{2}} dx = \frac{1}{\sqrt{2\pi}} \sqrt{2\pi} = 1. \quad \square$$

Proposition (Reparametrization Trick): Let $X \sim N(\mu, \sigma^2)$. Then $Z = \frac{X-\mu}{\sigma} \sim N(0, 1)$.

Proof: Notice $\frac{X-\mu}{\sigma}$ defines a bijection (in particular, a homeomorphism), so the inverse transform exists and is unique.

$$\text{We have } Z = \frac{X-\mu}{\sigma} \Rightarrow X = \sigma Z + \mu \text{ so } F_Z(z) = P(Z \leq z) \xrightarrow{\text{inverse transform}} = P\left(\frac{X-\mu}{\sigma} \leq z\right) = P(X \leq \sigma z + \mu).$$

$$\text{Therefore } f_Z(z) = \frac{d}{dz} F_Z(z) = \frac{d}{dz} P(X \leq \sigma z + \mu) = \frac{d}{dz} F_X(\sigma z + \mu) = f_X(\sigma z + \mu) \frac{d}{dz} (\sigma z + \mu) \xrightarrow{\text{how?}} = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{\sigma z + \mu - \mu}{\sigma}\right)^2} \cdot \sigma = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} \text{ a standard normal.} \quad \square$$

Corollary ($N(\mu, \sigma^2)$): Let $X \sim N(\mu, \sigma^2)$. Then $E[X] = \mu$ and $\text{Var}(X) = \sigma^2$.

Proof: Let $Y \sim N(0, 1)$ and consider $X = \sigma Y + \mu$ by the reparametrization trick.

$$\text{Then } E[X] = E[\sigma Y + \mu] = \sigma E[Y] + \mu = \mu.$$

$$\text{Var}(X) = \text{Var}(\sigma Y + \mu) = \sigma^2 \text{Var}(Y) = \sigma^2. \quad \square$$

Gamma Function Technique

Recall that if $X \sim \text{Gamma}(\alpha, \beta)$ where $\alpha, \beta > 0$, then $f_X(x) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta}$ for $x \in [0, \infty)$.

Since f_X defines a probability density, we have

$$\int_0^{\infty} \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta} dx = 1 \Rightarrow \int_0^{\infty} x^{\alpha-1} e^{-x/\beta} dx = \beta^\alpha \Gamma(\alpha). \text{ This defines the Gamma function technique.}$$

Corollary: Let $X \sim \text{Gamma}(\alpha, \beta)$. Then $E[X^p] = \frac{\beta^p \Gamma(p+\alpha)}{\Gamma(\alpha)}$.

Proof: We use the Γ -technique.

$$E[X^p] = \int_0^{\infty} x^p \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta} dx = \frac{1}{\beta^\alpha \Gamma(\alpha)} \int_0^{\infty} x^{p+\alpha-1} e^{-x/\beta} dx = \frac{1}{\beta^\alpha \Gamma(\alpha)} \beta^{p+\alpha} \Gamma(p+\alpha) = \frac{\beta^p \Gamma(p+\alpha)}{\Gamma(\alpha)}. \quad \square$$

Corollary: Let $X \sim \text{Gamma}(\alpha, \beta)$. Then $\text{Var}(X) = \alpha\beta^2$.

Proof: $\text{Var}(X) = E[X^2] - E[X]^2 = \beta^2 \frac{\Gamma(2+\alpha)}{\Gamma(\alpha)} - \left(\beta \frac{\Gamma(1+\alpha)}{\Gamma(\alpha)}\right)^2 = \beta^2 (\alpha+1)(\alpha) - (\alpha\beta)^2 = \alpha\beta^2 + \alpha^2\beta^2 - \alpha^2\beta^2 = \alpha\beta^2$. \square

* **Exercise:** Let $Z \sim N(\mu, \sigma^2)$. Use the reparametrization trick and Γ -technique to find $E[X]$ and $\text{Var}(X)$.

* **Exercise:** Use the Γ -technique to find $E[X]$ and $\text{Var}(X)$ where: 1) $X \sim \text{Unif}(a, b)$ 2) $X \sim \text{Beta}(a, b)$.

Moment Generating Functions

Previously, we defined the notion of a **moment generating function** for a random variable X .

We now develop a theory of moment generating functions.

* **Exercise:** Let $X \sim \text{Bin}(n, p)$. Find $M_X(t)$ and find the values of t such that $M_X(t)$ is well-defined.

Exercise: Let $X \sim \text{Gamma}(\alpha, \beta)$. Find $M_X(t)$ and find the values of t such that $M_X(t)$ is well-defined.

$$\begin{aligned} \text{We have } M_X(t) = E[e^{tx}] &= \int_0^\infty e^{tx} \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta} dx = \frac{1}{\beta^\alpha \Gamma(\alpha)} \int_0^\infty x^{\alpha-1} e^{-x/\beta + tx} dx \quad (\text{let } k = \frac{1}{\beta^\alpha \Gamma(\alpha)}) \\ &= k \int_0^\infty x^{\alpha-1} e^{-x(\frac{1}{\beta} - t)} dx \\ &= k \int_0^\infty x^{\alpha-1} \exp(-\frac{x}{(\frac{1}{\beta} - t)}) dx. \text{ Define } \tilde{\beta} = (\frac{1}{\beta} - t)^{-1}, \text{ so} \\ k \int_0^\infty x^{\alpha-1} \exp(-\frac{x}{(\frac{1}{\beta} - t)}) dx &= k \int_0^\infty x^{\alpha-1} e^{-x/\tilde{\beta}} dx \\ &\stackrel{\substack{\text{Γ-technique on } \tilde{\beta}}}{=} k \tilde{\beta}^\alpha \Gamma(\alpha) = \frac{1}{\beta^\alpha \Gamma(\alpha)} (\frac{1}{\beta} - t)^{-\alpha} \Gamma(\alpha) \\ &= \beta^{-\alpha} (\beta^{-1} - t)^{-\alpha} = (\beta \beta^{-1} - \beta t)^{-\alpha} = (1 - \beta t)^{-\alpha}. \end{aligned}$$

We used the Γ -technique above, which restricts our values of t , since we require $\alpha > 0$ and $\tilde{\beta} > 0 \Rightarrow (\frac{1}{\beta} - t) > 0 \Rightarrow t < \frac{1}{\beta}$.

Would we get weaker constraints if we did not use the Γ function? We verify that our bounds for t are necessary.

$$\begin{aligned} \text{Say } t \geq \frac{1}{\beta}. \text{ If } t = \frac{1}{\beta}, \text{ then } E[e^{tx}] &= \int_0^\infty e^{tx} k x^{\alpha-1} e^{-x/\beta} dx \\ &= k \int_0^\infty x^{\alpha-1} e^{-x/\beta + x/\beta} dx = k \int_0^\infty x^{\alpha-1} dx \rightarrow \infty. \end{aligned}$$

$$\begin{aligned} \text{If } t > \frac{1}{\beta}. \text{ Then } E[e^{tx}] &= \int_0^\infty e^{tx} k x^{\alpha-1} e^{-x/\beta} dx = k \int_0^\infty x^{\alpha-1} e^{x(t-\frac{1}{\beta})} dx \\ &\geq \underbrace{\min_{x \in [0, \infty]} e^{x(t-\frac{1}{\beta})}}_1 \int_0^\infty x^{\alpha-1} dx \rightarrow \infty. \end{aligned}$$

Therefore $M_X(t) = (1 - \beta t)^{-\alpha}$ for $t < \frac{1}{\beta}$ and since $\beta > 0$ is finite, we can find $\varepsilon > 0$ such that $M_X(t)$ is continuous on $B(0, \varepsilon) \Rightarrow E[X^n] = M_X^{(n)}(0)$ is well-defined.

Exercise: Let $X \sim N(0, 1)$. Find $M_X(t)$.

$$\begin{aligned} \text{We have } M_X(t) &= \int_{-\infty}^\infty e^{tx} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^\infty e^{-\frac{1}{2}(x-t)^2} e^{\frac{1}{2}t^2} dx = e^{\frac{t^2}{2}} \int_{-\infty}^\infty \underbrace{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-t)^2}}_1 dx = e^{\frac{t^2}{2}}. \end{aligned}$$

since pdf. of $N(t, 1)$

Notice the argument works for any $t \in \mathbb{R} \Rightarrow M_X(t)$ well-defined for all $t \in \mathbb{R}$.

Proposition (MGF of a linear function): Let X be a random variable and let $Y = aX + b$ for $a, b \in \mathbb{R}$.

Say $M_X(t)$ exists in $B(0, \varepsilon)$. Then $M_Y(t) = e^{bt} M_X(at)$ for $t \in B(0, \frac{\varepsilon}{|a|})$.

Proof: Follow the computation. $M_Y(t) = E[e^{tY}] = E[e^{t(ax+b)}] = E[e^{bt} e^{atx}] \xrightarrow{\text{homogeneity}} e^{bt} E[e^{atx}] = e^{bt} M_X(at)$.

Then $M_Y(t)$ exists if $M_X(at)$ exists $\Rightarrow at \in B(0, \varepsilon) \Rightarrow t \in B(0, \frac{\varepsilon}{|a|})$. \square

Example: Say $X \sim N(\mu, \sigma^2)$. Find $M_X(t)$ using the reparametrization trick.

Notice $X = \sigma Z + \mu$ where $Z \sim N(0, 1)$. By the above Proposition, we have $M_Z(t) = e^{\frac{t^2}{2}} \Rightarrow M_X(t) = e^{\mu t} M_Z(\frac{t}{\sigma}) = e^{\mu t} e^{\frac{t^2}{2(\sigma^2)}} \text{ for } t \in \mathbb{R}$.

Exercise: Say $X \sim \text{Gamma}(\alpha, \beta)$. Let $Y = X/\beta$. Find $M_Y(t)$ and state the values of t for which this is well-defined.

Let X and Y be random variables. We define the joint moment generating function $M_{(X,Y)}(t_1, t_2) = E[e^{(t_1 X + t_2 Y)}] = E[e^{t_1 X + t_2 Y}]$.

Given a joint MGF $M_{(X,Y)}(t_1, t_2)$, how can we marginalize to find $M_X(t_1)$ or $M_Y(t_2)$?

Corollary: If $M_{(X,Y)}(t_1, t_2)$ is a joint MGF, then $M_X(t_1) = M_{(X,Y)}(t_1, 0)$ and $M_Y(t_2) = M_{(X,Y)}(0, t_2)$.

Proof: Obvious. Since $M_{(X,Y)}: \mathbb{R}^2 \rightarrow \mathbb{R}$ well-defined in $B(0, \varepsilon) \subseteq \mathbb{R}^2$, we have $M_{(X,Y)}(t_1, 0)$ and $M_{(X,Y)}(0, t_2)$ well-defined. \square

The definition of a joint MGF (and the above Corollary) generalize nicely to an arbitrary family $(X_i)_{i \in \Lambda}$ of random variables.

Theorem (Joint MGFs): Let X, Y be random variables with joint MGF $M_{(X,Y)}(t_1, t_2)$. Then $E[X^j Y^k] = \left. \frac{\partial^{j+k}}{\partial t_1^j \partial t_2^k} M_{(X,Y)}(t_1, t_2) \right|_{(t_1, t_2)=0}$.

Proof: X, Y discrete. Then $M_{(X,Y)}(t_1, t_2) = E[e^{t_1 X + t_2 Y}] = \sum_x \sum_y e^{t_1 x + t_2 y} f(x, y)$

$$\Rightarrow \frac{\partial^{j+k}}{\partial t_1^j \partial t_2^k} M_{(X,Y)}(t_1, t_2) = \frac{\partial^{j+k}}{\partial t_1^j \partial t_2^k} \sum_x \sum_y e^{t_1 x + t_2 y} f(x, y) \xrightarrow{\substack{\text{finite sums} \\ \text{commute}}} \sum_x \sum_y \frac{\partial^{j+k}}{\partial t_1^j \partial t_2^k} e^{t_1 x + t_2 y} f(x, y) \xrightarrow{\substack{\text{independent of} \\ t_1 \text{ and } t_2}}$$

$$= \sum_x \sum_y x^j e^{t_1 x} y^k e^{t_2 y} f(x, y)$$

$$\xrightarrow{\substack{t_1=0 \\ t_2=0}} = \sum_x \sum_y x^j y^k f(x, y) = E[X^j Y^k]. \text{ Done. } \Delta$$

multivariate
series

X, Y continuous. Then $M_{(X,Y)}(t_1, t_2) = E[e^{t_1 X + t_2 Y}] = \iint_D e^{t_1 x} e^{t_2 y} f(x, y) dx dy$

$$\frac{\partial^{j+k}}{\partial t_1^j \partial t_2^k} M_{(X,Y)}(t_1, t_2) = \frac{\partial^{j+k}}{\partial t_1^j \partial t_2^k} \iint_D e^{t_1 x} e^{t_2 y} f(x, y) dx dy \xrightarrow{*} \iint_D \frac{\partial^{j+k}}{\partial t_1^j \partial t_2^k} e^{t_1 x} e^{t_2 y} f(x, y) dx dy$$

$$\xrightarrow{\substack{t_1=0 \\ t_2=0}} = \iint_D x^j y^k f(x, y) dx dy = E[X^j Y^k].$$

multivariate
series

Corollary: Let X, Y be random variables. Then $E[XY] = \left. \frac{\partial^2}{\partial t_1 \partial t_2} M_{(X,Y)}(t_1, t_2) \right|_{(t_1, t_2)=0}$.

Theorem (MGF Independence): X and Y are independent $\Leftrightarrow M_{X,Y}(t_1, t_2) = M_X(t_1)M_Y(t_2)$ for all $(t_1, t_2) \in B(0, \varepsilon) \subseteq \mathbb{R}^2$.

Proof (\Rightarrow): Let X and Y be independent. Then $M_{X,Y}(t_1, t_2) = E[e^{t_1 X} e^{t_2 Y}] \stackrel{\text{Independence}}{=} E[e^{t_1 X}] E[e^{t_2 Y}] = M_X(t_1) M_Y(t_2)$. Δ

(\Leftarrow): Now say X, Y are random variables such that $M_{X,Y}(t_1, t_2) = M_X(t_1)M_Y(t_2)$ for all $(t_1, t_2) \in B(0, \varepsilon)$.

Then $M_{X,Y}(t_1, t_2) = E[e^{t_1 X} e^{t_2 Y}] = E[e^{t_1 X}] E[e^{t_2 Y}] = M_X(t_1) M_Y(t_2)$.

Denote $g(X) = e^{t_1 X}$ and $h(Y) = e^{t_2 Y}$.

X, Y discrete. Then $E[g(X) h(Y)] = \sum_x \sum_y g(x) h(y) f(x, y)$

$\stackrel{\text{by assumption}}{=} (\sum_x g(x) f(x)) (\sum_y h(y) f(y))$. Suppose both sums are finite, so

$$\Rightarrow \sum_{i=1}^n \sum_{j=1}^m g(x_i) h(y_j) f(x_i, y_j) = \left(\sum_{i=1}^n g(x_i) f(x_i)\right) \left(\sum_{j=1}^m h(y_j) f(y_j)\right) = \sum_{i=1}^n \sum_{j=1}^m g(x_i) h(y_j) f(x_i) f(y_j)$$

$$\Rightarrow f(x_i, y_j) = f(x_i) f(y_j) \text{ for all } i, j.$$

Therefore $F_X(x_i) F_Y(y_j) = F_{X,Y}(x_i, y_j)$ so X and Y are independent (by definition). Δ

The other cases are out of scope. \square

This gives us an (easier) way to check for independence: find $M_{X,Y}(t_1, t_2) \Rightarrow$ derive M_X and M_Y , check factorization.

Theorem (Uniqueness of MGF): If random variables X and Y satisfy $M_X = M_Y$, then $X, Y \sim f$ are identical.

Example: Recall $Z \sim \text{Gamma}(\alpha, \beta) \Rightarrow M_Z(t) = (1 - \beta t)^{-\alpha}$. If $M_X(t_1) = (1 - t_1)^{-1}$ for $t_1 < 1 \Rightarrow X \sim \text{Gamma}(1, 1)$.

If $M_Y(t_2) = (1 - t_2)^{-2}$ for $t_2 < 1 \Rightarrow Y \sim \text{Gamma}(2, 1)$.

Notice $E[e^{t_1 X} e^{t_2 Y}] = M_{X+Y}(t)$ and $M_{X+Y}: \mathbb{R} \rightarrow \mathbb{R}$ is fundamentally different from $M_{X,Y}: \mathbb{R}^2 \rightarrow \mathbb{R}$.

Corollary: Let X, Y be independent random variables. Then $M_{X+Y}(t) = M_X(t)M_Y(t)$.

Proof: We have $M_{X+Y}(t) = E[e^{t(X+Y)}] = E[e^{tX} e^{tY}] = E[e^{tX}] E[e^{tY}] = M_X(t) M_Y(t)$. \square

Example: Say $X_1 \sim \text{Bin}(n, p)$ and $X_2 \sim \text{Bin}(m, p)$ are independent. Show $X_1 + X_2 \sim \text{Bin}(n+m, p)$ using MGF uniqueness.

We have $M_{X_1}(t) = (1 + pe^t - p)^n$ and $M_{X_2}(t) = (1 + pe^t - p)^m$

$$\Rightarrow M_{X_1 + X_2}(t) = M_{X_1}(t) M_{X_2}(t) = (1 + pe^t - p)^n (1 + pe^t - p)^m = (1 + pe^t - p)^{n+m} \stackrel{\text{MGF of } \text{Bin}(n+m, p)}{\Rightarrow} X_1 + X_2 \sim \text{Bin}(n+m, p).$$

Example: Say $X_1 \sim \text{Gamma}(\alpha_1, \beta)$ and $X_2 \sim \text{Gamma}(\alpha_2, \beta)$. Show that $X_1 + X_2 \sim \text{Gamma}(\alpha_1 + \alpha_2, \beta)$.
(done below)

Example: Say $X_i \sim N(\mu_i, \sigma_i^2)$ for $i = 1, \dots, n$ are all independent. Show that $c_0 + \sum_{i=1}^n c_i X_i \sim N(c_0 + \sum_{i=1}^n c_i \mu_i, \sum_{i=1}^n c_i^2 \sigma_i^2)$.
(done below)

Convergence

We introduce the notion of convergence in two ways.

Let $(X_i)_{i \in \mathbb{N}}$ be a sequence of random variables.

We say $(X_i)_{i \in \mathbb{N}}$ converges in probability to a random variable $X \Leftrightarrow \lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0$ } strong convergence

and write $X_n \xrightarrow{P} X$. An "unusual" outcome becomes less likely as the sequence progresses. } pointwise

We say $(X_i)_{i \in \mathbb{N}}$ converges in distribution to a random variable $X \Leftrightarrow \lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$ at all continuity points $x \in \mathbb{R}$ of F_X } weak convergence

and write $X_n \xrightarrow{D} X$.

Let $T \subseteq \mathbb{R}$ be an arbitrary set (interpreted as the set of time points). For each $t \in T$, define $X_t: \Omega \rightarrow \mathbb{R}$ a random variable.

Then $(X_t)_{t \in T}$ is called a stochastic process.

For $w \in \Omega$, the map $T \rightarrow \mathbb{R}$ is called a path.
 $t \mapsto X_t(w)$

Let $(X_t)_{t \in T}$ be a stochastic process with $T \subseteq \mathbb{Z}$ or $T \subseteq \mathbb{R}$.

We call $(X_t)_{t \in T}$ a Markov process (or a Markov chain) if for all $n \in \mathbb{N}$, $t_1, t_2, \dots, t_n, t \in T$ with $t_1 < t_2 < \dots < t_n < t$
 and $x_1, \dots, x_n, x \in \mathbb{R}$, we have $P(X_t=x | X_{t_1}=x_1, \dots, X_{t_n}=x_n) = P(X_t=x | X_{t_n}=x_n)$.
 set of time points
 ordered in time
 Markov property

In other words, the probability of the "future" only depends on the present.

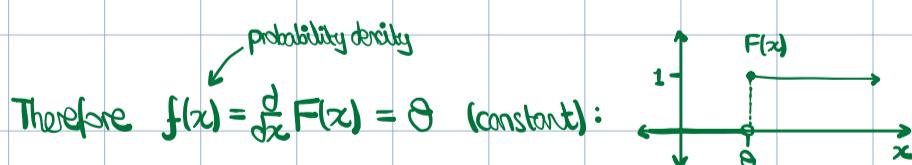
If $X_n \xrightarrow{\text{def}} F_n(x)$ and $X \xrightarrow{\text{def}} F(x)$ with $X_n \xrightarrow{D} X$, then $F_n(x) \approx F(x)$ at all continuity points of F . Useful approximation.

For notation, we write $X_i \xrightarrow{\text{IID}} F$ if X_i are independent and identically distributed random variables.

Example: Say $Y_i \xrightarrow{\text{IID}} \text{Unif}(0, \theta)$. Let $X_n = \max\{Y_1, \dots, Y_n\}$. Find X such that $X_n \xrightarrow{D} X$.

$$\begin{aligned} \text{Notice } X_1 &= \max\{Y_1\} \text{ . What's } \lim_{n \rightarrow \infty} F_n(x) \text{? Notice } F_n(x) = P(X_n \leq x) = P(\max\{Y_1, \dots, Y_n\} \leq x) \\ &\stackrel{\text{defn}}{=} P(\max\{Y_1, \dots, Y_n\} \leq x) \rightarrow P(Y_1 \leq x, \dots, Y_n \leq x) \\ &\stackrel{\text{independence}}{=} \prod_{i=1}^n P(Y_i \leq x) \stackrel{\text{identically distributed}}{=} \prod_{i=1}^n F_i(x) = (F_i(x))^n = \begin{cases} 0^n \text{ if } x \leq 0 \\ (\frac{x}{\theta})^n \text{ if } 0 < x < \theta \\ 1^n \text{ if } x \geq \theta \end{cases} \\ &\Rightarrow \lim_{n \rightarrow \infty} F_n(x) = \begin{cases} 0 \text{ for } x \leq 0 \\ 0 \text{ for } 0 < x < \theta \\ 1 \text{ for } x \geq \theta \end{cases} \end{aligned}$$

this convergence is not uniform!



and notice that $F_n(x) \rightarrow F(x)$ wherever F continuous $\Rightarrow F_n(x) \rightarrow F(x)$ on $\mathbb{R} \setminus \{0\}$. Thus $X_n \xrightarrow{D} X = \theta$.

Theorem: Let $(X_n)_{n \in \mathbb{N}}$ be a sequence/stochastic process. If $X_n \xrightarrow{P} X$, then $X_n \xrightarrow{D} X$.
 (strong convergence) (weak convergence)

Proof: We can only prove this for $X=a$ (constant) without heavy measure theory.

let $X=a \in \mathbb{R}$. Notice $F_x(x)$ is a step function continuous on $\mathbb{R} \setminus \{a\}$, and $F_x(x)$ is constant on $\mathbb{R} \setminus \{a\}$.

$\lim F_n(x) = F(x)$ on $\mathbb{R} \setminus \{a\}$.

We show $\forall \varepsilon > 0$, we have $F(x-\varepsilon) \leq \lim_{n \rightarrow \infty} F_n(x) \leq F(x+\varepsilon)$ for all $x \in \mathbb{R} \setminus \{0\}$.

$$\begin{aligned}
 (1) \quad F_n(x) &= P(X_n \leq x) = P(X_n \leq x, X \leq x+\varepsilon) + P(X_n \leq x, X > x+\varepsilon) \\
 &\stackrel{\text{upper bound}}{\leq} P(X \leq x+\varepsilon) + P(X - X_n \geq \varepsilon) \stackrel{\text{superset}}{\leq} P(X \leq x+\varepsilon) + P(|X_n - X| \geq \varepsilon) \\
 &= F(x+\varepsilon) + P(|X_n - X| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} F(x+\varepsilon). \\
 (2) \quad F(x-\varepsilon) &= P(X < x-\varepsilon, X_n \leq x) + P(X < x-\varepsilon, X_n > x) \\
 &\stackrel{\text{lower bound}}{\leq} P(X_n \leq x) + P(X_n - X \geq \varepsilon) \stackrel{\text{constant}}{\leq} P(X_n \leq x) + P(|X_n - X| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} \lim_{n \rightarrow \infty} F_n(x).
 \end{aligned}$$

Thus $F(x-\varepsilon) \leq \lim_{n \rightarrow \infty} F_n(x) \leq F(x+\varepsilon)$ for any $\varepsilon > 0$. Holds for all $x \in \mathbb{R} \setminus \{0\}$. \square

Theorem: Let $(X_n)_{n \in \mathbb{N}}$ be a sequence/stochastic process. If $X_n \xrightarrow{D} X = C$, then $X_n \xrightarrow{P} X = C$ at all continuity points of F_X .

Remark: We have $X_n \xrightarrow{P} X = C \Leftrightarrow X_n \xrightarrow{D} X = C$.

Proof: Notice $X = C \in \mathbb{R} \Rightarrow F_X$ continuous on $\mathbb{R} \setminus \{C\} \Rightarrow C$ not a continuity point of F_X .

Thus $\forall \varepsilon > 0$, we have $P(|X_n - X| \geq \varepsilon) = P(|X_n - C| \geq \varepsilon)$

$$\begin{aligned}
 &= P(X_n - C \geq \varepsilon) + P(X_n - C \leq -\varepsilon) \\
 &= P(X_n \geq C + \varepsilon) + P(X_n \leq C - \varepsilon) \quad \text{and we bound this above:} \\
 &\leq 1 - P(X_n \leq C + \frac{\varepsilon}{2}) \quad \begin{matrix} F_{X_n}(C - \varepsilon) \\ \text{by definition} \end{matrix} \\
 &\leq 1 - F_{X_n}(C + \frac{\varepsilon}{2}) + F_{X_n}(C - \varepsilon) \xrightarrow{n \rightarrow \infty} 1 - F_X(C + \frac{\varepsilon}{2}) + F_X(C - \varepsilon) \\
 &\quad \begin{matrix} \text{since } C + \frac{\varepsilon}{2} \text{ and } C - \varepsilon \text{ are} \\ \text{continuity points of } F_X \\ \text{and } X_n \xrightarrow{P} X. \end{matrix} \\
 &= 1 - 1 + 0 = 0.
 \end{aligned}$$

Therefore $0 \leq P(|X_n - X| \geq \varepsilon) \leq 1 - F_{X_n}(C + \frac{\varepsilon}{2}) + F_{X_n}(C - \varepsilon) \xrightarrow{n \rightarrow \infty} 0$

$\Rightarrow P(|X_n - X| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0$. Thus $X_n \xrightarrow{P} X$. \square

Remark: In general, note that $X_n \xrightarrow{D} X \not\Leftrightarrow X_n \xrightarrow{P} X$.

Example: Say $X \sim N(0, 1)$ and the sequence $(X_i)_{i \in \mathbb{N}}$ is defined by $X_i = -X$. We show $X_n \xrightarrow{D} X$ but $X_n \not\xrightarrow{P} X$.

Notice by symmetry, $X_i = -X \Rightarrow f_{X_i}(x) = f_X(-x) = \frac{1}{\sqrt{2\pi}} e^{\frac{(-x)^2}{2}} = \frac{1}{\sqrt{2\pi}} e^{\frac{x^2}{2}} = f_X(x)$ (could alternatively use MGF uniqueness to prove that $X_i \sim N(0, 1)$)

$\Rightarrow X_i \sim N(0, 1)$. Thus obviously $X_n \xrightarrow{D} X$, since $f_{X_i} = f_X \Rightarrow F_{X_i} = F_X$

$\Rightarrow F_{X_n}(x) \xrightarrow{n \rightarrow \infty} F_X(x)$ for all $x \in \mathbb{R}$.

However, $X_n \xrightarrow{P} X \Leftrightarrow \forall \varepsilon > 0$, we have $\lim_{n \rightarrow \infty} P(|X_n - X| \geq \varepsilon) = 0$, but

$= \lim_{n \rightarrow \infty} P(|-X - X| \geq \varepsilon) = P(2|X| \geq \varepsilon) \neq 0$ for $\varepsilon = 1$ (for example).

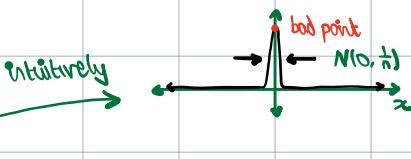
Therefore $X_n \not\xrightarrow{P} X$.

Notice that $X_n \xrightarrow{D} X$ restricts convergence to continuity points of F_X .

Example: Let $X_i \stackrel{\text{IID}}{\sim} N(0, 1)$ and define $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. Then $\bar{X}_n \xrightarrow{D} X = 0$ but $\lim_{n \rightarrow \infty} F_{\bar{X}_n}(0) \neq F_X(0)$

Let $F_n = F_{\bar{X}_n}$. Then $\lim_{n \rightarrow \infty} F_n(x) = \lim_{n \rightarrow \infty} P(\bar{X}_n \leq x)$.

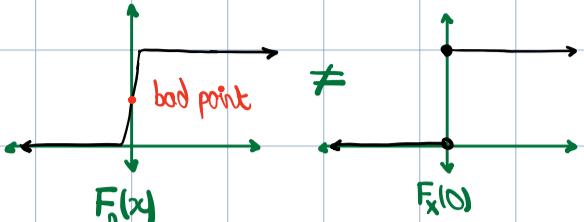
Notice $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$



(since 0 $\in \mathbb{R}$ not a continuity point).

$$\Rightarrow \bar{X}_n \sim N\left(\frac{1}{n} \sum_{i=1}^n 0, \frac{1}{n^2} \sum_{i=1}^n 1\right) = N\left(0, \frac{1}{n}\right). \text{ By the reparametrization trick, } \frac{\bar{X}_n - 0}{\sqrt{\frac{1}{n}}} \sim N(0, 1), \text{ so}$$

$$F_n(x) = P(\bar{X}_n \leq x) \stackrel{\text{reparam.}}{=} P\left(\frac{\bar{X}_n - 0}{\sqrt{\frac{1}{n}}} \leq \frac{x - 0}{\sqrt{\frac{1}{n}}}\right) = P(Z \leq \sqrt{n}x) \text{ where } Z \sim N(0, 1)$$



$$= F_Z(\sqrt{n}x) \xrightarrow{n \rightarrow \infty} \begin{cases} 0 & \text{if } x < 0 \\ \frac{1}{2} & \text{if } x = 0 \\ 1 & \text{if } x > 0 \end{cases} \text{ so } \lim_{n \rightarrow \infty} F_n(0) = \frac{1}{2} \neq 1 = F_x(0) \text{ where } X=0.$$

continuity points

$$\text{Notice } F_x(x) = \begin{cases} 0 & \text{if } x < 0 \\ \frac{1}{2} & \text{if } x = 0 \\ 1 & \text{if } x > 0 \end{cases} \text{ so } F_n(x) \xrightarrow{n \rightarrow \infty} F_x(x) \text{ for } x \in \mathbb{R} \setminus \{0\}.$$

Example: Let $Y_i \stackrel{\text{iid}}{\sim} \text{Unif}(0, \theta)$ and let $X_n = \max\{Y_1, \dots, Y_n\}$. Show $X_n \xrightarrow{P} \theta$ and find the distribution limit.

Intuitively, $\max\{Y_1, \dots, Y_n\} \xrightarrow{n \rightarrow \infty} \max[0, \theta] = \theta$. So let $X = \theta$. We show $X_n \xrightarrow{P} X = \theta$.

It suffices to prove $X_n \xrightarrow{P} X = \theta$ (strong convergence).

By definition, $\forall \varepsilon > 0$, we want $\lim_{n \rightarrow \infty} P(|X_n - \theta| \geq \varepsilon) = 0$. since $\theta \geq X_n$

$$\begin{aligned} \text{Notice } P(|X_n - \theta| \geq \varepsilon) &= P(|\max\{Y_1, \dots, Y_n\} - \theta| \geq \varepsilon) = P(\theta - \max\{Y_1, \dots, Y_n\} \geq \varepsilon) \\ &\stackrel{\text{independence}}{=} P(\max\{Y_1, \dots, Y_n\} \leq \theta - \varepsilon) = P(Y_1 \leq \theta - \varepsilon, \dots, Y_n \leq \theta - \varepsilon) \\ &\stackrel{\text{identical}}{=} \prod_{i=1}^n P(Y_i \leq \theta - \varepsilon) = (P(Y_1 \leq \theta - \varepsilon))^n \stackrel{n \rightarrow \infty}{=} \begin{cases} 0 & \text{if } \theta - \varepsilon < 0 \\ (\frac{\theta - \varepsilon}{\theta})^n & \text{if } 0 \leq \theta - \varepsilon < \theta \\ 1 & \text{if } \theta - \varepsilon \geq 0 \end{cases} \\ &\quad \downarrow \\ &\quad \lim_{n \rightarrow \infty} (1 - \frac{\varepsilon}{\theta})^n = 0 \text{ for } 0 \leq \theta - \varepsilon < \theta \end{aligned}$$

so $P(|X_n - \theta| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0$ for all $\varepsilon > 0$. Thus $X_n \xrightarrow{P} X = \theta \Rightarrow X_n \xrightarrow{P} X = \theta$. \square

Proposition (e-limit): If $b \in \mathbb{R}$ and $\lim_{n \rightarrow \infty} \gamma(n) = 0$, then $\lim_{n \rightarrow \infty} (1 + \frac{b}{n} + \frac{\gamma(n)}{n})^n = e^b$.

Example: Let $Y_i \stackrel{\text{iid}}{\sim} \text{Exp}(1)$. let $X_n = \max\{Y_1, \dots, Y_n\} - \log n$. Show that $X_n \xrightarrow{P} \theta$ and find the distribution limit.

Notice $f_{Y_i}(y) = e^{-y}$ for $y \geq 0$

$\Rightarrow F_{Y_i}(y) = \begin{cases} 0 & \text{for } y \leq 0 \\ 1 - e^{-y} & \text{for } y > 0 \end{cases}$. Let $F_n(x) = F_{X_n}(x)$ be the c.d.f. of X_n .

Then $F_n(x) = P(X_n \leq x) = P(\max\{Y_1, \dots, Y_n\} - \log n \leq x)$

$$\begin{aligned} &\stackrel{\text{independence}}{=} P(\max\{Y_1, \dots, Y_n\} \leq x + \log n) = P(Y_1 \leq x + \log n, \dots, Y_n \leq x + \log n) \\ &\stackrel{\text{identical}}{=} \prod_{i=1}^n P(Y_i \leq x + \log n) = (P(Y_i \leq x + \log n))^n \\ &= (F_{Y_i}(x + \log n))^n = \begin{cases} 0^n & \text{for } x + \log n \leq 0 \\ (1 - e^{-(x + \log n)})^n & \text{for } x + \log n > 0 \end{cases} \end{aligned}$$

$$\text{Notice } (1 - e^{-(x + \log n)})^n = (1 - e^{-x} e^{-\log n})^n = (1 - e^{-x} e^{\log n^{-1}})^n$$

$$= \left(1 - \frac{e^{-x}}{n}\right)^n \xrightarrow{n \rightarrow \infty} e^{-e^{-x}}. \text{ Notice } \{x \in \mathbb{R} \mid x + \log n > 0\} \xrightarrow{n \rightarrow \infty} \mathbb{R}.$$

$$\Rightarrow \lim_{n \rightarrow \infty} F_n(x) = e^{-e^{-x}} = F(x) \text{ for all } x \in \mathbb{R}.$$

We must verify that $F(x)$ defines a valid cumulative distribution function:

$$\lim_{x \rightarrow -\infty} F(x) = \lim_{x \rightarrow -\infty} e^{-e^{-x}} \rightarrow e^{-e^{\infty}} = e^{-\infty} = 0.$$

$$\lim_{x \rightarrow \infty} F(x) = \lim_{x \rightarrow \infty} e^{-e^{-x}} \rightarrow e^{-e^{-\infty}} = e^{-0} = 1.$$

$F(x)$ non-decreasing: satisfied by $-e^{-x}$ increasing.

Therefore $F(x)$ is a well-defined c.d.f. $\Rightarrow X_n \xrightarrow{D} X \stackrel{\text{def}}{\sim} e^{-e^{-x}}$ for $x \in \mathbb{R}$. \square

Let $(X_i)_{i \in \mathbb{N}}$ be a sequence/stochastic process with $(M_{X_i}(t))_{i \in \mathbb{N}}$ their moment generating functions.

Notice, for fixed $t \in \mathbb{R}$, that $(M_{X_i}(t))_{i \in \mathbb{N}} \subseteq \mathbb{R}$ is a real sequence.

Theorem (MGF convergence): Let $(X_i)_{i \in \mathbb{N}}$ be a sequence of random variables, and $(M_{X_i}(t))_{i \in \mathbb{N}}$ their MGFs, all well-defined in some common neighbourhood $B(0, \varepsilon) \subseteq \mathbb{R}$. Suppose $M_{X_i}(t) \rightarrow M_X(t)$ for $t \in B(0, \varepsilon)$. Then $X_n \xrightarrow{D} X$.

Proof: Heavy measure theory. \square

We will use this to prove the Central Limit Theorem.

Example: Let $X_n \sim \text{Bin}(n, p)$. Show that $X_n \xrightarrow{D}$ and find the limit such that $\lambda = \stackrel{\text{constant}}{np}$. Use the limit to approximate $F_{X_n}(x)$.

Since $np = \lambda$ is a constant, we have $X_n \sim \text{Bin}(n, \frac{\lambda}{n})$. We postulate that $X_n \xrightarrow{D} X$ where $X \sim \text{Poi}(\lambda)$.

$$\text{Notice } M_{X_n}(t) = (1 + pe^t - p)^n = \left(1 + \frac{\lambda}{n}e^t - \frac{\lambda}{n}\right)^n = \left(1 + \frac{\lambda}{n}(e^t - 1)\right)^n \xrightarrow[n \rightarrow \infty]{\text{e-limit}} e^{\lambda(e^t - 1)}$$

$$\text{and notice that } M_X(t) = e^{\lambda(e^t - 1)} \text{ for } t \in \mathbb{R}$$

$$\Rightarrow \lim_{n \rightarrow \infty} M_{X_n}(t) = M_X(t) \text{ for } t \in \mathbb{R} \Rightarrow X_n \xrightarrow{D} X \sim \text{Poi}(\lambda). \quad \square$$

How can we approximate $F_{X_n}(x)$? Notice $X_n \xrightarrow{D} X$ are both discrete, so F_X has continuity points $\mathbb{R} \setminus \mathbb{Z}_{\geq 0}$.

Therefore, we don't necessarily have $\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$ for $x \in \mathbb{Z}_{\geq 0}$.

However, for $0 < \varepsilon < 1$ we have that $x - \varepsilon, x + \varepsilon \in \mathbb{R} \setminus \mathbb{Z}_{\geq 0}$ for $x \in \mathbb{Z}_{\geq 0}$.

$$\Rightarrow \lim_{n \rightarrow \infty} F_{X_n}(x \pm \varepsilon) = F_X(x \pm \varepsilon) \text{ so } P(X_n = x) = F_{X_n}(x + \varepsilon) - F_{X_n}(x - \varepsilon) \xrightarrow{n \rightarrow \infty} F_X(x + \varepsilon) - F_X(x - \varepsilon) = P(X = x).$$

continuity points

Theorem (Central Limit): Let X_i be independent, identically distributed random variables with $E[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2 < \infty$.

Then $\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{D} N(0, 1)$ where $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ (sample average).

Proof: For this course, we must add one (weak) assumption: $M_{X_i}(t)$ exist in some $B(0, \varepsilon)$. We use MGF convergence.

Define $Z_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}$. We want $Z_n \xrightarrow{D} Z \sim N(0, 1)$. Denote $M_n = M_{Z_n}$ and $M = M_Z$.

Suffices to show $M_n(t) \rightarrow M(t)$ for $t \in B(0, \varepsilon^*)$.

We have $M(t) = e^{\frac{1}{2}t^2}$ for all $t \in \mathbb{R}$, so we want $\lim_{n \rightarrow \infty} M_n(t) = e^{\frac{1}{2}t^2}$.

$$\begin{aligned} \text{Notice } Z_n &= \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} = \frac{\sqrt{n}\left(\frac{1}{n} \sum_{i=1}^n X_i - \mu\right)}{\sigma} = \frac{\sqrt{n}}{n} \left(\sum_{i=1}^n X_i - n\mu\right) \sigma^{-1} \\ &= \frac{1}{\sqrt{n}} \left(\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma}\right)\right) \text{ so } Z_n \text{ is a linear combination of standardized } X_i \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i \text{ where } Y_i = \frac{X_i - \mu}{\sigma} \sim N(0, 1). \end{aligned}$$

Important: $Y_i \sim N(0, 1) \Rightarrow E[Y_i] = 0$ and $E[Y_i^2] = \text{Var}(Y_i) + E[Y_i]^2 = \text{Var}(Y_i) = 1$.

$$\Rightarrow M'_{Y_i}(0) = 0 \text{ and } M''_{Y_i}(0) = 1.$$

Computing $M_n(t)$.

$$\text{We have } Z_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i \Rightarrow M_n(t) = E[\exp\left(\frac{t}{\sqrt{n}} \sum_{i=1}^n Y_i\right)] \text{ where } X_i \stackrel{\text{ iid}}{\sim} N(\mu, \sigma^2) \Rightarrow Y_i \stackrel{\text{iid}}{\sim} N(0, 1).$$

$$= E\left[\prod_{i=1}^n \exp\left(\frac{t}{\sqrt{n}} Y_i\right)\right]$$

$$= \prod_{i=1}^n E[\exp(\frac{t}{\sqrt{n}} Y_i)] = \left(M_{Y_i}(\frac{t}{\sqrt{n}})\right)^n \text{ where } M_{Y_i}(\frac{t}{\sqrt{n}}) = E[\exp(\frac{t}{\sqrt{n}} Y_i)] = E[\exp(\frac{t}{\sqrt{n}}(X_i - \mu))]$$

exists for $\frac{t}{\sqrt{n}} \in B(0, \varepsilon)$
 $\Rightarrow t \in B(0, \varepsilon \sqrt{n}) \xrightarrow{n \rightarrow \infty} \mathbb{R}$.

Taylor Expansion of $M_{Y_i}(\frac{t}{\sqrt{n}})$ at $t=0$.

$$M_{Y_i}(\frac{t}{\sqrt{n}}) = M_{Y_i}(0) + M'_{Y_i}(0)(\frac{t}{\sqrt{n}}) + \frac{1}{2!} M''_{Y_i}(0)(\frac{t}{\sqrt{n}})^2 + O(n^{-\frac{3}{2}}).$$

1 0 1

higher order terms

$$= 1 + \frac{1}{n} (\frac{1}{2} t^2) + O(n^{-\frac{3}{2}}) \Rightarrow M_n(t) = \left(M_{Y_i}(\frac{t}{\sqrt{n}})\right)^n = \left(1 + \frac{1}{n} \frac{t^2}{2} + O(n^{-\frac{3}{2}})\right)^n$$

around $t=0$ where $O(n^{-\frac{3}{2}}) = \frac{O(n^{-\frac{1}{2}})}{n}$ and $O(n^{-\frac{1}{2}}) \xrightarrow{n \rightarrow \infty} 0$.

e-limit

$$e^{\frac{t^2}{2}} \text{ for } t \in B(0, \varepsilon^*).$$

Therefore $M_n(t) \rightarrow M(t)$ for $t \in B(0, \varepsilon^*)$, so by **MGF convergence**, we have $\bar{Z}_n \xrightarrow{D} Z \sim N(0, 1)$. \square

We define the χ distribution as follows: say $X_i \stackrel{\text{iid}}{\sim} \text{Gamma}(1, \beta)$. Define χ^2 as $\frac{2 \sum_{i=1}^n X_i}{\beta} \sim \chi^2(2n)$.

If $X \sim \chi(n)$, then $f_X(x) = \frac{x^{k-1} e^{-\frac{x^2}{2}}}{2^{k/2-1} \Gamma(k/2)}$ for $x \geq 0$.

Example: Say $Y_n \sim \chi(n)$ for $n \in \mathbb{N}_{\geq 0}$. Show that $\bar{Z}_n = \frac{Y_n - n}{\sqrt{2n}} \xrightarrow{D} Z \sim N(0, 1)$.

We have $Y_n = \sum_{i=1}^n X_i$ where $X_i \stackrel{\text{iid}}{\sim} \text{Gamma}(\frac{1}{2}, 2)$.

Therefore since $E[X_i] = \frac{1}{2} \cdot 2 = 1 < \infty$ and $\text{Var}(X_i) = \frac{1}{2} \cdot 2^2 = 2 < \infty$.

$$\xrightarrow{\text{CLT}} \frac{\sqrt{n}(\bar{X}_n - 1)}{\sqrt{2}} \xrightarrow{D} N(0, 1) \quad \text{so} \quad \frac{\sqrt{n}(\frac{1}{n} \sum_{i=1}^n X_i - 1)}{\sqrt{2}} = \frac{1}{\sqrt{2n}} (\sum_{i=1}^n X_i - n) = \frac{Y_n - n}{\sqrt{2n}} \xrightarrow{D} N(0, 1).$$

Example: Fires are reported to a station following a Poisson process, with an average of 1 fire per 4 hours.

Find the probability that the 500th fire is reported by the end of the 84th day.

Let X_i be the waiting time between the $(i-1)^{\text{th}}$ and $(i)^{\text{th}}$ fire.

Then $X_i \stackrel{\text{iid}}{\sim} \text{Exp}(\frac{1}{6})$ and $f_{X_i}(x) = f_i(x) = \frac{1}{6} e^{-\frac{x}{6}}$ for $x > 0$.

Notice $E[X_i] = \frac{1}{6} < \infty$ and $\text{Var}(X_i) = \frac{1}{36} < \infty$ so

$$\xrightarrow{\text{CLT}} \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{D} N(0, 1) \quad \text{so with } n=500, \mu=\frac{1}{6}, \sigma=\frac{1}{6}, \text{ we want } \sum_{i=1}^{500} X_i \leq 84.$$

$$\text{We have } \sqrt{500} \left(\frac{1}{500} \sum_{i=1}^{500} X_i - \frac{1}{6} \right) / \left(\frac{1}{6} \right) = \sqrt{500} \left(\frac{\sum X_i}{500} - \frac{1}{6} \right) \cdot 6 \approx Z \sim N(0, 1)$$

$$\text{and want } P\left(\sqrt{500} \left(\frac{\sum X_i}{500} - \frac{1}{6} \right) \leq \sqrt{500} \left(\frac{84}{500} - \frac{1}{6} \right)\right) \text{ by reparametrization.}$$

$$\begin{array}{ccc} \xrightarrow{s} & & \\ N(0, 1) & & 0.18 \\ \xrightarrow{z} & & 0.18 \\ & & \approx P(Z \leq 0.18). \text{ Use a lookup table.} \end{array}$$

Proposition (Markov's Inequality): For any random variable X with $E[|X|^k] < \infty$, we have $P[|X| \geq c] \leq \frac{E[|X|^k]}{c^k}$ for all $c > 0$, for all $k \in \mathbb{Z}_{>0}$.

Proof: Say X_i are continuous. Similar proof for X_i discrete.

$$\begin{aligned} \text{We have } \frac{E[|X|^k]}{c^k} &= \frac{1}{c^k} \int_{-\infty}^{\infty} |x|^k f_X(x) dx \text{ and } R = \{|x| \geq c\} \cup \{|x| < c\}: \\ &= \frac{1}{c^k} \left(\underbrace{\int_{|x| \geq c} |x|^k f_X(x) dx + \int_{|x| < c} |x|^k f_X(x) dx}_{\geq 0} \right) \\ &\geq \frac{1}{c^k} \int_{|x| \geq c} |x|^k f_X(x) dx \geq \frac{1}{c^k} \min_{\substack{x \in \{|x| \geq c\} \\ x=c}} |x|^k \int_{|x| \geq c} f_X(x) dx = \int_{|x| \geq c} f_X(x) dx = P(|X| \geq c). \quad \square \end{aligned}$$

Example: If $Y_n = \bar{X}_n \sim N(\mu, \frac{\sigma^2}{n})$, then $Y_n \xrightarrow{P} \mu$.

Notice by Markov's Inequality that $E[|Y_n|^k] c^{-k} \geq P[|Y_n| \geq c]$ for all $c > 0$.

$$\begin{aligned} &\Rightarrow P(|Y_n - \mu| \geq c) \leq E[|Y_n - \mu|^k] c^{-k} \\ &\Rightarrow P(|Y_n - \mu| \geq \varepsilon) \leq E[|Y_n - \mu|^k] \varepsilon^{-k} = M_{Y_n - \mu}^{(k)}(0) \text{ for all } k \in \mathbb{Z}_{>0}. \text{ Let } k=2: \\ &\Rightarrow P(|Y_n - \mu| \geq \varepsilon) \leq E[|Y_n - \mu|^2] \varepsilon^{-2} = E[(Y_n - \mu)^2] \varepsilon^{-2} = (\text{Var}(Y_n - \mu) + E[Y_n - \mu]^2) \varepsilon^{-2} = \text{Var}(Y_n) \varepsilon^{-2} = \frac{\sigma^2}{n\varepsilon^2}. \\ \text{Thus } P(|Y_n - \mu| \geq \varepsilon) &\leq \frac{\sigma^2}{n\varepsilon^2} \xrightarrow{n \rightarrow \infty} 0 \text{ so } Y_n \xrightarrow{P} \mu \text{ by definition.} \end{aligned}$$

Chebychev's inequality

Theorem (Weak Law of Large Numbers): Suppose X_i are IID with $E[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2 < \infty$, then $\bar{X}_n \xrightarrow{P} \mu$

where $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ (average).

Proof: By independence, we have $\text{Var}(\bar{X}_n) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}$.

Using Chebychev's inequality: $P(|\bar{X}_n - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{n\varepsilon^2} \xrightarrow{n \rightarrow \infty} 0$

so $\bar{X}_n \xrightarrow{P} \mu$ by definition. Done. \square

Example: Say $X_1 \sim \text{Gamma}(\alpha_1, \beta) \perp X_2 \sim \text{Gamma}(\alpha_2, \beta)$. Show that $X_1 + X_2 \sim \text{Gamma}(\alpha_1 + \alpha_2, \beta)$.

We have that $M_{X_1+X_2}(t) = M_{X_1}(t)M_{X_2}(t)$ by independence.

$$\Rightarrow M_{X_1+X_2}(t) = M_{X_1}(t)M_{X_2}(t) = (1-\beta t)^{-\alpha_1} (1-\beta t)^{-\alpha_2} = (1-\beta t)^{-(\alpha_1+\alpha_2)} = M_Y(t) \text{ where } Y \sim \text{Gamma}(\alpha_1 + \alpha_2, \beta)$$

$$\Rightarrow X_1 + X_2 \sim \text{Gamma}(\alpha_1 + \alpha_2, \beta) \text{ by MGF uniqueness. } \square$$

Example: Say $X_i \sim N(\mu_i, \sigma_i^2)$ for $i=1, \dots, n$ are all independent. Show that $c_0 + \sum_{i=1}^n c_i X_i \sim N(c_0 + \sum_{i=1}^n c_i \mu_i, \sum_{i=1}^n c_i^2 \sigma_i^2)$.

We use MGF uniqueness. Notice that

$$M_{c_0 + \sum_{i=1}^n c_i X_i}(t) = M_{c_0}(t) \prod_{i=1}^n M_{c_i X_i}(t). \text{ We have } M_{c_0}(t) = E[e^{c_0 t}] = e^{c_0 t} \text{ and}$$

$$M_{c_i X_i}(t) = E[e^{c_i X_i t}] = E[e^{X_i(c_i t)}] = M_{X_i}(c_i t).$$

Write $X_i = \sigma_i Z + \mu_i$ where $Z \sim N(0, 1)$. Then $M_{X_i}(c_i t) = M_{\sigma_i Z + \mu_i}(c_i t) \stackrel{\substack{\text{Prop.} \\ \mu_i t}}{=} e^{\mu_i t} M_Z(\sigma_i c_i t) = e^{\mu_i t} e^{\frac{1}{2}(\sigma_i c_i t)^2}$.

$$\Rightarrow M_{c_0 + \sum_{i=1}^n c_i X_i}(t) = M_{c_0}(t) \prod_{i=1}^n M_{c_i X_i}(t) = e^{c_0 t} \prod_{i=1}^n e^{\mu_i t} e^{\frac{1}{2}(\sigma_i c_i t)^2} = e^{t(c_0 + \sum_{i=1}^n \mu_i)} e^{\frac{1}{2} \sum_{i=1}^n (\sigma_i c_i t)^2} = e^{t(c_0 + \sum_{i=1}^n \mu_i)} e^{\frac{1}{2} t^2 \sum_{i=1}^n (\sigma_i c_i)^2} \stackrel{\text{MGF uniqueness}}{\sim} N(c_0 + \sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2 c_i^2).$$

Notice if $Y \sim N(\mu, \sigma^2)$, then $Y = \sigma Z + \mu \Rightarrow M_Y(t) = M_{\sigma Z + \mu}(t) = e^{\mu t} M_Z(\sigma t) = e^{\mu t} e^{\frac{1}{2}(\sigma t)^2}$.

Example: Say $Y_n \sim \text{Bin}(n, \theta)$. Show that $\frac{Y_n}{n} \xrightarrow{P} \theta$.

Method 1: Markov's Inequality. We show $P(|\frac{Y_n}{n} - \theta| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0$.

Define $X_n = \frac{Y_n}{n}$. Then $P(|X_n - \theta| \geq \varepsilon) \leq \frac{E[(X_n - \theta)^k]}{\varepsilon^k}$ for any $\varepsilon > 0$, $k \in \mathbb{Z}_{>0}$. Choose $\varepsilon = 2$ and $k = 2$.

$$(E[X_n] = E[\frac{Y_n}{n}] = n\theta \frac{1}{n} = \theta)$$

$$\Rightarrow P(|X_n - \theta| \geq \varepsilon) \leq \frac{E[(X_n - \theta)^2]}{\varepsilon^2} = \frac{E[(X_n - \theta)^2]}{\varepsilon^2} \varepsilon^{-2} = \frac{\text{Var}(X_n)}{\varepsilon^2} \varepsilon^{-2} = \frac{\theta(1-\theta)}{n} \varepsilon^{-2} \xrightarrow{n \rightarrow \infty} 0.$$

Var(X_n-θ) = E(X_n-θ)² - (E(X_n-θ))²

Thus $X_n \xrightarrow{P} \theta$.

Method 2: Indicator functions and WLLN.

Define $Y_n = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$ where $X_i \sim \text{Bin}(1, \theta)$ are indicator functions, so $Y_n = \sum_{i=1}^n X_i$ and $\frac{Y_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i$.

Notice X_i are IID with $E[X_i] = \theta$ and $\text{Var}(X_i) = \theta(1-\theta) < \infty$.

$$\xrightarrow{\text{WLLN}} \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \theta. \text{ Thus } \frac{Y_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \theta.$$

Exercise: Say $X \sim \text{Gamma}(\alpha, \beta)$. Let $Y = X/\beta$. Find $M_Y(t)$ and state the values of t for which this is well-defined.

$$\text{We have } M_Y(t) = M_{X/\beta}(t) = E[e^{tX/\beta}] = M_X\left(\frac{t}{\beta}\right) = \left(1 - \beta\left(\frac{t}{\beta}\right)\right)^{-\alpha} = (1-t)^{-\alpha}.$$

$M_X(t)$ defined for $t < \frac{1}{\beta} \Rightarrow M_Y(t/\beta)$ defined for $t/\beta < \frac{1}{\beta} \Rightarrow t < 1$.

Exercise: Let $Z \sim N(\mu, \sigma^2)$. Use the reparametrization trick and Γ -technique to find $E(X)$ and $\text{Var}(X)$.

$$\text{Easy. } Z = \sigma X + \mu \text{ where } X \sim N(0, 1) \Rightarrow M_Z(t) = M_{\sigma X + \mu}(t) = e^{\mu t} M_X(\sigma t) = e^{\mu t} e^{\frac{\sigma^2 t^2}{2}}$$

Exercise: Use the Γ -technique to find $E[X]$ and $\text{Var}(X)$ where: 1) $X \sim \text{Unif}(a, b)$ 2) $X \sim \text{Beta}(a, b)$.

1) Let $X \sim \text{Unif}(a, b)$. Then $E[X] = \int_a^b x f_x(x) dx = \int_a^b x \frac{1}{b-a} dx = \frac{1}{b-a} \left(\frac{1}{2}x^2\right)_a^b = \frac{1}{b-a} \frac{1}{2}(b^2 - a^2) = \frac{b+a}{2}$ (as expected). How could we have used Γ -technique?

$$\text{Var}(X) = \int_a^b x^2 f_x(x) dx - E[X]^2 = \frac{1}{b-a} \left(\frac{1}{3}x^3\right)_a^b.$$

2) Let $X \sim \text{Beta}(a, b)$. Then $f_x(x) = \frac{1}{B(a,b)} x^{a-1} (1-x)^{b-1}$ for $x \in [0, 1]$

$$\Rightarrow E[X] = \int_0^1 x f_x(x) dx = \frac{1}{B(a,b)} \int_0^1 x^a (1-x)^{b-1} dx. \text{ Let } \alpha_1 = a+1 \text{ and } \alpha_2 = b, \text{ so}$$

$$= \frac{1}{B(\alpha_1-1, \alpha_2)} \int_0^1 x^{\alpha_1-1} (1-x)^{\alpha_2-1} dx = \frac{\Gamma(\alpha_1-1+\alpha_2)}{\Gamma(\alpha_1-1)\Gamma(\alpha_2)} \frac{\Gamma(\alpha_1)\Gamma(\alpha_2)}{\Gamma(\alpha_1+\alpha_2)}$$

Γ -technique

$$= \frac{\Gamma(\alpha_1+\alpha_2-1)}{\Gamma(\alpha_1-1)} \frac{(\alpha_1-1)\Gamma(\alpha_1-1)}{(\alpha_1+\alpha_2-1)\Gamma(\alpha_1+\alpha_2-1)} = \frac{(\alpha_1-1)}{(\alpha_1+\alpha_2-1)} = \frac{a}{a+b}.$$

$$\Rightarrow \text{Var}(X) = \int_0^1 x^2 f_x(x) dx - E[X]^2 \text{ where } \int_0^1 x^2 f_x(x) dx = \frac{1}{B(a,b)} \int_0^1 x^{a+1} (1-x)^b dx$$

$$= B(a,b)^{-1} \int_0^1 x^{\alpha_1-1} (1-x)^{\alpha_2} dx \text{ where } \alpha_1 = a+2 \text{ and } \alpha_2 = b.$$

$$= B(\alpha_1-2, \alpha_2)^{-1} \frac{\Gamma(\alpha_1)\Gamma(\alpha_2)}{\Gamma(\alpha_1+\alpha_2)} \text{ by the } \Gamma\text{-technique, and so}$$

$$= \frac{\Gamma(\alpha_1-2+\alpha_2)}{\Gamma(\alpha_1-2)\Gamma(\alpha_2)} \frac{\Gamma(\alpha_1)\Gamma(\alpha_2)}{\Gamma(\alpha_1+\alpha_2)} = \frac{\Gamma(\alpha_1-2+\alpha_2)}{\Gamma(\alpha_1-2)} \frac{(\alpha_1-1)(\alpha_1-2)\Gamma(\alpha_1-2)}{(\alpha_1+\alpha_2-1)(\alpha_1+\alpha_2-2)\Gamma(\alpha_1+\alpha_2-2)}$$

$$= \frac{(\alpha_1-1)(\alpha_1-2)}{(\alpha_1+\alpha_2-1)(\alpha_1+\alpha_2-2)} = \frac{(a+1)(a)}{(a+b+1)(a+b)} *$$

so $\text{Var}(X) = \frac{(a+1)(a)}{(a+b+1)(a+b)} - \left(\frac{a}{a+b}\right)^2$.

*: Can we use MGF to find $E[X^2]$ instead? We have $E[X^2] = M_X^{(2)}(0)$

and $M_X(t) = E[e^{tx}] = \int_0^1 e^{tx} f_x(x) dx$ leibniz by continuity on $[0,1] \Rightarrow$ uniform continuity

$$\Rightarrow M_X^{(2)}(t) = \frac{d^2}{dt^2} \int_0^1 e^{tx} f_x(x) dx = \int_0^1 \frac{\partial^2}{\partial t^2} e^{tx} f_x(x) dx$$

$$= \int_0^1 x^2 e^{tx} f_x(x) dx$$

$$= B(a,b)^{-1} \int_0^1 e^{tx} x^{a+1} (1-x)^b dx \text{ and it's unclear how to continue.}$$

Exercise: let $X \sim \text{Bin}(n, p)$. Find $M_X(t)$ and find the values of t such that $M_X(t)$ is well-defined.

We have $M_X(t) = E[e^{tx}] = \sum_{x=0}^n e^{tx} P_X(x) = \sum_{x=0}^n e^{tx} \cdot \binom{n}{x} p^x (1-p)^{n-x}$

$$= \sum_{x=0}^n \binom{n}{x} (pe^t)^x (1-p)^{n-x}$$

$$= (pe^t + (1-p))^n = (pe^t + 1 - p)^n \text{ by Binomial formula, for } t \in \mathbb{R}.$$

10.4 12) Say X has a discrete uniform distribution on $\{a, a+1, \dots, b\}$ with $P_X(x) = \frac{1}{b-a+1}$ for $x \in \{a, \dots, b\}$.

a) Find $M_X(t)$.

We have $M_X(t) = E[e^{tx}] = \sum_{x=a}^b e^{tx} P_X(x) = \sum_{x=a}^b e^{tx} \frac{1}{b-a+1} = \frac{1}{b-a+1} \sum_{x=a}^b e^{tx} = \frac{1}{b-a+1} \frac{e^{at}(1-e^{(b-a+1)t})}{1-e^t}$

$$= \frac{1}{b-a+1} \frac{e^{at}-e^{(b+1)t}}{1-e^t}.$$

b) Use $M_X(t)$ to find $E[X]$ and $E[X^2]$.

$E[X] = M'_X(0) = \frac{d}{dt} \left(\frac{1}{b-a+1} \frac{e^{at}-e^{(b+1)t}}{1-e^t} \right) \Big|_{t=0}$ by l'Hopital.

Similarly for $M''_X(0)$.

10.4 13) $\text{Im}(X) = \{0, 1, 2\}$, $E[X] = 1$, $E[X^2] = 1.5$.

a) We have $E[X] = \sum_{x=0}^2 x P_x(x) = P_x(1) + 2P_x(2) = 1$.

$$E[X^2] = \sum_{x=0}^2 x^2 P_x(x) = P_x(1) + 4P_x(2) = 1.5. \text{ Thus } 2P_x(2) = 0.5 \Rightarrow P_x(2) = \frac{1}{4} \text{ and } P_x(1) = \frac{1}{2} \\ \Rightarrow P_x(0) = \frac{1}{4}.$$

$$\text{Thus } M_X(t) = E[e^{tx}] = \sum_{x=0}^2 e^{tx} P_x(x) = e^{0t} \cdot \frac{1}{4} + e^{1t} \cdot \frac{1}{2} + e^{2t} \cdot \frac{1}{4}.$$

$$b) E[X^3] = M_X^{(3)}(0) = \frac{d^3}{dt^3} \left(\frac{1}{4}e^{2t} + \frac{1}{2}e^t + \frac{1}{4} \right)_{t=0} = (2e^{2t} + \frac{1}{2}e^t)_{t=0} = 2 + \frac{1}{2} = \frac{5}{2}.$$

$$E[X^4] = M_X^{(4)}(0) = \frac{d^4}{dt^4} \left(\frac{1}{4}e^{2t} + \frac{1}{2}e^t + \frac{1}{4} \right)_{t=0} = (4e^{2t} + \frac{1}{2}e^t)_{t=0} = \frac{9}{2}.$$

c) If $E[X]=a$ and $E[X^2]=b$, then $E[X] = a \cdot P_x(1) + 2 \cdot P_x(2)$ and $E[X^2] = b \cdot P_x(1) + 4 \cdot P_x(2)$

$$\Rightarrow P_x(2) = \frac{1}{2}(b-a) \text{ and } P_x(1) = a - (b-a) = 2a-b \text{ and } P_x(0) = 1 - \left(\frac{1}{2}(b-a) + 2(a-b) \right).$$

Therefore $P_x(x)$ fully determined.

10.4 16) Let $X_i \stackrel{\text{ind}}{\sim} N(0, 1)$.

c) Let $S_n = \sum_{i=1}^n X_i$. Find $M_{S_n}(t)$, and thus the distribution.

$$\text{We have } M_{S_n}(t) = E[e^{tS_n}] = E[e^{t \sum_{i=1}^n X_i}] = E[\prod_{i=1}^n e^{tX_i}] \\ \xrightarrow{\text{independence}} = \prod_{i=1}^n E[e^{tX_i}] \text{ where } E[e^{tX_i}] = M_{X_i}(t) = e^{\frac{t^2}{2}} \text{ by previous results} \\ \xrightarrow{\text{identical}} = (e^{\frac{t^2}{2}})^n = e^{nt^2}.$$

Notice $Z \sim N(\mu, \sigma^2) \Rightarrow M_Z(t) = M_{\mu + \sigma Z}(t) = e^{\mu t} M_\sigma(\sigma t) = e^{\mu t} e^{\frac{\sigma^2 t^2}{2}}$ so $\mu=0$ and $\sigma^2=n$, so

by uniqueness of MGF, we have $S_n \sim N(0, n)$.

d) $Z = n^{-1/2}(S_n - n)$. Then $Z = \frac{1}{\sqrt{n}}S_n - \sqrt{n}$ and $S_n = nX$ where $X \sim N(0, 1)$

$$\Rightarrow Z = \sqrt{n}X - \sqrt{n}. \text{ Therefore } M_Z(t) = M_{\sqrt{n}X - \sqrt{n}}(t) = e^{\sqrt{n}t} M_X(\sqrt{n}t) = e^{\sqrt{n}t} e^{\frac{n t^2}{2}} \text{ so } \mu = \sqrt{n} \text{ and } \sigma^2 = n$$

by uniqueness of MGF, we have $Z \sim N(\sqrt{n}, n)$.

10.4 18) Say X continuous with $f_X(x) = \frac{1}{\theta^2} x e^{-x/\theta}$ for $x > 0$, $\theta > 0$.

a) Find $M_X(t)$.

$$\text{We have } M_X(t) = E[e^{tx}] = \int_0^\infty e^{tx} \frac{1}{\theta^2} x e^{-x/\theta} dx = \frac{1}{\theta^2} \int_0^\infty x e^{tx} e^{-x/\theta} dx \\ = \theta^{-2} \int_0^\infty x e^{-x(\theta-t)} dx \quad (\int_0^\infty x^{\alpha-1} e^{-x/\beta} dx = \beta^\alpha \Gamma(\alpha)) \\ \xrightarrow{\alpha=2, \beta=(\theta-t)^{-1}} = \theta^{-2} \int x^{\alpha-1} e^{-x/\beta} dx = \theta^{-2} \beta^\alpha \Gamma(\alpha) \text{ for } \beta > 0 \Rightarrow (\theta-t)^{-1} > 0 \\ \Rightarrow \theta-t > 0 \\ \Rightarrow t < \theta \text{ and } \theta > 0 \\ \Rightarrow M_X(t) \text{ exists in } B(0, \theta). \\ = \theta^{-2} (\theta-t)^{-2} \Gamma(2) \\ = \theta^{-2} (\theta-t)^{-2}$$

b) Say $X \sim \text{Exp}(\theta) \perp Y \sim \text{Exp}(\theta)$. Find the distribution of $S = X+Y$.

We have $M_{X+Y}(t) = M_X(t)M_Y(t)$ by independence.

$$\begin{aligned}
 M_x(t) &= E[e^{tx}] = \int_0^\infty e^{tx} \theta e^{-\theta x} dx \\
 &= \theta \int_0^\infty e^{-x(\theta-t)} dx \\
 \stackrel{\alpha=1}{\sim} \stackrel{\beta=(\theta-t)^{-1}}{\sim} &= \theta \int_0^\infty x^{\alpha-1} e^{-x/\beta} dx = \theta \beta^\alpha \Gamma(\alpha) = \theta (\theta-t)^{-1} \Gamma(1)^{-1} = \theta (\theta-t)^{-1} \\
 \Rightarrow M_{x+y}(t) &= \theta^2 (\theta-t)^{-2}
 \end{aligned}$$

From (a) we define $f(x; \theta) = \frac{1}{\theta^2} x e^{-x/\theta}$ with MGF $\theta^{-2} (\theta-t)^{-2}$

$\Rightarrow f(x; \theta)$ defines $2\text{Exp}(\theta)$ since $2\text{Exp}(\theta)$ has MGF $\theta (\theta-t)^{-1}$.

$$\Rightarrow f_s(x) = \theta^2 x e^{-x/\theta}$$

10.4 19) Recall the MGF of a Bernoulli random variable X is $M_x(t) = (1-p+pe^t)^n$ for $t \in \mathbb{R}$.

Find the MGF of $Z_n = \frac{X-np}{\sqrt{np(1-p)}}$ as $n \rightarrow \infty$ (where p fixed).

We have:

$$\begin{aligned}
 M_{Z_n}(t) &= M_{\frac{X-np}{\sqrt{np(1-p)}}}(t) = e^{-\sqrt{np(1-p)}^{-1/2} t} M_x(\sqrt{np(1-p)} t) \\
 &= e^{-\sqrt{np(1-p)}^{-1/2} t} (1-p+pe^{\sqrt{np(1-p)} t})^n \\
 &= ((e^{\frac{1}{\sqrt{n}} \sqrt{np(1-p)}^{-1/2} t}) / (1-p+pe^{\sqrt{np(1-p)} t}))^n
 \end{aligned}$$

10.4 20) The number of trades N of stock XXX in a day follows $\text{Poi}(\lambda)$. At each trade (i^{th} trade), the change in stock price is X_i , where $X_i \sim N(0, \sigma^2)$. Say X_i are independent of each other and of N . Find the MGF of the total change in stock price over the day.

$$\begin{aligned}
 \text{We want } M_X(t) \text{ where } X &= \sum_{i=1}^N X_i. \text{ Thus } M_X(t) = E[e^{tX}] = E[e^{t \sum_{i=1}^N X_i}] \stackrel{\text{independence}}{=} \prod_{i=1}^N E[e^{tX_i}] \stackrel{\text{identical}}{=} (M_{X_i}(t))^n \\
 &= (e^{\frac{\sigma^2 t^2}{2}})^n = e^{\frac{n\sigma^2 t^2}{2}}
 \end{aligned}$$

so $X \sim N(0, n\sigma^2)$ by MGF uniqueness.

9.7.1) In a row of 25 switches, each is considered to be "on" or "off". Define $P(\text{on}) = 0.6$ for a switch, independently of other switches. Find the mean and variance of #unlike pairs among the 24 adjacent switch pairs.

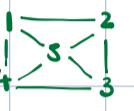
In the desired outcome, there are only two possible configurations: ON, OFF, ON, ..., ON so say \mathcal{U} = all pairs unlike.

Define $X_i = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ switch on} \\ 0 & \text{if } i^{\text{th}} \text{ switch off.} \end{cases}$ Say $Y_i = \begin{cases} 1 & \text{if } X_i \neq X_{i+1} \\ 0 & \text{if } X_i = X_{i+1}. \end{cases}$

$$\text{Then } P(\mathcal{U}) = P(X_1=1, X_2=0, \dots, X_{25}=1) + P(X_1=0, X_2=1, \dots, X_{25}=0)$$

$$= P(X_1=1)P(X_2=0)\dots P(X_{25}=1) + P(X_1=0)P(X_2=1)\dots P(X_{25}=0) \text{ by independence}$$

$$= 0.6^{13} 0.4^{12} + 0.6^{12} 0.4^{13}.$$

9.7.36)  Define $X_i = 1$ if i^{th} island isolated. We have $X_i = 1$ with probability p^3 for $i=1\dots, 4$ and p^4 for $i=5$.

Let $Y = \sum_{i=1}^5 X_i$ be the number of isolated islands.

$$\text{Then } E[Y] = E\left[\sum_{i=1}^5 X_i\right] = \sum_{i=1}^5 E[X_i] = 4E[X_1] + E[X_5] = 4p^3 + p^4.$$

$$\text{Var}(Y) = \text{Var}\left(\sum_{i=1}^5 X_i\right) = \sum_{i=1}^5 \text{Var}(X_i) + 2 \sum_{i < j} \text{Cov}(X_i, X_j). \text{ Really annoying.}$$

Exercise: Say X, Y are continuous random vars. with joint function $f_{x,y}(x, y) = f(x, y) = e^{-x-y}$ for $x, y > 0$.

a) Find $M_{x,y}(t_1, t_2) = M(t_1, t_2)$ and specify the ranges of t_1, t_2 for which it exists.

$$\begin{aligned} \text{We have } M_{x,y}(t_1, t_2) &= E[e^{tx+ty}] = \iint_{\mathbb{R}^2} e^{tx+ty} f_{x,y}(x, y) dx dy = \iint_{\mathbb{R}^2} e^{tx} e^{ty} e^{-x} e^{-y} dx dy = \int_0^\infty e^{tx} e^{-x} \int_0^\infty e^{ty} e^{-y} dy dx \\ &= \int_0^\infty e^{x(t_1-1)} dx \int_0^\infty e^{y(t_2-1)} dy \\ &\stackrel{t_1 < 1, t_2 < 1}{=} \left[\frac{e^{x(t_1-1)}}{t_1-1} \right]_0^\infty \left[\frac{e^{y(t_2-1)}}{t_2-1} \right]_0^\infty \\ &= \frac{1}{(1-t_1)(1-t_2)}. \end{aligned}$$

b) What is $M_x(t)$ and what is its marginal distribution?

$$M_x(t) = M_{x,y}(t, 0) = \frac{1}{1-t}. \text{ Since Gamma}(\alpha, \beta) \text{ has MGF } \frac{1}{(1-\beta t)^\alpha} \text{ we have } X \sim \text{Gamma}(1, 1).$$

c) Are X and Y independent?

We have $X \perp Y \Leftrightarrow M_x M_y = M_{x,y}$. We have this property clearly, so $X \perp Y$.

Exercise: Say $X \sim N(0, 1)$.

a) Find $M_x(t)$.

$$\begin{aligned} \text{We have } M_x(t) &= E[e^{tx}] = \int_{\mathbb{R}} e^{tx} f_x(x) dx \\ &= \int_{-\infty}^{\infty} e^{tx} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{tx - \frac{x^2}{2}} dx = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{\frac{t^2}{2} + \frac{x^2}{2} - tx - \frac{t^2}{2}} dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{\frac{t^2}{2}} e^{-\frac{1}{2}(x-t)^2} dx = e^{\frac{t^2}{2}} \underbrace{\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-t)^2}{2}} dx}_{\text{since pdf of } N(t, 1)} = e^{\frac{t^2}{2}} \text{ for all } t \in \mathbb{R}. \end{aligned}$$

b) Say $Z \sim N(\mu, \sigma^2)$. Find M_z using M_x .

$$Z = \tau X + \mu \text{ by reparametrization. Then } M_z(t) = M_{\tau X + \mu}(t) = e^{\mu t} M_x(\tau t) = e^{\mu t} e^{\frac{(\tau t)^2}{2}} = e^{\mu t} e^{\frac{\sigma^2 t^2}{2}}.$$

Exercise: If $X_n \xrightarrow{P} a$ and $Y_n \xrightarrow{P} b$, then $X_n Y_n \xrightarrow{P} ab$. (Hint: Show $Y_n \xrightarrow{P} b \Rightarrow \exists M > 0$ such that $\lim_{n \rightarrow \infty} P(|Y_n| < M) = 1$.)

We prove the hint. Say $Y_n \xrightarrow{P} b$, so by definition, for all $\varepsilon > 0$, we have $\lim_{n \rightarrow \infty} P(|Y_n - b| \geq \varepsilon) = 0$.

Notice $Y_n \xrightarrow{P} b \Rightarrow Y_n \xrightarrow{D} b$, so $F_{Y_n}(x) \rightarrow F_b(x)$ at all continuity points of $F_b(x)$ (so $x \in \mathbb{R} \setminus \{b\}$).

$$\Rightarrow \lim_{n \rightarrow \infty} F_{Y_n}(x) = F_b(x) \text{ for } x \in \mathbb{R} \setminus \{b\}.$$

Thus we have $F_{Y_n}(M) \rightarrow F_b(M) = \begin{cases} 0 & \text{for } M < b \\ 1 & \text{for } M \geq b \end{cases}$ so pick $M > b$. Satisfied.

$$F_{Y_n}(-M) \rightarrow F_b(-M) = \begin{cases} 0 & \text{for } -M < b \\ 1 & \text{for } -M \geq b \end{cases}$$

Therefore $\exists M, N$ such that $\lim_{n \rightarrow \infty} P(|Y_n| < N) = 1$ and $\lim_{n \rightarrow \infty} P(|X_n| < M) = 1$.

We want $\lim_{n \rightarrow \infty} P(|X_n Y_n - ab| \geq \varepsilon) = 0$

$$\text{We have } P(|X_n Y_n - ab| \geq \varepsilon) = P(|X_n Y_n + aY_n - aY_n - ab| \geq \varepsilon)$$

$$\leq P(|(X_n - a)Y_n| + |a(Y_n - b)| \geq \varepsilon) \text{ since } |\alpha + \beta| \leq |\alpha| + |\beta|$$

$$\leq P(|(X_n - a)Y_n| \geq \frac{\varepsilon}{2}) + P(|a(Y_n - b)| \geq \frac{\varepsilon}{2})$$

... algebra.

Exercise: Say $(X_i)_{i \in \mathbb{N}}$ is a sequence such that $E[X_i] = \mu_i$ and $\text{Var}(X_i) = \sigma_i^2$. If $\mu_i \rightarrow \mu$ and $\sigma_i^2 \rightarrow 0$, then $X_n \xrightarrow{P} \mu$.

By Markov's inequality, we have $P(|X_n - \mu| \geq C) \leq \frac{E[|X_n - \mu|^k]}{C^k}$ for any $C, k > 0$.

Therefore $P(|X_n - \mu| \geq \varepsilon) \leq \frac{E[|X_n - \mu|^2]}{\varepsilon^2} = \frac{\text{Var}(X_n)}{\varepsilon^2} = \frac{\sigma_n^2}{\varepsilon^2} \xrightarrow{n \rightarrow \infty} 0$. Therefore $X_n \xrightarrow{P} \mu$ by definition.

Exercise: Say $X \sim \text{Poi}(\mu)$ where $\mu \in \mathbb{Z}$ large. Show that $\frac{X-\mu}{\sqrt{\mu}} \xrightarrow{D} N(0, 1)$ and use the approximation to find $P(X > 12)$.

We use the CLT. Notice that if $X_1 \sim \text{Poi}(\lambda_1) \perp X_2 \sim \text{Poi}(\lambda_2)$ then $X_1 + X_2 \sim \text{Poi}(\lambda_1 + \lambda_2)$.

Therefore since $\mu \in \mathbb{Z}$ large, consider $X_i \xrightarrow{D} \text{Poi}(1)$ so $\sum_{i=1}^n X_i \sim \text{Poi}(\sum_{i=1}^n 1) = \text{Poi}(\mu)$.

Since $E[X_i] = 1 < \infty$ and $\text{Var}(X_i) = 1 < \infty$ we have

$$\begin{aligned} &\xrightarrow{\text{CLT}} \frac{\sqrt{\mu}(\bar{X}_n - 1)}{\sqrt{\mu}} \xrightarrow{D} N(0, 1) \\ &\quad \begin{matrix} \xrightarrow{\text{mean}} \\ \xrightarrow{\text{standard deviation}} \end{matrix} \sqrt{\mu}\left(\frac{1}{\sqrt{\mu}}\sum_{i=1}^n X_i - 1\right) = \sqrt{\mu}\left(\frac{1}{\sqrt{\mu}}X - 1\right) = \frac{1}{\sqrt{\mu}}(X - \mu) \xrightarrow{D} N(0, 1). \end{aligned}$$

We want $P(X > 12)$. Notice $P(X > 12) = 1 - P(X \leq 12) = 1 - F_X(12)$ and we have $F_X(x) \rightarrow F_z(x)$ at all $x \in \mathbb{R} \setminus \mathbb{Z}$.

Thus by using a continuity correction, we have $1 - F_X(12) = 1 - F_X(12.5)$

$$\begin{aligned} &\approx 1 - F_z(12.5) \quad \text{where } \xrightarrow{\text{CLT}} X \xrightarrow{D} Z \sim N(\mu, \sigma^2) \\ &= 1 - F_y\left(\frac{12.5 - \mu}{\sqrt{\mu}}\right) \quad \text{where } Y \sim N(0, 1). \end{aligned}$$