

Lecture 1
(31 July 2023)

Probability & Random processes

Grading Plan (Tentative)

Assignments - 15 %.

Quiz 1 - 15 %.

Quiz 2 - 15 %.

Mid-sem - 20 %.

End-sem - 35 %.

Module 1 (Sets and Basics of Probability).

Sets, Probabilistic models, Conditional probability, Bayes' Rule

Module 2 (Discrete Random Variables).

Module 3 (Continuous Random variables).

Module 4 (Tail Bounds and Limit Theorems).

Module 5 (Random processes).

Textbook : Introduction to Probability

by

Bertsekas and Tsitsiklis, 2nd Edition

P.T.O

Introduction to Probability

Why study Probability:

- Randomness and uncertainty exist in our daily lives.
- Probability theory is a mathematical framework that allows us to describe and analyze random phenomena (i.e., events or experiments whose outcomes we cannot predict with certainty).
- Probability (roughly) means possibility. It helps us to predict how likely or unlikely an event will occur.
[Example: Flipping a fair coin]

Different Approaches to Probability

A. Probability as the Ratio of Favourable to Total Outcomes (classical Approach)

Probability of an event E
= No. of ways E can occur

Total no. of possible outcomes

Example. Suppose we throw a pair of unbiased dice.

- 1) What is the probability of getting a sum of 7?
- 2) What is the probability of getting a sum of 10?

Ans. 1) $6/36 = 1/6$, 2) $3/36 = 1/12$.

Example. Suppose we throw a fair coin thrice. What is the

Probability of getting at most 2 Heads?

Ans. The set of possible outcomes is {TTT, TTH, HTT, THT, THH, HTH, HHT, HHH}.

$$P(\text{at most 2 Heads}) = \frac{7}{8}$$
$$(= 1 - \frac{1}{8})$$

This approach suffers from at least two significant problems.

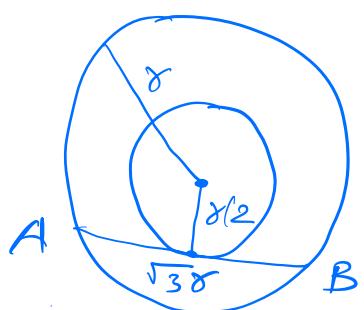
- 1) It cannot deal with outcomes that are not equally likely; and
- 2) it cannot handle uncountably infinite outcomes without ambiguity.

(when the no. of possible outcomes is infinite, we must use length, area, or some other measure of infinity for determining the ratio)

Example (Bertrand Paradox). we are given a circle C of radius δ and we wish to determine the probability P that the length l of a 'randomly selected' chord AB is greater than the length $\sqrt{3}\delta$ (length of the side of an equilateral triangle).

We show that this problem has at least two reasonable and different solutions,

I. If the center M of the chord AB lies inside the circle C_1 of radius $\delta/2$:



- then $l > \sqrt{3}\delta$.

Favourable outcomes for chord center

= all Points inside the circle C_1 ,

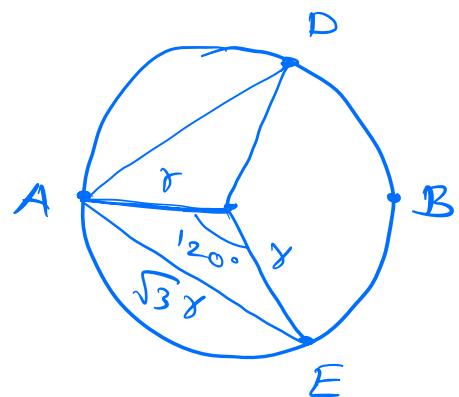
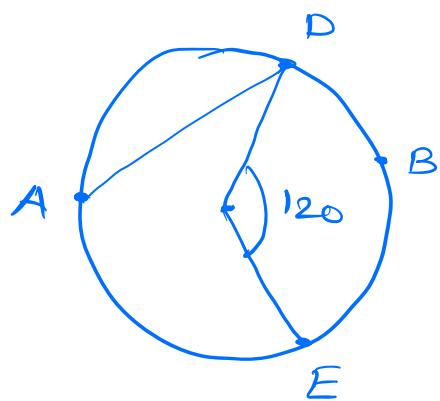
Total possible outcomes for chord center

= all Points inside C

$$P = \frac{\text{area}(C_1)}{\text{area}(C)} = \frac{1}{4}.$$

(using area as the measure of infinite points)

B. we assume that the end A of chord AB is fixed. This reduces the no. of possibilities but it has no effect on the value of P because the no. of favourable locations of B is reduced proportionately,



If B is on the 120° arc DBE, then
 $l > \sqrt{3}\sigma$.

Favourable outcomes

= Points on this arc,

Total outcomes

= all Points on the circumference of C

$$P = \frac{2\pi\sigma/3}{2\pi\sigma} = \frac{1}{3}.$$

This demonstrates the ambiguities associated with the classical definition,

B. Probability as a Measure of Frequency of Occurrence

- Define the probability of an event E by performing the experiment n times. No. of times E occurs is denoted by n_E .

$$P(E) = \lim_{n \rightarrow \infty} \frac{n_E}{n}.$$

clearly since $n_E \leq n$, we must have
 $0 \leq P(E) \leq 1$.

Disadvantages

- 1) We can never perform the experiment infinite number of times
- 2) The definition does not capture belief factor

Despite the problems with the frequency definition of probability, the relative freq. concept is essential in applying probability theory to the real world.

C. Probability Based on Axiomatic Theory

We need to introduce

- Random experiment - It is simply an experiment in which outcomes are non-deterministic, that is probabilistic.
- Sample space - set of all outcomes of the experiment.
- Event - Any subset of the sample space.

Sample Space and Probability

Sets

- A set is a collection of objects, which are the elements of the set.
- A set with no elements is called the empty set, denoted by ϕ or $\{\}$.
- Finite set: A set with a finite no. of elements:

$$S = \{x_1, x_2, \dots, x_n\}.$$

Eg. set of possible outcomes of a coin toss
 $= \{H, T\} = \{\text{heads, tails}\}.$

set of possible outcomes of a die roll

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

- If a set contains infinitely many elements x_1, x_2, \dots , which can be enumerated in a list (i.e., a bijective mapping with naturals), we write $S = \{x_1, x_2, \dots\}$ and call S as a countably infinite set.

Eg. set of even integers = $\{0, \pm 2, 4, -4, \dots\}$

- A set is uncountable if its elements cannot be enumerated in a list.

- Subset notation: $A \subseteq B \Leftrightarrow (x \in A \Rightarrow x \in B)$.
- Universal set Ω contains all objects that could be of interest in a particular context.

Set Operations.

- Complement of a set S , $S^c = \{x \in \Omega : x \notin S\}$.

$$\Omega^c = \emptyset, \emptyset^c = \Omega.$$

- Union of two sets A and B ,

$$A \cup B = \{x \in \Omega : x \in A \text{ or } x \in B\}.$$

- Intersection of two sets A and B ,

$$A \cap B = \{x \in \Omega : x \in A \text{ and } x \in B\}.$$

- Infinite union.

If for every $n \in \mathbb{N}$ we are given S_n ,

$$\bigcup_{n=1}^{\infty} S_n = S_1 \cup S_2 \cup \dots = \{x \in \Omega : x \in S_n \text{ for some } n\}$$

- Infinite intersection

$$\bigcap_{n=1}^{\infty} S_n = S_1 \cap S_2 \cap \dots = \{x \in \Omega : x \in S_n \text{ for all } n\}$$

- Two sets are said to be disjoint if their intersection is empty (i.e., there is no element in common).
More generally, given multiple no. of sets, they are said to be pairwise disjoint if every pair of those sets is disjoint

- Real numbers \mathbb{R}

set of pairs of reals \mathbb{R}^2

set of triplets of reals \mathbb{R}^3 .

Some Properties

$$- A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$$

$$- A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$$

- De Morgan's laws

$$(A \cup B)^c = A^c \cap B^c, \quad (A \cap B)^c = A^c \cup B^c$$

Exercise 1.1. Prove that $\mathbb{Q} \cap [0, 1]$ is a countably infinite set.

Exercise 1.2. Prove that $\{0, 1\}^\infty$ is an uncountably infinite set.

[Cantor's diagonalization argument]

Exercise 1.3.

(i) For $n \in \mathbb{N}$, show that

$$\left(\bigcup_{i=1}^n A_i \right)^c = \bigcap_{i=1}^n A_i^c.$$

[Hint: Use mathematical induction]

(ii) Show that

$$\left(\bigcup_{i=1}^{\infty} S_i \right)^c = \bigcap_{i=1}^{\infty} S_i^c,$$

$$\left(\bigcap_{i=1}^{\infty} S_i \right)^c = \bigcup_{i=1}^{\infty} S_i^c.$$

[Note that induction does not directly give such a statement for infinite number.]

Reason (via a simple counter example):

Let $P(n)$ be the statement " n is finite".

$P(n)$ is true for every $n \in \mathbb{N}$,

$P(\infty)$ is false.]

Lecture 2 (3 August 2023)

Probabilistic Models

- A probabilistic model is a mathematical description of an uncertain situation or a random experiment.

Elements of a Probabilistic model:

- Sample space Ω , the set of all possible outcomes of an experiment.
A subset of a sample space is called an event.
- Probability law, which assigns a non-negative number $P(A)$ to an event A that encodes our knowledge or belief about the collective "likelihood" of the elements of A .

Sample space and Events

- Sample space is the set of all possible outcomes of a random experiment.
(The random experiment produces exactly one out of all the possible outcomes)

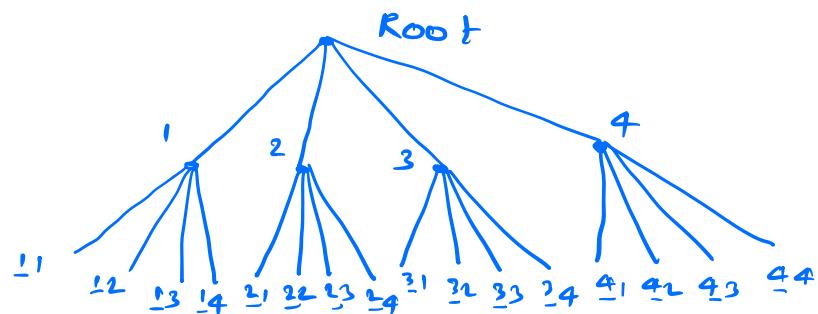
- The elements of the sample space should be
 - mutually exclusive
 - collectively exhaustive

Consider rolling a die. Is $\{1, 2, 3, 4, 5\}$ the sample space?

[No, it is not collectively exhaustive]

Example. Consider rolling a 4-sided die twice (a single random experiment).

11	12	13	14
21	22	23	24
31	32	33	34
41	42	43	44

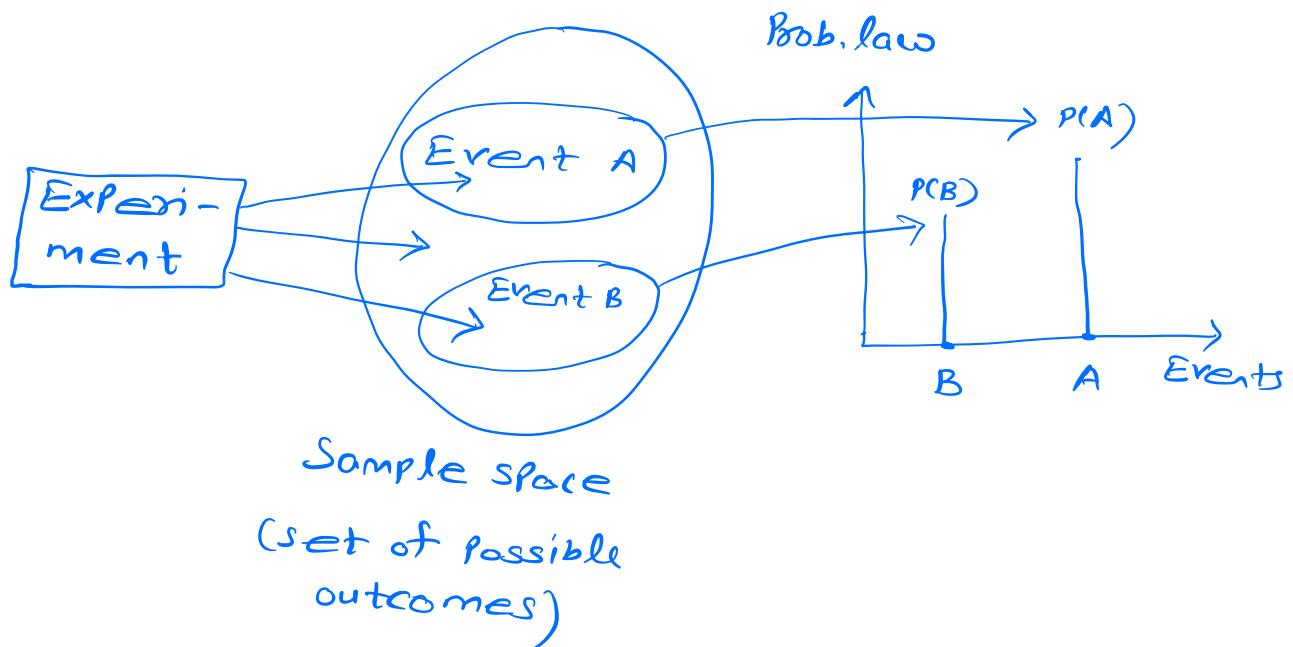
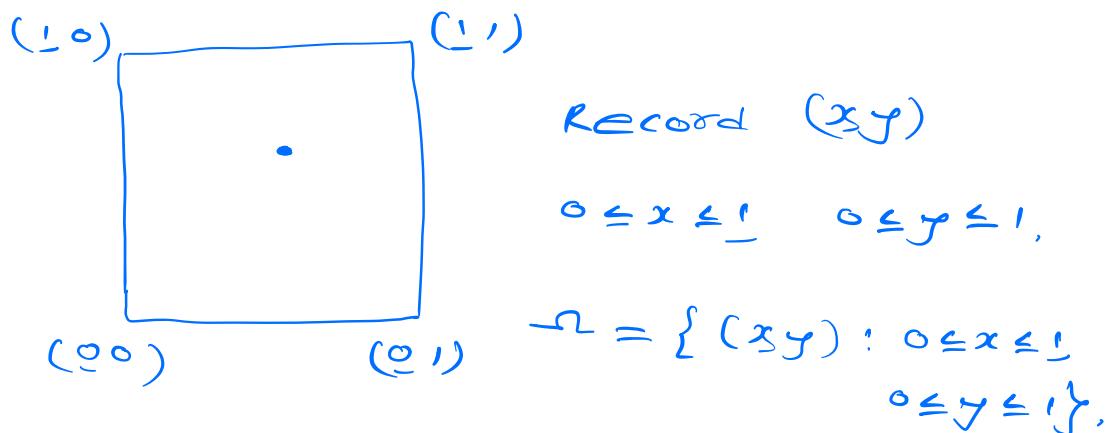


These are two equivalent descriptions of the sample space of a random experiment involving two rolls of a 4-sided die.

- Sample space of an experiment may consist of a finite or an infinite number of possible outcomes. Finite sample spaces are conceptually

and mathematically simpler, still sample spaces with an infinite number of elements are quite common,

Example. Continuous sample space. Consider throwing a dart on a 1×1 square target.



Pictorial view of a Probabilistic model

- The sample space should have enough detail to distinguish between all outcomes of interest, while avoiding the irrelevant details.

Example. Consider a single random experiment that involves 3 successive coin tosses, and there are two different scenarios of interest.

Game 1: we receive Rs. 1 each time a head comes up.

Game 2: we receive Rs. 1 for every coin toss up to and including the first time a head comes up. Then, we receive Rs. 2 for every coin toss up to the second time a head comes up, Rs. 4 up to the third time.

$$\Omega_1 = \{0, 1, 2, 3\}$$

$$\Omega_2 = \left\{ \begin{array}{c} \text{HHH, HHT, THH, HTH, TTH, THT, HTT, TTT} \\ 1 \ 2 \ 4 \ 124 \ 112 \ 122 \ 111 \end{array} \right\}$$

Probability Laws

- Suppose we have settled on the sample space Ω associated with a random experiment. To complete the probabilistic model, we must introduce a probability law.
- Intuitively, this specifies the "likelihood" of any outcome, or any set of possible outcomes (an event). More precisely, we assign a number $P(A)$ to every event A satisfying the following axioms.

Axioms of Probability:

1. (Nonnegativity) $P(A) \geq 0$ for every event A ,
2. (Additivity) If A and B are two disjoint events, then $P(A \cup B) = P(A) + P(B)$.

More generally, if A_1, A_2, \dots is a sequence of disjoint events, then

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

3. (Normalization) $P(\Omega) = 1$.

Properties of Probability Laws

(1) $P(A) + P(A^c) = 1.$

(2) $P(A) \leq 1.$

(3) $P(\emptyset) = 0.$

(4) For pairwise disjoint events $A, B,$ and C

$$P(A \cup B \cup C) = P(A) + P(B) + P(C).$$

(Similarly for $n \in \mathbb{N}$ events)

(5) If $A \subseteq B$, then $P(A) \leq P(B).$

(6) $P(A \cup B) = P(A) + P(B) - P(A \cap B).$

(7) $P(A \cup B) \leq P(A) + P(B)$ [Union Bound].

(8) $P(A \cup B \cup C) = P(A) + P(B \cap A^c) + P(A^c \cap B^c \cap C).$

$$A \cap A^c = \emptyset \Rightarrow P(A \cup A^c) = P(\Omega) = 1 = P(A) + P(A^c)$$

[This proves (1)]

$$P(A) = 1 - P(A^c) \leq 1 \text{ since } P(A^c) \geq 0 \text{ by nonnegativity}$$

[This proves (2)]

$$\begin{aligned} \Omega \cap \emptyset &= \emptyset \Rightarrow P(\Omega) = P(\Omega) + P(\emptyset) \\ &\Rightarrow P(\emptyset) = 0. \quad [\text{This proves (3)}] \end{aligned}$$

For disjoint sets $A, B,$ and C

$$P(A \cup B \cup C) = P(A) + P(B \cup C) = P(A) + P(B) + P(C)$$

Similarly $P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i)$ for pairwise disjoint sets A_1, A_2, \dots, A_n ,
 [This proves (4)]

$$\begin{aligned} A \subseteq B &\quad B = A \cup (B \setminus A) \\ &= A \cup (B \cap A^c) \end{aligned}$$

since A and $B \cap A^c$ are disjoint sets

$$\begin{aligned} P(B) &= P(A \cup (B \cap A^c)) = P(A) + P(B \cap A^c) \\ &\geq 0. \end{aligned}$$

[This proves (5)]

$$\begin{aligned} P(A \cup B) &= P(A \cup (B \cap A^c)) = P(A) + P(B \cap A^c) \\ &= P(A) + P(B) - P(A \cap B) \end{aligned}$$

$$\begin{aligned} \text{since } P(B) &= P((A \cap B) \cup (B \cap A^c)) \\ &= P(A \cap B) + P(A^c \cap B) \end{aligned}$$

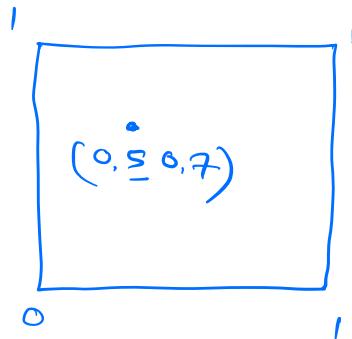
[This proves (6)]

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\ &\leq P(A) + P(B) \quad (\because \text{Nonnegativity}) \end{aligned}$$

$$P(A \cup B) \geq \max\{P(A), P(B)\} \quad [\text{This proves (7)}]$$

Example. (Continuous model)

Probabilistic models with continuous sample spaces differ from their discrete counterparts in that the probabilities of the single-element events may not be sufficient to characterize the probability law. Consider throwing a dart on 1×1 square target.



$$A \subseteq [0,1]^2 \quad P(A) = \text{area of } A.$$

$$P(\{(0.5, 0.7)\}) = 0$$

$$P([0, 0.5]^2) = \frac{1}{4}$$

Lecture 3

(7 August 2023)

Conditional Probability

- Provides us with a way to reason about the outcome of an experiment based on Partial information.
- More precisely given an experiment - a corresponding sample space, and a probability law, suppose that we know that the outcome is within some given event B , we wish to quantify the likelihood that the outcome also belongs to some other given event A .
- Conditional probability of A given B denoted by $P(A|B)$.

Example,

On rolling a die, what is the probability that the outcome is 2 given that the outcome is even?

Given an experiment, a corresponding sample space, and a probability law, suppose that we know that the outcome is in B . We want to quantify the likelihood that the outcome also belongs to some other event A .

$$P(\text{outcome is 1} \mid \text{the outcome is even}) \\ = 0$$

2	1
4	3
6	5

$$P(A \mid B) = 0 \quad \text{if } A \cap B = \emptyset$$

$$\propto P(A \cap B)$$

$$P(\text{outcome is 2} \mid \text{the outcome is even})$$

$$= \frac{1}{3} = \left(\frac{1}{6}\right) / \left(\frac{1}{2}\right)$$

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

Conditional probability of A given B s.t. $P(B) > 0$

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

- This specifies a new probability law on the same sample space Ω .

Exercise 3.1. For a fixed $B \subseteq \Omega$, prove that $P(\cdot|B)$ satisfies the probability axioms.

- If the possible outcomes are finitely many and equally likely, then

$$P(A|B) = \frac{\text{no. of elements in } A \cap B}{\text{no. of elements in } B}$$

Exercise 3.3 A fair 4-sided die is rolled twice and we assume that all sixteen possible outcomes are equally likely. Let x and y be the results of the 1st and the 2nd roll, respectively. Find $P(\max\{x,y\}=m | \min\{x,y\}=z)$ for $m=1, 2, 3, 4$.

Bayes' rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \text{ for } P(B) > 0$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}, \text{ for } P(A) > 0$$

$$\Rightarrow P(A|B)P(B) = P(B|A)P(A) \text{ for } P(A)P(B) > 0.$$

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

Multiplication Rule

$$P\left(\bigcap_{i=1}^n A_i\right) = P(A_1)P(A_2|A_1) \dots P(A_n | \bigcap_{i=1}^{n-1} A_i).$$

$$\begin{aligned} P(A_1 \cap A_2 \cap A_3) &= P(A_1)P(A_2 | A_1)P(A_3 | A_1 \cap A_2) \\ &= P(A_1)P(A_2 | A_1)P(A_3 | A_1 \cap A_2) \end{aligned}$$

Proof of the multiplication rule follows by mathematical induction.

Total Probability Theorem

Let A_1, A_2, \dots, A_n be mutually exclusive and collectively exhaustive events of the sample space (each possible outcome is included in exactly one of the events A_1, A_2, \dots, A_n) and assume that $P(A_i) > 0 \forall i$.

Then for any event B we have

$$P(B) = \sum_{i=1}^n P(A_i \cap B)$$

$$= \sum_{i=1}^n P(A_i) P(B|A_i).$$

Proof. $P(B) = P(B \cap \Sigma)$

$$= P(B \cap (A_1 \cup A_2 \cup \dots \cup A_n))$$

$$= P\left(\bigcup_{i=1}^n (B \cap A_i)\right)$$

$$= \sum_{i=1}^n P(A_i \cap B) = \sum_{i=1}^n P(A_i) P(B|A_i),$$

Bayes' Rule (Refined version)

Let A_1, A_2, \dots, A_n be disjoint events that form a partition of the sample space and assume that $P(A_i) > 0 \forall i$. Then, for any event B such that $P(B) > 0$ we have

$$P(A_i|B) = \frac{P(A_i) P(B|A_i)}{P(B)}$$

$$= \frac{P(A_i) P(B|A_i)}{\sum_{j=1}^n P(A_j) P(B|A_j)},$$

Example. (Exercise 3,4 in class)

Let $A = \{\text{an aircraft is present}\}$,

$B = \{\text{the radar generates an alarm}\}$.

We are given that

$$P(A) = 0.05 \quad P(B|A) = 0.99 \quad P(B|A^c) = 0.1.$$

$P(\text{aircraft is present} | \text{radar generates alarm})$

$$= P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(A)P(B|A) + P(A^c)P(B|A^c)}$$

$$= \frac{0.05 \times 0.99}{0.05 \times 0.99 + 0.95 \times 0.1}.$$

Independence. Two events A & B are called independent events if $P(A \cap B) = P(A)P(B)$.

$$P(A|B) = P(A) \quad \text{for } P(B) > 0, \quad \text{or}$$

$$P(B|A) = P(B) \quad \text{for } P(A) > 0,$$

Exercise 3.2. If A and B are independent events, show that

- (i) A^c and B are independent
- (ii) A^c and B^c are independent.

Lecture 4
(17 August 2023)

Independence of several events

We say that the events A_1, A_2, \dots, A_n are independent if

$$P(\bigcap_{i \in S} A_i) = \prod_{i \in S} P(A_i)$$

for every subset S of $\{1, 2, \dots, n\}$.

For three events A, B and C :

$$P(A \cap B) = P(A) P(B) \rightarrow ①$$

$$P(B \cap C) = P(B) P(C) \rightarrow ②$$

$$P(A \cap C) = P(A) P(C) \rightarrow ③$$

$$P(A \cap B \cap C) = P(A) P(B) P(C) \rightarrow ④$$

①-③ do not imply ④

④ does not imply ①-③

Exercise 4.1. Consider two independent fair coin tosses, and the following events:

$$H_1 = \{HT, HH\},$$

$$H_2 = \{TH, HH\},$$

$$D = \{TT, HH\}.$$

Show that H_1, H_2 and D are pairwise independent but not independent.

Exercise 4.2. Consider two independent rolls of a fair six-sided die and the following events:

$$A = \{1^{\text{st}} \text{ roll is } 1, 2 \text{ or } 3\}$$

$$B = \{1^{\text{st}} \text{ roll is } 3, 4 \text{ or } 5\}$$

$$C = \{\text{the sum of the two rolls is } 9\}.$$

Show that ④ holds but ①-③ do not hold,

Conditional Independence

Given an event c with $P(c) > 0$ the events A and B are called independent if

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

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

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

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

Assume $P(B|c) \neq 0$ the above definition can be written as

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

In words this states that if c is known to have occurred, an additional knowledge that B also occurred does not change the probability of A .

Interestingly, independence of two events A and B does not imply conditional independence and vice versa.

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

$$\Leftrightarrow P(A \cap B|c) = P(A|c) P(B|c)$$

in general.

Exercise 4.3. Consider two independent fair coin tosses, and the events same as in Exercise 4.1.

$$H_1 = \{ HT, HH \}$$

$$H_2 = \{ TH, HH \}$$

$$D = \{ HT, TH \}.$$

$$P(H_1 \cap H_2) = P(H_1)P(H_2) \text{ and } P(H_1 \cap H_2 | D) \neq P(H_1 | D)P(H_2 | D),$$

Exercise 4.4. There are two coins: a fair coin and a fake two-headed coin ($P(H) = 1$). We choose one of the two coins at random, each being chosen with probability $\frac{1}{2}$ and toss it twice.

H_1, H_2 same as Exercise 4.3.

$$F = \{ \text{fair coin is selected} \}.$$

$$P(H_1 \cap H_2 | F) = P(H_1 | F)P(H_2 | F) \text{ and}$$

$$P(H_1 \cap H_2) \neq P(H_1)P(H_2).$$

Review of Counting.

Permutations: Given n distinct objects, and let $k \leq n$, we wish to count the number of different ways that we can pick k out of these n objects and arrange them in a sequence, i.e., the number of distinct k -object sequences.

$${}^n P_k = \frac{n!}{(n-k)!} = n(n-1)\dots(n-k+1).$$

Combinations: counting the number of k -element subsets of a given n -element set. Notice that forming a combination is different than forming a permutation, because in a combination there is no ordering of the selected elements.

$${}^n C_k = {}^n P_k = \frac{n!}{(n-k)! k!}$$

Partitions: consider n and n_1, \dots, n_k s.t.
 $n = n_1 + n_2 + \dots + n_k$.

No. of partitions of n distinct elements
 into k disjoint subsets with the i th subset
 containing exactly n_i elements

$$= \frac{n!}{n_1! n_2! \dots n_k!}$$

Continuity of Probability

Notation. $\bigcup_{i=1}^{\infty} A_i \triangleq \lim_{n \rightarrow \infty} \bigcup_{i=1}^n A_i$,

$$\sum_{i=1}^{\infty} P(A_i) \triangleq \lim_{n \rightarrow \infty} \sum_{i=1}^n P(A_i)$$

Theorem.

$$\lim_{n \rightarrow \infty} P\left(\bigcup_{i=1}^n A_i\right) = P\left(\lim_{n \rightarrow \infty} \bigcup_{i=1}^n A_i\right).$$

Proof. $B_i = A_i$,

$$B_i = A_i \setminus \bigcup_{j=1}^{i-1} A_j, \quad i = 2, 3, \dots$$

(i) B_i 's are disjoint (justify)

(ii) $\bigcup_{i=1}^n A_i = \bigcup_{i=1}^n B_i$, $n \in \mathbb{N}$ and

$$\bigcup_{i=1}^{\infty} A_i = \bigcup_{i=1}^{\infty} B_i.$$

Let $c_n = \bigcup_{i=1}^n A_i$. $c_2 = B_1 \cup B_2$
(base case)

Assume $c_{n-1} = \bigcup_{i=1}^{n-1} B_i$.

$$\begin{aligned} c_n &= c_{n-1} \cup A_n = c_{n-1} \cup (A_n \setminus c_{n-1}) \\ &= c_{n-1} \cup (A_n \setminus \bigcup_{i=1}^{n-1} A_i) \\ &= c_{n-1} \cup B_n = \bigcup_{i=1}^n B_i. \end{aligned}$$

$$x \in \bigcup_{i=1}^{\infty} A_i \Rightarrow \exists n \text{ s.t. } x \in A_n \Rightarrow x \in \bigcup_{i=1}^n B_i$$

$$\Rightarrow x \in \bigcup_{i=1}^{\infty} B_i$$

likewise $x \in \bigcup_{i=1}^{\infty} B_i \Rightarrow x \in \bigcup_{i=1}^{\infty} A_i$.

Now we have

$$\begin{aligned}\lim_{n \rightarrow \infty} P\left(\bigcup_{i=1}^n A_i\right) &= \lim_{n \rightarrow \infty} P\left(\bigcup_{i=1}^n B_n\right) \\&= \lim_{n \rightarrow \infty} \sum_{i=1}^n P(B_i) \quad (\text{by additivity}) \\&= \sum_{i=1}^{\infty} P(B_i) \\&= P\left(\bigcup_{i=1}^{\infty} B_i\right) \quad (\text{by additivity}) \\&= P\left(\bigcup_{i=1}^{\infty} A_i\right) \\&= P\left(\lim_{n \rightarrow \infty} \bigcup_{i=1}^n A_i\right).\end{aligned}$$

□

Consequences.

1) Let A_1, A_2, \dots be a monotonically increasing sequence, i.e., $A_i \subseteq A_{i+1}$, $i \in \mathbb{N}$.

Then $P\left(\bigcup_{i=1}^{\infty} A_i\right) = \lim_{n \rightarrow \infty} P(A_n)$.

2) Let A_1, A_2, \dots be a monotonically decreasing sequence, i.e., $A_{i+1} \subseteq A_i$, $i \in \mathbb{N}$.

Then

$$P\left(\bigcap_{i=1}^{\infty} A_i\right) = \lim_{n \rightarrow \infty} P(A_n).$$

This follows from the previous consequence by considering complements of the events and using De Morgan's laws.

Lecture 5

(19 August 2023)

Continuity of probability

$$\lim_{n \rightarrow \infty} P\left(\bigcup_{i=1}^n A_i\right) = P\left(\lim_{n \rightarrow \infty} \bigcup_{i=1}^n A_i\right).$$

$$B_i = A_i \setminus \bigcup_{j=1}^{i-1} A_j, \quad i \in \mathbb{N}.$$

B_i and $B_{i'}$ are disjoint for $i \neq i'$.

Let $i < i'$ without loss of generality,

Let $x \in B_i$,

$$\Rightarrow x \in A_i$$

$$\Rightarrow x \notin A_{i'} \setminus \bigcup_{j=1}^{i'-1} A_j \quad \text{since } A_i \subseteq \bigcup_{j=1}^{i'-1} A_j$$

$$\Rightarrow x \notin B_{i'}.$$

$$\text{Let } x \in B_{i'} \Rightarrow x \in A_{i'} \setminus \bigcup_{j=1}^{i'-1} A_j$$

$$\Rightarrow x \notin A_i \Rightarrow x \notin B_i$$

Therefore B_i and $B_{i'}$ are disjoint sets.

$$\bigcup_{i=1}^n A_i = \bigcup_{i=1}^n B_i .$$

We prove this using mathematical induction,

$$\text{Let } C_n = \bigcup_{i=1}^n A_i, \quad n = 1, 2, \dots$$

$$C_2 = A_1 \cup A_2 = B_1 \cup B_2 \quad (\text{base case})$$

$$\text{Assume that } C_{n-1} = \bigcup_{i=1}^{n-1} B_i .$$

Consider

$$\begin{aligned} C_n &= C_{n-1} \cup A_n \\ &= C_{n-1} \cup (A_n \setminus C_{n-1}) \\ &= C_{n-1} \cup \left(A_n \setminus \bigcup_{i=1}^{n-1} B_i \right) \\ &= C_{n-1} \cup B_n = \bigcup_{i=1}^n B_i . \end{aligned}$$

$$\text{Also, } \bigcup_{i=1}^{\infty} A_i = \bigcup_{i=1}^{\infty} B_i .$$

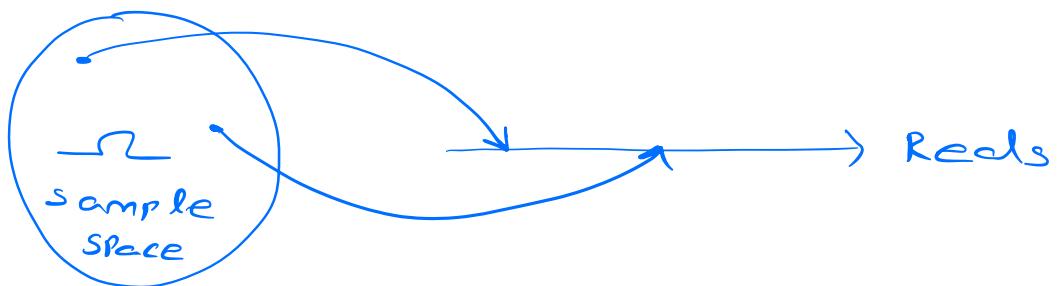
$$\text{Let } x \in \bigcup_{i=1}^{\infty} A_i \Rightarrow \exists n \in \mathbb{N} \text{ s.t. } x \in A_n$$

$$\Rightarrow x \in \bigcup_{i=1}^n A_i \Rightarrow x \in \bigcup_{i=1}^n B_i \Rightarrow x \in \bigcup_{i=1}^{\infty} B_i .$$

$$\text{Similarly } x \in \bigcup_{i=1}^{\infty} B_i \Rightarrow x \in \bigcup_{i=1}^{\infty} A_i .$$

Random Variables

In many probabilistic models, the outcomes are not numerical (real values), but they may be associated with some numerical values of interest.



Formally, a random variable is a real valued function of the experimental outcome.

$$X : \Omega \rightarrow \mathbb{R}.$$

Coin toss $\Omega = \{H, T\}$

$$X(H) = 1, X(T) = 0.$$

Roll a die twice

$X = \text{maximum of two rolls}$

Definition. A discrete random variable is a real-valued function of the outcome of the experiment that can take a finite or countably infinite number of values.

- A discrete random variable has an associated probability mass function (PMF), which gives the probability of each numerical value that the random variable can take.

Definition. Consider a probabilistic model with sample space Ω and probability law P . Let X be a random variable, $X: \Omega \rightarrow \mathbb{R}$. Probability mass of $x \in \mathbb{R}$ is defined as

$$P_X(x) = P(\{\omega \in \Omega : X(\omega) = x\}).$$

Notation. $\{X = x\} \triangleq \{\omega \in \Omega : X(\omega) = x\}$.

Example. Consider two independent tosses of a fair coin, and let X be no. of heads obtained.

$$P_X(x) = \begin{cases} \frac{1}{4} & \text{if } x=0 \text{ or } x=2 \\ \frac{1}{2} & \text{if } x=1 \\ 0 & \text{o.w.} \end{cases}$$

Let $X \subseteq \mathbb{R}$ be the range of the function
 $X: \Omega \rightarrow \mathbb{R}$.

Claim. $\sum_{x \in X} P_X(x) = 1$.

$$\text{Proof. } \sum_{x \in X} P_X(x) = \sum_{x \in X} P(X=x)$$

$$= \sum_{x \in X} P\left(\underbrace{\{\omega : X(\omega)=x\}}_{A_x}\right)$$

[$A_x, x \in X$ are disjoint events and form a partition of Ω]

$$= P\left(\bigcup_{x \in X} A_x\right) = P(\Omega) = 1.$$

(by additivity) (normalization)

- we denote the probability that X takes a value within a set $S \subseteq \mathbb{R}$ by

$$P(X \in S) \triangleq \sum_{x \in S} P_X(x).$$

Example. If X is the no. of heads obtained in two independent tosses of a fair coin,

the probability of at least one head is

$$P(X > 0) = P(X \in \{1, 2\})$$

$$= \frac{1}{2} + \frac{1}{4} = \frac{3}{4}.$$

The Bernoulli Random Variable

Consider the toss of a coin which comes up a head with probability p and a tail with a probability $1-p$.

$$X(H) = 1 \quad X(T) = 0.$$

$$P_X(1) = p, \quad P_X(0) = 1-p.$$

This is a very important r.v. (random variable). In practice it is used to model generic probabilistic situations with just two outcomes.

By combining multiple Bernoulli r.v.'s we get Binomial random variable.

Binomial Random Variable

A coin is tossed n times (independently)

$$P(H) = p \quad P(T) = 1-p.$$

Let X be the no. of heads in the n -toss sequence \rightarrow Binomial RV,

$$P_x(k) = P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k=0, 1, \dots, n.$$

$$\sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} = 1,$$

Lecture 6
(21 August 2023)

Geometric Random Variable

Suppose that we repeatedly and independently toss a coin until a head comes up for the first time. Let $P(H) = p$, $0 < p < 1$.

Geometric r.v. is the no. of tosses needed for a head to come up for the first time.

$$P_X(k) = (1-p)^{k-1} p \quad k = 1, 2, \dots,$$

Since $(1-p)^{k-1} p$ is the probability of the sequence consisting of $k-1$ successive tails followed by a head.

This is a valid PMF because

$$\sum_{k=1}^{\infty} P_X(k) = \sum_{k=1}^{\infty} (1-p)^{k-1} p$$

$$= p, \sum_{k=0}^{\infty} (1-p)^k$$

$$= p, \frac{1}{1-(1-p)} = 1,$$

Poisson Random variable

$$P_X(k) = e^{-\lambda} \cdot \frac{\lambda^k}{k!} \quad k=0, 1, 2, \dots$$

where λ is a positive parameter characterizing the PMF. This is a valid PMF because

$$\sum_{k=0}^{\infty} e^{-\lambda} \cdot \frac{\lambda^k}{k!} = e^{-\lambda} \cdot \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{-\lambda} \cdot e^{\lambda} = 1.$$

Theorem. Consider a binomial distribution with parameters n and p . As $n \rightarrow \infty$ and $p = \gamma_n$, while keeping λ constant, we have

$$\lim_{n \rightarrow \infty} \frac{n!}{(n-k)! k!} p^k (1-p)^{n-k} = e^{-\lambda} \cdot \frac{\lambda^k}{k!}.$$

Proof

$$\begin{aligned} & \frac{n!}{(n-k)! k!} p^k (1-p)^{n-k} = \frac{n(n-1)\dots(n-k+1)}{k!} \left(\frac{\lambda}{n}\right)^k \left(1-\frac{\lambda}{n}\right)^{n-k} \\ &= \underbrace{\frac{n(n-1)\dots(n-k+1)}{n^k}}_{\rightarrow 1} \cdot \underbrace{\frac{\lambda^k}{k!} \left(1-\frac{\lambda}{n}\right)^{-k} \left(1-\frac{\lambda}{n}\right)^n}_{\rightarrow 1} \end{aligned}$$

$$\lim_{n \rightarrow \infty} (1 - \frac{\lambda}{n})^n \rightarrow e^{-\lambda}.$$

This gives $\frac{e^{-\lambda} \cdot \lambda^k}{k!}$.

When n is very large and p is very small
Poisson PMF is a good approximation to
binomial PMF.

Example. $n=100$ $p=0.01$. The probability of
5 successes in 100 trials

$$= \binom{100}{5} (0.01)^5 (0.99)^{95} = 0.00290,$$

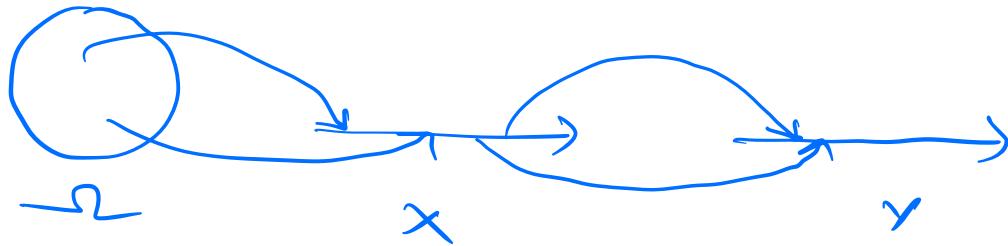
Using the Poisson PMF as an approximation with
 $\lambda=np=1$, the probability equals

$$e^{-1} \frac{1}{5!} = 0.00306,$$

Functions of Random Variables

Let X be a random variable, i.e.,
 $X: \Omega \rightarrow \mathbb{R}$, and y be a

function of x , i.e., $y = g(x)$, $g: \mathbb{R} \rightarrow \mathbb{R}$.



y is also a random variable.

$$P_y(y) = \sum_{x: g(x)=y} P_x(x),$$

Exercise 6.1
Prove this

Example. Let $y = |x|$, and

$$P_x(x) = \begin{cases} \frac{1}{9}, & \text{if } x \text{ is an integer in } [-4, 4] \\ 0, & \text{o.w.} \end{cases}$$

Find the pmf of y .

Proof. $X = [-4:4]$,

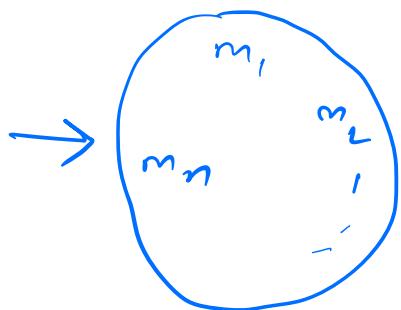
$$Y = [0:4],$$

$$P_y(0) = \frac{1}{9}, \quad P_y(i) = \frac{2}{9} \quad i \in \{1:4\}$$

$$P_y(y) = 0 \quad \text{o.w.}$$

Expectation, The PMF of x provides us with several numbers, the probabilities of all the possible values of x . It is often desirable, however, to summarize this information in a single representative number. This is accomplished by the expectation of x , which is a weighted average of the possible values of x .

Motivation. wheel of fortune



m_i comes up with probability p_i ,
 $i \in [1:n]$

Suppose you spin the wheel k times,
let k_i be the no. of times that the outcome is m_i - $k = \sum_{i=1}^n k_i$.

Total reward = $\sum_{i=1}^n m_i k_i$.

Total reward per spin = $\frac{\sum_{i=1}^n m_i k_i}{k}$.

If we interpret probabilities as relative

frequencies $\frac{k_i}{K} \rightarrow p_i$ as $K \rightarrow \infty$,

$$\sum_{i=1}^n \frac{m_i k_i}{K} \rightarrow \sum_{i=1}^n m_i p_i,$$

Expectation (or Expected value or mean),

Expectation of a RV x with PMF p_x is defined by

$$E[x] = \sum_{x \in X} x p_x(x),$$

where X is the set of all possible values then x can take.

Mean of a Bernoulli RV with $p_x(1) = p$ is

$$E[x] = 1.p + 0.(1-p) = p.$$

Consider a RV x with $p_x(2^k) = \frac{1}{2^k}$ $k = 1, 2, \dots$

$$E[x] = \sum_{k=1}^{\infty} 2^k \cdot \frac{1}{2^k} = \infty$$

Thus, expectation may not be always well defined. However there exists RVs

taking infinite number of values with finite mean.

Exercise 6.2. Find the expectation of a geometric RV with probability of heads equal to p , where $0 < p < 1$.

Lecture 7
(24 August 2023)

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

since $(X - E[X])^2$ can only take nonnegative values, the variance is always nonnegative.

The variance provides a measure of dispersion of X around its mean.

$\sigma_X = \sqrt{\text{Var}(X)}$ is called the standard deviation,

-Can variance be zero?

Yes!

We will get back to this soon.

P.T.O

$$\text{Example. } P_x(x) = \begin{cases} \frac{1}{9} & x \in [-4:4] \\ 0 & \text{o.w.} \end{cases}$$

Compute the variance $\text{Var}(x) = E[(x - E[x])^2]$.

$$E[x] = 0.$$

Let $y = (x - E[x])^2 \in \{0, 4, 9, 16\}$.

$$P_y(0) = \frac{1}{9}, \quad P_y(i^2) = \frac{2}{9}, \quad i \in [1:4].$$

$$\begin{aligned} \text{Var}(x) = E[y] &= 0 \cdot \frac{1}{9} + 1 \cdot \frac{2}{9} + 4 \cdot \frac{2}{9} + 9 \cdot \frac{2}{9} + 16 \cdot \frac{2}{9} \\ &= \frac{60}{9}. \end{aligned}$$

Expected value of Functions of Random Variables

Theorem.

Let x be a RV and $g(x)$ be a function of x . Then the expected value of $g(x)$ is given by

$$E[g(x)] = \sum_{x \in \Omega} x P_x(x).$$

Proof. Let $y = g(x)$, $y = \{y \in \mathbb{R} : \exists x \text{ s.t. } g(x) = y\}$.

$$\begin{aligned}
 E[g(x)] &= E[y] = \sum_{y \in Y} y p_y(y) \\
 &= \sum_{y \in Y} y \sum_{x: g(x)=y} p_x(x) \\
 &= \sum_{y \in Y} \sum_{x: g(x)=y} y p_x(x) \\
 &= \sum_{y \in Y} \sum_{x: g(x)=y} g(x) p_x(x) \\
 &= \sum_{x \in X} g(x) p_x(x).
 \end{aligned}$$

□

Using this

$$\text{Var}(x) = E[(x - E[x])^2] = \sum_x (x - E[x])^2 p_x(x)$$

For the example above $\text{Var}(x) = \sum_{x \in \{-4:4\}} (x-0)^2 p_x(x) = \frac{60}{9}$,
 $\rightarrow \text{Var}(x) = 0 \Rightarrow x = E[x]$, for all x . X is a constant.

Properties of mean and variance.

$$- y = ax + b$$

$$E[y] = a E[x] + b$$

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

$$= E[(x - \mu - \sigma E[x])^2]$$

$$= E[\sigma^2(x - E[x])^2] = \sigma^2 \text{Var}(x)$$

Variance in Terms of Moments

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

Mean and Variance of some common RVs

Bernoulli RV:

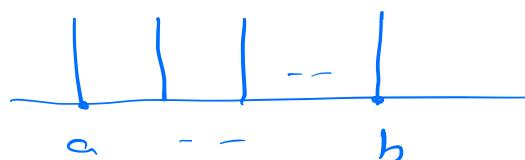
$$P_x(1) = p = 1 - P_x(0)$$

$$E[x] = p, \quad E[x^2] = p, \quad \text{Var}(x) = p - p^2 = p(1-p)$$

Discrete Uniform random variable:

$$P_x(k) = \frac{1}{b-a+1},$$

$$k \in [a; b]$$



$$E[x] = \frac{a+b}{2}$$

$$\text{Let } y \sim \text{uniform}\{1:n\}, \quad y = x - a + 1 \Rightarrow n = b - a + 1$$

$$\text{Var}(y) = \frac{n^2-1}{12}, \quad E[y^2] = \sum_{k=1}^n k^2 / n = \frac{1}{6} (n+1)(2n+1)$$

General case: $\text{Var}(y) = \frac{(b-a+1)^2 - 1}{12}$ for $n = b-a+1$,
 since $\text{Var}(x) = \text{Var}(x-a+1) = \text{Var}(y)$,
 $x-a+1 \in [1:n]$ for $n = b-a+1$,

Binomial random variable:

$$P_x(k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k \in [0:n],$$

$$E[x] = \sum_{k=0}^n k \binom{n}{k} p^k (1-p)^{n-k}$$

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

$$= n \binom{n-1}{k-1}.$$

$$E[x] = np \sum_{k=1}^n \binom{n-1}{k-1} p^{k-1} (1-p)^{n-k}$$

$$n' = n-1 \quad k' = k-1$$

$$= np \sum_{k'=0}^{n'} \binom{n'}{k'} p^{k'} (1-p)^{n'-k'}$$

$$= np (p+1-p)^{n'} = np$$

$$\begin{aligned}
 E[X^2] &= \sum_{k=0}^n k^2 \binom{n}{k} p^k (1-p)^{n-k} \\
 &= np \sum_{k=1}^n k \binom{n-1}{k-1} p^{k-1} (1-p)^{n-k} \\
 &\quad \left[n' = n-1 \quad k' = k-1 \right] \\
 &= np \sum_{k'=0}^{n-1} (k'+1) \binom{n'}{k'} p^{k'} (1-p)^{n'-k'} \\
 &= np(1 + (n-1)p)
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(x) &= E[X^2] - E[X]^2 \\
 &= np(1 + (n-1)p) - n^2 p^2 \\
 &= np(1 + np - p - np) = np(1-p).
 \end{aligned}$$

Exercise. Find the mean and variance of geometric random variable with probability of success p .

$$P_X(k) = (1-p)^{k-1} p \quad k \in \mathbb{N},$$

Exercise. Find the mean and variance of Poisson RV with PMF $P_X(k) = e^{-\lambda} \frac{\lambda^k}{k!}$, $k \in \mathbb{N}_0$.

Lecture 8
(31 August 2023)

Joint PMFs of Multiple Random Variables

Consider two RVs associated with the same experiment. The probabilities of the values that x and y can take are captured by the joint PMF of x and y , denoted P_{xy} .

$$X : \Omega \rightarrow \mathbb{R},$$

$$Y : \Omega \rightarrow \mathbb{R}.$$

$$\begin{aligned} P_{xy}(x, y) &= P(\{X=x\} \cap \{Y=y\}) \\ &\stackrel{\text{def}}{=} P(X=x, Y=y). \end{aligned}$$

$$P((X, Y) \in A) = \sum_{(x, y) \in A} P_{xy}(x, y).$$

Let $R_{xy} = \{(x, y) : P_{xy}(x, y) > 0\}$

$$R_x = \{x : P_x(x) > 0\}$$

$$R_y = \{y : P_y(y) > 0\}.$$

In fact, we can calculate the pmfs of x and y by using the formulas

$$P_x(x) = \sum_{y \in R_x} P_{xy}(x, y), \quad P_y(y) = \sum_{x \in R_y} P_{xy}(x, y),$$

$$\begin{aligned} P_x(x) &= P(x=x) \\ &= P(x=x \cap \Omega) \end{aligned}$$

$$= P(x=x \cap \bigcup_{y \in R_y} y=y)$$

$$= \sum_{y \in R_y} P(x=x \cap y=y)$$

$$= \sum_{y \in R_y} P_{xy}(x, y).$$

$$\text{Also } \sum_{(x, y) \in R_{xy}} P_{xy}(x, y) = \sum_{(x, y) \in R_{xy}} P(x=x, y=y)$$

$$= \sum_{(x, y) \in R_{xy}} P(A_{xy})$$

$$= P\left(\bigcup_{(x, y) \in R_{xy}} A_{xy}\right) = P(\Omega) = 1,$$

$A_{xy} - (x, y) \in R_{xy}$ forms a partition of Ω .

Functions of multiple Random Variables

Let $Z = g(\underline{x}, \underline{y})$ be a function of RVS \underline{x} and \underline{y} .

$$P_Z(z) = \sum_{(\underline{x}, \underline{y}): g(\underline{x}, \underline{y})=z} P_{\underline{x}, \underline{y}}(\underline{x}, \underline{y}).$$

Exercise • Show that

$$E[g(\underline{x}, \underline{y})] = \sum_{(\underline{x}, \underline{y})} g(\underline{x}, \underline{y}) P_{\underline{x}, \underline{y}}(\underline{x}, \underline{y}).$$

Linearity of Expectation

$$g(\underline{x}, \underline{y}) = \underline{x} + \underline{y}$$

$$E[\underline{x} + \underline{y}] = \sum_{(\underline{x}, \underline{y})} (\underline{x} + \underline{y}) P_{\underline{x}, \underline{y}}(\underline{x}, \underline{y})$$

$$= \sum_x x P_X(x) + \sum_y y P_Y(y) = E[X] + E[Y].$$

More than Two Random Variables

$$P_{XYZ}(x, y, z) = P(X=x, Y=y, Z=z).$$

$$P_X(x) = \sum_{(y, z) \in R_{YZ}} P_{XYZ}(x, y, z),$$

$$P_{YZ}(y, z) = \sum_{x \in R_X} P_{XYZ}(x, y, z),$$

$$E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i]$$

Example. Mean of the Binomial distribution

$$X = X_1 + X_2 + \dots + X_n \quad (\text{Binomial})$$

$$P(X_i = 1) = p = 1 - P(X_i = 0) \quad (\text{Bernoulli})$$

$$E[X] = np,$$

Conditioning on RV on an Event

The conditional PMF of a RV x conditioned on a particular event A with $P(A) > 0$ is defined by

$$\begin{aligned} P_{X|A}(x) &= P(X=x|A) \\ &= P(\{X=x\} \cap A) / P(A). \end{aligned}$$

$$\sum_x P_{X|A}(x) = \sum_x \frac{P(\{X=x\} \cap A)}{P(A)} = \frac{P(A)}{P(A)} = 1.$$

Example, Let X be the roll of a fair six-sided die and let A be the event that the roll is an even number.

$$\begin{aligned} P_{X|A}(k) &= \frac{P(X=k \text{ and } X \text{ is even})}{P(\text{roll is even})} \\ &= \begin{cases} \frac{1}{3}, & \text{if } k=2, 4, 6 \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

$$P_X(x) = \sum_{i=1}^n P(A_i) P_{X|A_i}(x) \text{ if } A_1, \dots, A_n \text{ form a partition of } \Omega.$$

Conditioning a RV on another RV

Consider two RVs x and y associated with the same experiment.

If we know that the value of y is some particular y with $P_y(y) > 0$, this provides partial knowledge about the value of x . This knowledge is captured by the conditional PMF $P_{x|y}$ of x given y ,

$$\begin{aligned} P_{x|y}(x|y) &= P(x=x|y=y) \\ &= \frac{P(x=x, y=y)}{P(y=y)} = \frac{P_{xy}(x,y)}{P_y(y)}. \end{aligned}$$

$$\sum_x P_{x|y}(x|y) = 1,$$

$$P_{xy}(x,y) = P_x(x) P_{y|x}(y|x)$$

$$= P_x(y) P_{x|y}(x|y)$$

Example. In each lecture a professor is asked 0 or 2 questions with equal probability $\frac{1}{3}$. He answers each question incorrectly with probability $\frac{1}{4}$ independent of other questions. Let x and y be the no. of questions the professor is asked and the no. of questions he answers wrong in a given lecture, respectively. Find $P_{xy}(x,y)$.

$$P_{xy}(0,k) = \begin{cases} \frac{1}{3} \cdot 1 = \frac{16}{48} & \text{if } k=0 \\ 0 & \text{if } k=1, 2. \end{cases}$$

$$P_{y|x}(0|1) = \frac{3}{4}, P_{y|x}(1|1) = \frac{1}{4}, P_{y|x}(2|1) = 0,$$

$$P_{y|x}(0|2) = \frac{9}{16}, P_{y|x}(1|2) = \frac{3}{8}, P_{y|x}(2|2) = \frac{1}{16}.$$

P.T.O

Conditional Expectation

For an event A with $P(A) > 0$, the condition expectation $E[X|A]$ is defined by

$$E[X|A] = \sum_x x P_{X|A}(x),$$

$$E[g(x)|A] = \sum_x g(x) P_{X|A}(x),$$

$$E[X|Y=y] = \sum_x x P_{X|Y}(x|y)$$

Total Expectation Theorem:

If A_1, \dots, A_n be disjoint events that form a partition of the sample space Ω , with $P(A_i) > 0$ for all i , then

$$E[X] = \sum_{i=1}^n P(A_i) E[X|A_i],$$

Proof.

$$\begin{aligned} E[X] &= \sum_x x P_X(x) \\ &= \sum_x x \left[\sum_{i=1}^n P(A_i) P_{X|A_i}(x) \right] \end{aligned}$$

$$= \sum_{i=1}^n P(A_i) \left[\sum_x x P_{x|A_i}(x) \right]$$

$$= \sum_{i=1}^n P(A_i) E[x|A_i],$$

Also

$$E[X] = \sum_y p_y(y) E[X|Y=y].$$

Similarly

$$E[X|B] = \sum_{i=1}^n P(A_i|B) E[X|A_i \cap B],$$

Lecture 9

(4 September 2023)

Conditional Expectation as a RV

$$E[x|y=y] = \sum_x x P_{x|y}(x|y)$$

Let $g(y) = E[x|y=y]$.

Consider the RV $g(y)$.

Define $E[x|y] := g(y)$.

Law of Iterated Expectations:

Let x and y be two random variables distributed according to a joint PMF P_{xy} .

Then $E[E[x|y]] = E[x]$.

Proof. $E[E[x|y]] = E[g(y)]$

$$= \sum_y g(y) P_y(y)$$

$$= \sum_y E[x|y=y] P_y(y)$$

$$= \sum_y \sum_x x P_{x|y}(x|y) P_y(y)$$

$$= \sum_x x \sum_y P_{x|y}(x|y) = \sum_x x P_x(x) = E[X],$$

Minimum Mean Square Error Estimator

$(x, y) \sim p_{xy}$. Given an observation y we need to estimate x .

Let $\hat{x} = f(y)$,

Theorem . $f(y) = g(y) = E[x|y=y]$ minimizes the expected squared error

$$E[(x - \hat{x})^2] = E[(x - f(y))^2].$$

Proof . $E[(x - f(y))^2]$

$$= \sum_y E[(x - f(y))^2 | y=y] P_y(y)$$

$$= \sum_y E[\tilde{x} + f(y) - 2f(y)x | y=y] P_y(y)$$

$$= \sum_y (f(y) - 2f(y)E[x|y=y] + E[\tilde{x}|y=y]) P_y(y)$$

Achieves the minimum at $f(y) = E[x|y=y]$.

Conditional Variance

$$\text{Var}(x|y=y) = E[x^2|y=y] - E[x|y=y]^2$$

$\text{Var}(x|y)$ is defined as a function of RV, y which takes a value $\text{Var}(x|y=y)$ when $y=y$.

Exercise. Prove that

$$\text{Var}(x) = E[\text{Var}(x|y)] + \text{Var}(E[x|y]).$$

Independence

- RV x is independent of the event A if $P(x=x \text{ and } A) = P_x(x)P(A)$, for all x ,
- (or) $P_{x|A}(x) = P_x(x)$, for all x if $P(A) > 0$.

Independence of RVs.

Two RVs x and y are independent if the events $\{x=x\}$ and $\{y=y\}$ are independent for all xy , i.e.,

$$P_{xy}(x,y) = P_x(x)P_y(y),$$

(conditional indep.) $P_{xy|A}(x,y) = P_{x|A}(x)P_{y|A}(y)$,

Example,

		y				
		1	2	3	4	
x		1	0	$\frac{1}{20}$	0	0
		2	0	$\frac{1}{20}$	$\frac{3}{20}$	$\frac{1}{20}$
		3	$\frac{1}{20}$	$\frac{4}{20}$	$\frac{1}{20}$	$\frac{2}{20}$
		4	$\frac{1}{20}$	$\frac{1}{20}$	$\frac{1}{20}$	0

P_{XY}

Are x and y
independent?

$$P_{XY}(1) = 0 \neq P_X(1)P_Y(1) > 0,$$

Are x and y independent conditioned
on $A = \{X \geq 3, Y \leq 2\}$?

$$P_{XY|A}(3,1) = \frac{2}{9} - P_{XY|A}(3,2) = \frac{4}{9},$$

$$P_{XY|A}(4,1) = \frac{1}{9} - P_{XY|A}(4,2) = \frac{2}{9},$$

$$P_{X|A}(3) = \frac{2}{3} = 1 - P_{X|A}(4),$$

$$P_{Y|A}(1) = \frac{1}{3} = 1 - P_{Y|A}(2).$$

So x and y are not independent.

But x and y are independent conditioned
on A.

Example. $x, y \in \{0, 1\}$. P_{xy} is joint PMF.

Suppose $P_{xy}(0, 1) = P_x(0)P_y(1)$, Are x and y independent?

$$\begin{aligned}P_{xy}(1, 1) &= P_y(1) - P_{xy}(0, 1) \\&= P_y(1) - P_x(0)P_y(1) \\&= P_y(1)P_x(1).\end{aligned}$$

Similarly $P_{xy}(x, y) = P_x(x)P_y(y)$ for all $x, y \in \{0, 1\}$.

Theorem. If x and y are independent random variables, then $E[xy] = E[x]E[y]$.

Proof. $E[xy] = \sum_{x,y} xy P_{xy}(x, y)$

$$\begin{aligned}&= \sum_{x,y} xy P_x(x)P_y(y) \\&= \left(\sum_x xP_x(x)\right)\left(\sum_y yP_y(y)\right) \\&= E[x]E[y].\end{aligned}$$

Exercise. If x and y are independent,

- Show that $g(x)$ and $h(y)$ are independent.
- Prove that $E[g(x)h(y)] = E[g(x)]E[h(y)]$.

Indicator Random Variables

For $A \subseteq \Omega$, $I_A(\omega) = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases}$.

Exercise

Two events A and B are independent if and only if I_A and I_B are independent RVs.

Independence of several Random Variables

Three RVs x, y and z are said to be independent if

$$P_{xyz}(x, y, z) = P_x(x) P_y(y) P_z(z) \text{ for all } x, y, z.$$

Independence of n RVs:

$$P_{x_1, x_2, \dots, x_n}(x_1, \dots, x_n) = \prod_{i=1}^n P_{x_i}(x_i),$$

Variance of sum of Independent Rvs

$$\text{Var}(x+y)$$

$$= \text{Var}(\underbrace{x - E[x]}_{\tilde{x}} + \underbrace{y - E[y]}_{\tilde{y}})$$

$$= \text{Var}(\tilde{x} + \tilde{y})$$

$$= E[(\tilde{x} + \tilde{y})^2]$$

$$= E[\tilde{x}^2 + \tilde{y}^2 + 2\tilde{x}\tilde{y}]$$

$$= \text{Var}(x) + \text{Var}(y) + 2\text{cov}(x, y),$$

where $\text{cov}(x, y) = E[(x - E[x])(y - E[y])]$.

If x and y are independent, $\text{cov}(x, y) = 0$.

$$\text{Var}(x+y) = \text{Var}(x) + \text{Var}(y).$$

- If x_1, x_2, \dots, x_n are independent Rvs

$$\text{Var}(\sum_{i=1}^n x_i) = \sum_{i=1}^n \text{Var}(x_i),$$

- Consider n independent coin tosses with probability of heads p .

Let $x_i = \mathbb{1}\{\text{i}^{\text{th}} \text{ coin toss is heads}\}$,
 $i \in [1:n]$,

$$\sum_{i=1}^n x_i = \text{no. of heads}$$

$$\text{Var}(\sum_{i=1}^n x_i) = np(1-p),$$

Lecture 10

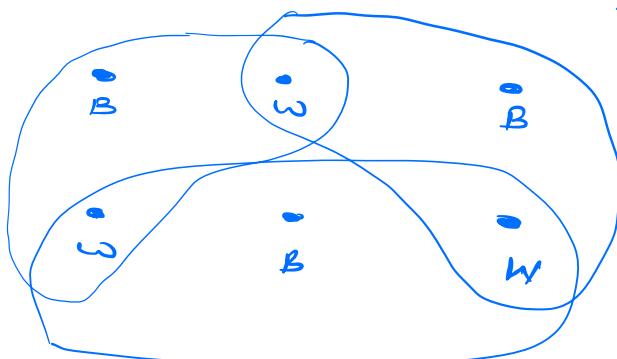
(7 September 2023)

Some Applications

Combinatorics and Graph Theory

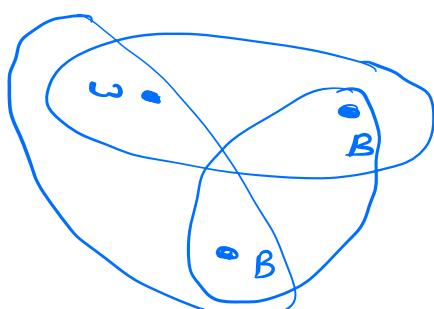
2-coloring: [Existence proof using union bound]

Let S be a set of some elements and $T_1, T_2, \dots, T_m \subseteq S$ be subsets s.t. $|T_i| = l$ for $i \in [1:m]$. Can we 2-color S (meaning assign each element of S a color) such that each T_i has elements of both colors (i.e., not monochromatic)?



$$l = 3 \quad m = 3$$

2-coloring exists



$$l = 3 \quad m = 3$$

No 2-coloring

Theorem - Given S and $T_1 T_2 \dots T_m \subseteq S$
 s.t. $|T_i| = l$ for $i \in [1:m]$, there exists a
 valid 2-coloring of S such that no T_i is
 monochromatic if $m < 2^{l-1}$.

For the examples above

$$l=3 \quad m=3 \Rightarrow 3 < 2^{3-1} = 7$$

$$l=2 \quad m=3 \Rightarrow 3 < 2^{2-1} = 3$$

Let $S = \{x_1, x_2, \dots, x_n\}$,

Randomly color each element of α black
 or white, independently and identically
 distributed, each with probability $\frac{1}{2}$.
 Let E_i be the event that T_i is
 monochromatic.

$$P(E_i) = \frac{1}{2^l} + \frac{1}{2^l} = \frac{1}{2^{l-1}}$$

$$P(\bigcup_{i=1}^m E_i) \leq \sum_{i=1}^m P(E_i) = \frac{m}{2^{l-1}} < 1,$$

$P(\text{each } T_i \text{ is not monochromatic})$

$$= 1 - P\left(\bigcup_{i=1}^m E_i\right)$$

$$\geq 1 - \frac{m}{2^{l-1}} > 0.$$

Because we have a non-zero probability this implies that there exists a 2-coloring of S that gives all m valid non-monochromatic sets T_i .

Let $A = \{\omega \in \Omega \text{ satisfying some property}\}$

$P(A) > 0 \Rightarrow \exists \omega \in \Omega \text{ satisfying that property.}$

\Rightarrow This is called the Probabilistic method.

Bipartite Subgraphs: [Existence proof using Linearity of Expectation]

Theorem, let G be a simple graph with vertex set $[1:n]$ and m edges. Then G contains a bipartite subgraph with more than $m/2$ edges.

Proof. Let us split the vertices of G into two disjoint nonempty sets A and B , ($2^{n-1}-1$ such partitions)

Then A and B span a bipartite subgraph H of G (we remove the edges within A and within B).

$$\Omega = \{ \text{all partitions } (A, B) \}$$

$$|\Omega| = 2^{n-1} - 1.$$

Let $X(H) = \text{no. of edges in } H$. $H \in \Omega$.

Number the edges from 1 through m and

let $X_i(A, B) = \begin{cases} 1 & \text{if the edge } i \text{ has one vertex} \\ & \text{in } A \text{ and one in } B \\ 0 & \text{otherwise.} \end{cases}$

$$P(X_i = 1) = \frac{2^{n-2}}{2^{n-1} - 1}, \text{ because we can}$$

get a subdivision of $[n]$ leading to $X_i = 1$ by first putting the two endpoints of the edge i to different subsets, then splitting

the remaining $(n-2)$ -element vertex set in any of 2^{n-2} ways.

$$E[x_i] = 1 \cdot P_{x_i}(1) + 0 \cdot P_{x_i}(0)$$

$$= \frac{2^{n-2}}{2^{n-1}-1} > \frac{1}{2} .$$

$$x = \sum_{i=1}^n x_i \Rightarrow E[x] = \sum_{i=1}^n E[x_i] > \frac{n}{2} .$$

$$\Rightarrow \exists x \text{ s.t. } x > \frac{n}{2} .$$

$$\Rightarrow \exists H \in \Omega \text{ s.t. } x(H) = x > \frac{n}{2} ,$$

Minimum Mean Squared Error (MMSE)

The MMSE estimate of the random variable x , given that we have observed y is given by $\hat{x}(y) = E[x|y]$.

In particular, the estimate function $\hat{x}(y) = E[x|y=y]$ achieves the minimum in

$$\min_{f(\cdot)} E[(x-f(y))^2] .$$

Proved in the last lecture.

Entropy (Uncertainty)

Consider a RV x_1 with PMF

$$P_{x_1}(0) = \frac{1}{2} = P_{x_1}(1),$$

consider another r.v. x_2 with PMF

$$P_{x_2}(0) = 0.9 - P_{x_2}(1) = 0.1,$$

which random variable's realization is hard to guess x_1 or x_2 ?

It appears that the uncertainty in x_1 is higher than that of x_2 . This uncertainty is exactly captured by Entropy.

For a RV, x with $P_x(1) = p = 1 - P_x(0)$

Entropy is defined as

$$H(x) = -P_x(1) \log P_x(1) - P_x(0) \log P_x(0)$$

$$= -p \log p - (1-p) \log (1-p),$$

$$=: h(p).$$

Is $H(x_2) \leq H(x_1)$? Yes!

$$H(X_1) = h(p) = -p \log p - (1-p) \log(1-p)$$

$$h'(p) = -p \cdot \frac{1}{p} - \log p + \frac{1-p}{1-p} + \log(1-p) = 0$$

$$\Rightarrow p = 1-p \Rightarrow p = \frac{1}{2}$$

$$h''(p) = \frac{-1}{p^2} - \frac{1}{(1-p)^2} < 0$$

$h(\frac{1}{2}) = 1$ is the maximum value of $h(p)$ according to the intuition X_1 has more uncertainty.

- Entropy is the fundamental quantity in information theory.

Lecture 11

(11 September 2023)

Recall that a random variable is a function from the sample space of a random experiment to the real numbers.

$$X : \Omega \rightarrow \mathbb{R},$$

Let Ω be an uncountable set and P be the associated probability law.

For a set $B \subseteq \mathbb{R}$, $\{X \in B\} = \{\omega : X(\omega) \in B\}$,
i.e., $P(X \in B) = P(\{\omega : X(\omega) \in B\})$.

Continuous Random Variable

A RV X is called continuous if there is a nonnegative function f_X , called the probability density function (PDF) of X such that

$$P(X \in B) = \int_B f_X(x) dx, \text{ for every } B \subseteq \mathbb{R}.$$

In particular, the probability that the value of x falls within an interval is

$$P(a \leq x \leq b) = \int_a^b f_x(x) dx.$$

(area under the graph of the PDF f_x)

$$P(X=a) = \int_a^a f_x(x) dx = 0$$

For this reason including or excluding the endpoints of an interval has no effect on its probability:

$$\begin{aligned} P(a \leq x \leq b) &= P(a < x \leq b) = P(a \leq x < b) \\ &= P(a < x < b). \end{aligned}$$

$$- f_x(x) \geq 0 \quad \forall x \in \mathbb{R}$$

$$- \int_{-\infty}^{\infty} f_x(x) dx = 1,$$

This means that the entire area under the graph of the PDF must be equal to 1.

To interpret the PDF note that for an interval $[x, x+\delta]$ with very small δ , we have

$$P(x \in [x, x+\delta]) = \int_x^{x+\delta} f_x(t) dt \approx f_x(x) \delta$$

$$f_x(x) \approx \frac{P(x \in [x, x+\delta])}{\delta}.$$

So we can view $f_x(x)$ as the 'probability mass per unit length' near x . It is important to realize that even though a PDF is used to calculate event probabilities, $f_x(x)$ is not the probability of any particular event. In particular, it is not restricted to be less than or equal to 1.

Example. Uniform RV.

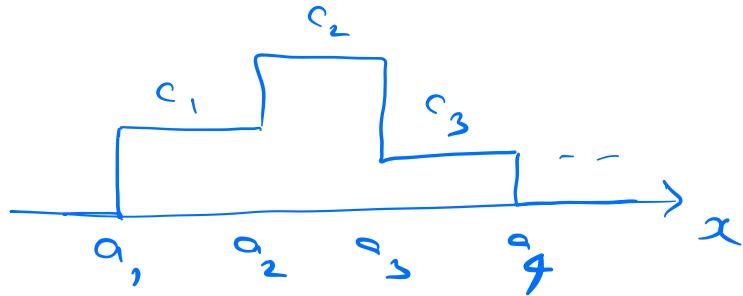
Consider a PDF

$$f_x(x) = \begin{cases} c & \text{if } a \leq x \leq b \\ 0 & \text{o.w.} \end{cases}$$

$$\int_a^b f_x(x) dx = 1 \Rightarrow c = \frac{1}{b-a}.$$

Example. piecewise constant PDF.

$$f_X(x) = \begin{cases} c_i & \text{if } a_i \leq x \leq a_{i+1}, \quad i=1, 2, \dots, n-1 \\ 0 & \text{o.w.} \end{cases}$$



$$\int_{-\infty}^{\infty} f_X(x) dx = 1 \Rightarrow \sum_{i=1}^{n-1} c_i (a_{i+1} - a_i) = 1,$$

— A PDF can take arbitrarily large values.
For example, uniform RV on [a, b] with $b-a < 1$.

Another example - $f_X(x) = \begin{cases} \frac{1}{\sqrt{x}} & \text{if } 0 < x \leq 1 \\ 0 & \text{o.w.} \end{cases}$

Even though $f_X(x)$ becomes infinitely large as x approaches zero, this is still a valid PDF because

$$\int_{-\infty}^{\infty} f_X(x) dx = \int_0^1 \frac{1}{\sqrt{x}} dx = 1.$$

Expectation,

The expected value of a continuous RV X is defined by

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx,$$

- A function of continuous RV is also a RV. It can be continuous or discrete.
 $y = g(x)$.

$$y = g(x) = x \text{ continuous RV}$$

$$g(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}, \quad y = g(x) \text{ is discrete.}$$

Theorem (Expected value rule for functions of RVs)

$$E[g(x)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx,$$

Proof. We will present a proof assuming that $g(x)$ is non-negative.

For a nonnegative continuous RV y

$$E[y] = \int_0^{\infty} P(y > y) dy,$$

$$\begin{aligned}
 \int_0^\infty P(Y > y) dy &= \int_0^\infty \int_y^\infty f_Y(t) dt dy \\
 &= \int_{t=0}^\infty \left(\int_{y=0}^t dy \right) f_Y(t) dt \\
 &= \int_{t=0}^\infty t f_Y(t) dt = E[Y],
 \end{aligned}$$

Now, for any function g s.t. $g(x) \geq 0$

$$\begin{aligned}
 E[g(X)] &= \int_0^\infty P(g(X) > y) dy \\
 &= \int_0^\infty \int_{x: g(x) > y} f_X(x) dx dy \\
 &= \int_{x=-\infty}^\infty +_X(x) \int_{y=0}^{g(x)} dy dx \\
 &= \int_{x=-\infty}^\infty g(x) f_X(x) dx,
 \end{aligned}$$

Exercise. Prove the above for a general

real-valued function g .

Hint. Show that $E[x] = \int_0^\infty p(x > x) dx - \int_0^\infty p(x < -x) dx$.

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

Properties

$$- Y = ax + b$$

$$E[y] = aE[x] + b$$

$$\text{Var}(y) = a^2 \text{Var}(x)$$

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

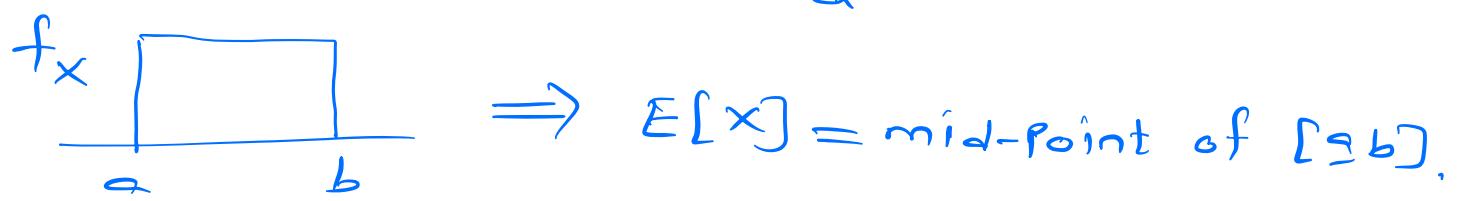
$$= \int_{-\infty}^{\infty} (x - E[x])^2 f_x(x) dx$$

$$= E[x^2] - (E[x])^2$$

Example, Mean and variance of the Uniform Random variable.

$$f_x(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{o.w.} \end{cases}$$

$$E[x] = \int_{-\infty}^{\infty} x f_X(x) dx = \int_a^b \frac{1}{b-a} x dx = \frac{a+b}{2}.$$



$$E[x^2] = \int_a^b x^2 \cdot \frac{1}{b-a} dx = \frac{a^2 + b^2 + ab}{3}$$

$$\begin{aligned} \text{Var}(x) &= E[x^2] - (E[x])^2 = \frac{a^2 + b^2 + ab}{3} - \frac{(a+b)^2}{4} \\ &= \frac{(b-a)^2}{12}. \end{aligned}$$

Exponential Random Variable

An exponential RV has a PDF of the form

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{if } x \geq 0 \\ 0, & \text{otherwise,} \end{cases}$$

where λ is a positive parameter characterizing the PDF.

This is a valid PDF because

$$\int_0^{\infty} \lambda e^{-\lambda x} dx = \lambda \left[\frac{e^{-\lambda x}}{-\lambda} \right]_0^{\infty} = 1,$$

- Exp. RV can model the amount of time until an incident of interest takes place.
 The probability that x exceeds a certain value decreases exponentially.
 Indeed, for any $a \geq 0$, we have

$$P(x \geq a) = \int_a^{\infty} \lambda e^{-\lambda x} dx = e^{-\lambda a}.$$

$$\begin{aligned} E[x] &= \int_0^{\infty} x \lambda e^{-\lambda x} dx \\ &= x \Big| \frac{e^{-\lambda x}}{-\lambda} \Big|_0^{\infty} + \int_0^{\infty} e^{-\lambda x} dx \\ &= 0 + \frac{e^{-\lambda x}}{-\lambda} \Big|_0^{\infty} = \frac{1}{\lambda}. \end{aligned}$$

$$\begin{aligned} E[x^2] &= \int_0^{\infty} x^2 \lambda e^{-\lambda x} dx \\ &= x^2 \Big(-e^{-\lambda x} \Big) \Big|_0^{\infty} + \int_0^{\infty} 2x \cdot e^{-\lambda x} dx \\ &= 0 + \frac{2}{\lambda} E[x] = \frac{2}{\lambda^2}. \end{aligned}$$

$$\text{Var}(x) = \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{1}{\lambda^2},$$

Lecture 12

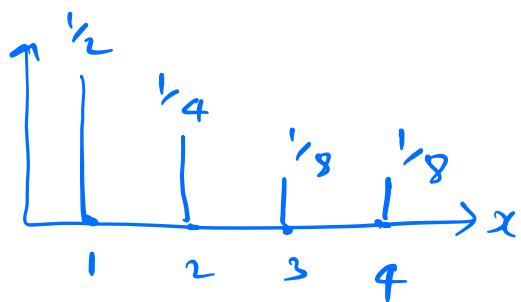
(14 September 2023)

Cumulative Distribution Functions

We would like to describe all kinds of RVs with a single mathematical concept. This is accomplished with CDF.

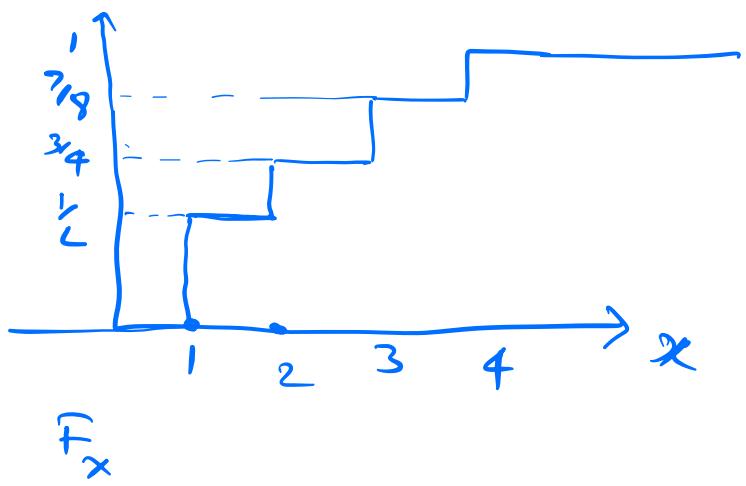
$$F_X(x) = \begin{cases} \sum_{k \leq x} P_X(k) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^x f_X(t) dt & \text{if } X \text{ is continuous} \end{cases}$$

- Any random variable associated with a given probability model has a CDF regardless of whether it is discrete or continuous. This is because $\{X \leq x\}$ is always an event and therefore has a well-defined probability.

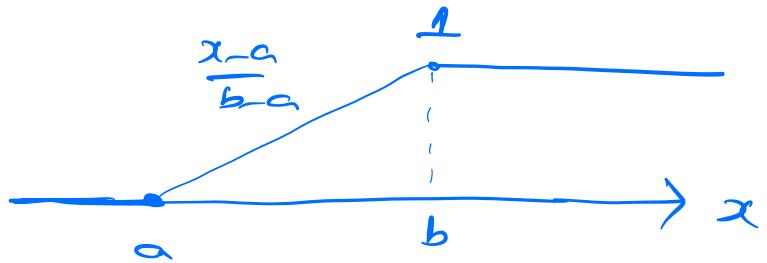
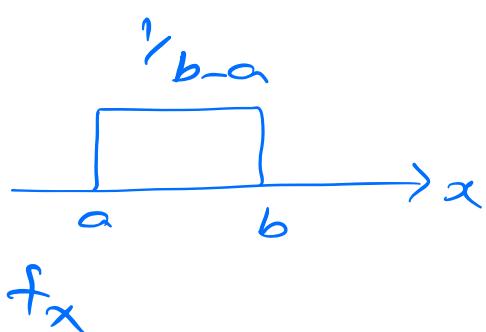


P_X

Discrete RV



F_X



Continuous RV

- CDF of discrete RV has jumps occurring at the values of positive probability mass. However F_x is continuous from the right.

Properties of a CDF

1) F_x is monotonically non-decreasing:

if $x \leq y$, then $F_x(x) \leq F_x(y)$.

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

$$\leq P(X \leq y) = F_y(y).$$

2) $\lim_{x \rightarrow -\infty} F_x(x) = 0$ $\lim_{x \rightarrow \infty} F_x(x) = 1$,

$$A_n = \{X \leq n\} \quad A_1 \subseteq A_2 \subseteq \dots$$

$$\bigcup_{i=1}^{\infty} A_i = \Omega.$$

$$\lim_{x \rightarrow \infty} F_x(x) = \lim_{n \rightarrow \infty} P(X \leq n)$$

Every sub-sequence of a monotonic sequence converges to the same limit

$$= P\left(\bigcup_{i=1}^{\infty} (X \leq i)\right)$$

$$= P(\Omega) = 1.$$

Let $A_n = \{X \leq -n\}$, $A_1 \supseteq A_2 \supseteq \dots$

$$\bigcap_{i=1}^{\infty} A_i = \emptyset$$

$$\lim_{x \rightarrow -\infty} F_x(x) = \lim_{x \rightarrow -\infty} P(X \leq x)$$

$$= \lim_{n \rightarrow -\infty} P(X \leq n)$$

$$= \lim_{n \rightarrow \infty} P(X \leq -n)$$

$$= P\left(\bigcap_{n=1}^{\infty} (X \leq -n)\right) = P(\emptyset) = 0,$$

3) F_x is right continuous.

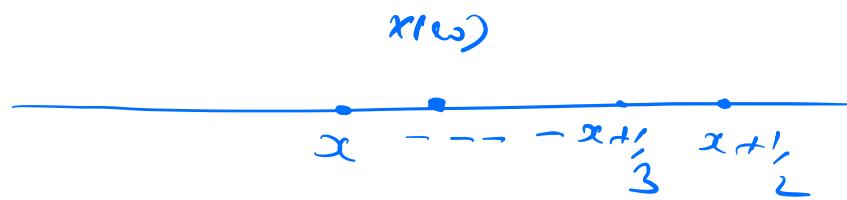
$$\lim_{\varepsilon \rightarrow 0^+} F_x(x+\varepsilon) = F_x(x),$$

$$\lim_{\varepsilon \rightarrow 0^+} F_x(x+\varepsilon) = \lim_{n \rightarrow \infty} F_x(x+\frac{1}{n})$$

$$= \lim_{n \rightarrow \infty} P(X \leq x + \frac{1}{n})$$

$$A_n = \{X \leq x + \frac{1}{n}\}, \quad A_1 \supseteq A_2 \supseteq \dots$$

$$\bigcap_{i=1}^{\infty} A_i = \bigcap_{i=1}^{\infty} \{X \leq x + \frac{1}{i}\} = \{X \leq x\}$$



$$x(\omega) \leq x + \frac{1}{i} \quad \forall i \in \mathbb{N}[n] \Rightarrow x(\omega) \leq x,$$

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0^+} F_x(x+\varepsilon) &= P\left(\bigcap_{i=1}^{\infty} A_i\right) \\ &= P(X \leq x) = f_x(x) \end{aligned}$$

- If X is discrete, then F_X is a piecewise constant function of x ,
- If X is continuous, then F_X is a continuous function of x ,

$$\lim_{\varepsilon \rightarrow 0^-} F_x(x+\varepsilon) = \lim_{\varepsilon \rightarrow 0^-} P(X \leq x+\varepsilon)$$

$$= \lim_{n \rightarrow \infty} P(x \leq x - \epsilon_n)$$

$$A_1 \subseteq A_2 \subseteq \dots \subseteq A_n = \{x \leq x - \epsilon_n\}$$

$$= P(\bigcup_{i=1}^{\infty} A_i)$$

$$= P(x < x) = P(x \leq x) - P(x = x) \\ = F_x(x), \quad = e$$

4) If x is discrete and takes integer values

$$F_x(k) = \sum_{i=-\infty}^k P_x(i)$$

$$P_x(k) = P(x \leq k) - P(x \leq k-1)$$

$$= F_x(k) - F_x(k-1), \text{ for all } k \in \mathbb{Z}.$$

If x is continuous

$$F_x(x) = \int_{-\infty}^x f_x(t) dt$$

$$f_x(x) = \frac{d F_x(x)}{dx} \text{ for all } x \text{ s.t. } f \text{ is continuous at } x,$$

[By Fundamental theorem of Calculus]

Lecture 13
 (18 September 2023)

- Sometimes, in order to calculate the PMF or PDF of a discrete or continuous r.v., respectively, it is more convenient to first calculate the CDF.

Example. x_1, x_2 and x_3 are independent r.v.s with PMF $P_{x_i}(k) = \frac{1}{10}$, $k \in \{1:10\}$. Find the PMF of $x = \max\{x_1, x_2, x_3\}$.

$$\begin{aligned}
 \text{Proof. } F_x(k) &= P(\max\{x_1, x_2, x_3\} \leq k) \\
 &= P(x_1 \leq k, x_2 \leq k, x_3 \leq k) \\
 &= P(x_1 \leq k) P(x_2 \leq k) P(x_3 \leq k) \\
 \text{since } x_1, x_2 \text{ and } x_3 \text{ are independent r.v.s} \\
 \Leftrightarrow P(x_1 \leq k, x_2 \leq k, x_3 \leq k) &= \sum_{a,b,c \leq k} P_{x_1, x_2, x_3}(a, b, c) = \sum_{a,b,c \leq k} P_{x_1}(a) P_{x_2}(b) P_{x_3}(c) \\
 &= P(x_1 \leq k) P(x_2 \leq k) P(x_3 \leq k).
 \end{aligned}$$

$$\text{so } P(x \leq k) = \left(\frac{k}{10}\right)^3.$$

$$P_x(k) = F_x(k) - F_x(k-1) = \left(\frac{k}{10}\right)^3 - \left(\frac{k-1}{10}\right)^3$$

Q) Let A and B be events with $P(A) > 0$ and $P(B) > 0$. B suggests A if $P(A|B) > P(A)$ and B does not suggest A if $P(A|B) < P(A)$.

(a) Show that B suggests A if and only if A suggests B.

$$P(A|B) > P(A) \Rightarrow \frac{P(B|A) P(A)}{P(B)} > \cancel{P(A)}$$

$$\Rightarrow P(B|A) > P(B),$$

(b) If $P(B^c) > 0$, show that B suggests A if and only if B^c does not suggest A.

$$P(A|B) > P(A) \quad P(A|B^c) = \frac{P(A \cap B^c)}{P(B^c)}$$

$$= \frac{P(A) - P(A \cap B)}{P(B^c)}$$

$$< P(A) P(B^c) / \cancel{P(B^c)}$$

Q) $X = \sum_{i=1}^n x_i$, x_i 's are Bernoulli RVs (need not be independent). Show that

Solution, $E[X^2] = \sum_{i=1}^n p_{x_i}(1) E[X|x_i=1]$,

$$E[X^2] = E\left[\left(\sum_{i=1}^n x_i\right)^2\right]$$

$$= \sum_{i=1}^n E[x_i^2] + \sum_{\substack{i,j \\ i \neq j}} E[x_i x_j]$$

$$= \sum_{i=1}^n E[x_i] + \sum_{i \neq j} E[x_i x_j]$$

$$= \sum_{i=1}^n E\left[x_i \left(1 + \sum_{j \neq i} x_j\right)\right]$$

$$= \sum_{i=1}^n E\left[x_i \left(1 + (x - x_i)\right)\right]$$

$$= \sum_{i=1}^n E\left[x_i + x x_i - x_i^2\right]$$

$$= \sum_{i=1}^n E[xx_i] + \sum_{i=1}^n \underbrace{\left(E[x_i] - E[x_i^2]\right)}_{=0}$$

$$= \sum_{i=1}^n p_{x_i}(1) E[X|x_i=1].$$

Q) Suppose that we flip a coin n times to obtain n random bits. Consider all $m = \binom{n}{2}$ pairs of these bits in some order. Let y_i be the exclusive-or of the i th pair of bits and let $y = \sum_{i=1}^m y_i$ be the number of y_i 's that equal 1.

(a) Show that each y_i is 0 with probability $\frac{1}{2}$ and 1 with probability $\frac{1}{2}$.

(b) Show that y_i 's are not mutually independent.

(c) Show that $E[y_i y_j] = E[y_i]E[y_j]$ if $i \neq j$.

(d) Find $\text{Var}(y)$.

Solution, x_1, x_2, \dots, x_n be the fair coin flips.

(a) $P(x_1 \oplus x_2 = 0) = P(x_1 = x_2)$

$$= P_{x_1, x_2}(00) + P_{x_1, x_2}(11) = \frac{1}{2}$$

$$P(x_1 \oplus x_2 = 1) = \frac{1}{2}.$$

so, $P(y_i = 0) = P(y_i = 1) = \frac{1}{2}$ since

$y_i = x_j + x_k$ for some $j \neq k$.

(b) Consider $y_i = x_1 \oplus x_2, y_j = x_1 \oplus x_3, y_k = x_1 \oplus x_3$.

y_i, y_j and y_k are not independent.

$$P(y_i = 1, y_j = 1, y_k = 1)$$

$$= 0 \neq \frac{1}{8} = P(y_i = 1)P(y_j = 1)P(y_k = 1).$$

(c) $E[y_i] = \frac{1}{2}$.

$$E[y_i y_j] = E[(x_i \oplus x_j)(x_k \oplus x_l)]$$

for $i \neq j, k \neq l, \{i, j\} \neq \{k, l\}$

If $\{i, j\} \cap \{k, l\} = \emptyset, E[y_i y_j] = E[x_i \oplus x_j]E[x_k \oplus x_l]$

$$= \frac{1}{4}.$$

SUPPOSE $|\{i, j\} \cap \{k, l\}| = 1$.

$$E[(x_1 \oplus x_2)(x_2 \oplus x_3)]$$

$$= \sum_{x_1 x_2 x_3} (x_1 \oplus x_2)(x_2 \oplus x_3) P_{x_1}(x_1) P_{x_2}(x_2) P_{x_3}(x_3)$$

$$= \frac{1}{8} \sum_{x_1 x_2 x_3} (x_1 \oplus x_2)(x_2 \oplus x_3)$$

0 1 0
1 0 1

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

$$\therefore E[x_i x_j] = E[x_i] E[x_j] = \frac{1}{4}.$$

$$(d) \text{Var}(y) = E[y^2] - (E[y])^2$$

$$= \sum_{i=1}^m E[x_i^2] + \sum_{i \neq j} E[x_i x_j] - (m \cdot \frac{1}{2})^2$$

$$= \frac{m}{2} + \frac{m(m-1)}{4} - \frac{m^2}{4} = \frac{n(n-1)}{8}.$$

(Q) We say that α is a median of a RV x if $P(x \leq \alpha) \geq \frac{1}{2}$ and $P(x \geq \alpha) \geq \frac{1}{2}$.

It is possible for the median to be non-unique, with all values in an interval satisfying the definition,

(a) Let $x \in \{0, 1, 2\}$ with probabilities p_0, p_1 and p_2 , respectively. Find the median of x for each of the cases below.

$$(i) P_0 = 0.2 \quad P_1 = 0.4 \quad P_2 = 0.4$$

$$(ii) P_0 = 0.2 \quad P_1 = 0.2 \quad P_2 = 0.6$$

$$(iii) P_0 = 0.2 \quad P_1 = 0.3 \quad P_2 = 0.5$$

(b) Suppose x is a continuous RV with PDF f_x s.t. $f_x(x) = \begin{cases} 0.5 & \text{for } 0 \leq x \leq 0.5 \\ 0 & \text{for } 0.5 < x \leq 1. \end{cases}$

We know that $f_x(x) > 0$ for all $x > 1$, and $f_x(x) = 0$ for all $x < 0$, but is otherwise unknown. Find the median.

Solution, (a) (i) $\alpha = 1$ is the median.

$$P(X \leq 1) = 0.6 \geq \frac{1}{2}$$

$$P(X \geq 1) = 0.8 \geq \frac{1}{2}$$

(ii) $\alpha = 2$ is the median.

$$P(X \leq 2) = 1 \geq \frac{1}{2}$$

$$P(X \geq 2) = 0.6 \geq \frac{1}{2}$$

(iii) Let $\alpha \in [1, 2]$.

$$P(X \leq \alpha) = \begin{cases} 0.5 & \text{if } \alpha \in [1, 2], \\ 1 & \text{if } \alpha = 2 \end{cases}$$
$$\geq \frac{1}{2}$$

$$P(X \geq \alpha) = \begin{cases} 0.5 & \text{if } \alpha \in (1, 2] \\ 0.8 & \text{if } \alpha = 1 \end{cases} \geq \frac{1}{2}$$

(b) $P(X \leq 0.5) = 0.5$

Let $\alpha \in [0.5, 1]$.

$$P(X \leq \alpha) = 0.5 \geq \frac{1}{2}$$

$$P(X \geq \alpha) = \int_{\alpha}^{\infty} f(x) dx = \frac{1}{2} \geq \frac{1}{2}$$

Q) Suppose x and y are two independent RVs such that

$$E[x^4] = 2 \quad E[y^2] = 1 \quad E[x^2] = 1 \quad \text{and} \quad E[y] = 0.$$

Find $\text{Var}(x\tilde{y})$.

$$\begin{aligned}\text{Var}(x\tilde{y}) &= E[x^4 y^2] - (E[x\tilde{y}])^2 \\ &= 2 \cdot 1 - (1 \cdot 0)^2 = 2.\end{aligned}$$

Lecture 14
(25 September 2023)

Geometric and Exponential CDFs

CDF, defined for any type of RV, provides a convenient means for exploring the relations between continuous and discrete random variables. We explore the relation between geometric and exponential RVs.

Let X be a geometric RV with success probability p , i.e., X is the no. of trials until the first success in a sequence of independent Bernoulli trials, where the probability of success at each trial is p .

$$P(X=k) = \underbrace{(1-p)}_n^{k-1} p \quad k=1, 2, 3, \dots$$

$$F_x^G(n) = \sum_{k=1}^n P_X(k) = p \cdot \frac{1 - (1-p)^n}{p} = 1 - (1-p)^n.$$

If X is exponential RV

$$F_x^E(x) = \begin{cases} \int_0^x e^{-\lambda t} dt = 1 - e^{-\lambda x} & \text{for } x > 0, \\ 0 & \text{for } x \leq 0. \end{cases}$$

Consider the values of $F_x^E(x)$ at $x = n\delta$, $n = 1, 2, \dots$

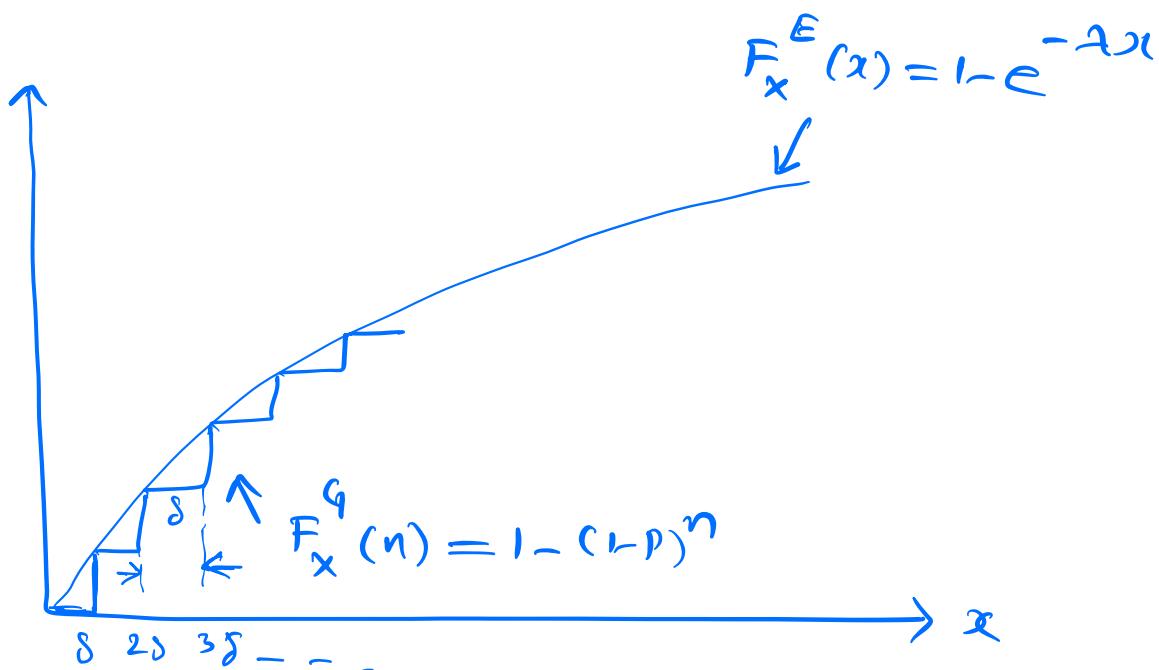
$$1 - e^{-\lambda n\delta} = 1 - (1-p)^n$$

$$e^{-\lambda\delta}$$

$$= 1-p \Rightarrow \delta = -\ln(1-p)/\lambda.$$

Then we see that the values of the exponential and the geometric CDFs are equal whenever $x = n\delta$, with $n = 1, 2, \dots$ i.e.,

$$F_x^E(n\delta) = F_x^G(n), \quad n = 1, 2, \dots$$



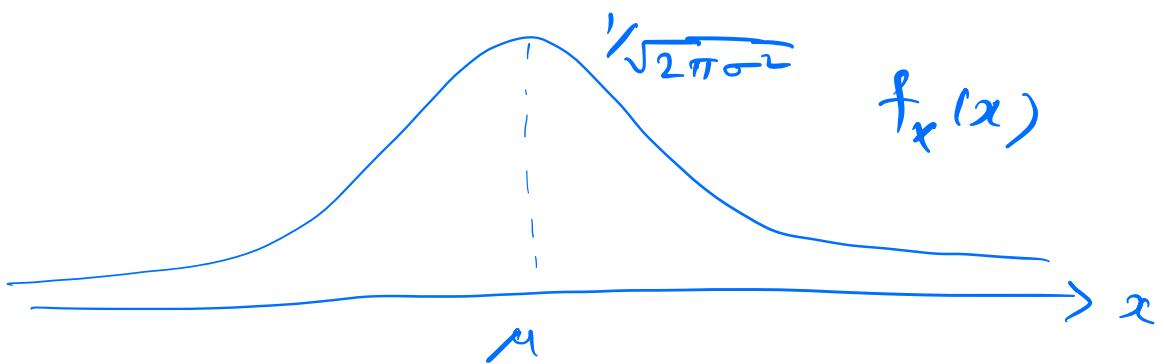
Suppose now we toss a biased coin very quickly (every δ seconds where $\delta \ll 1$) with a small probability of heads (equal to $p = 1 - e^{-1/\delta}$). Then the first time to obtain a head (a geometric random variable with parameter p) is a close approximation to an exponential r.v. with parameter λ , in the sense that the corresponding CDFs are very close to each other as shown in the above figure.

Gaussian Random Variables

A continuous RV x is said to be Gaussian or normal if it has a PDF of the form

$$f_x(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2},$$

where $\mu \in \mathbb{R}$, $\sigma \in (0, \infty)$.



Is it a valid PDF?

$$\int_{-\infty}^{\infty} f_x(x) dx = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2} dx$$

$$\text{Let } t = \frac{x-\mu}{\sigma}$$

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt = 1$$

$$1 = \int_{x=-\infty}^{\infty} \int_{y=-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2} e^{-y^2/2} dx dy$$

$$= \frac{1}{2\pi} \int_{x=-\infty}^{\infty} \int_{y=-\infty}^{\infty} e^{-(x^2+y^2)/2} dx dy$$

$x = r \cos \theta$, $y = r \sin \theta$ gives

$$\frac{1}{2\pi} \int_0^\infty \int_0^{2\pi} e^{-r^2/2} r d\theta dr$$

$r=0 \quad \theta=0$

$$= \frac{1}{2\pi} \cancel{(2\pi)} \int_0^\infty e^{-r^2/2} r dr$$

$$= -e^{-r^2/2} \Big|_{r=0}^\infty = 1.$$

$$\therefore I = \int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = 1,$$

$$E[x] = \int_{-\infty}^\infty x \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2} dx$$

$$= \int_{-\infty}^\infty (\sigma t + \mu) \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$$

$$= \mu + \int_{-\infty}^\infty \sigma t \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$$

$$\underbrace{\qquad}_{=\frac{\sigma}{\sqrt{2\pi}}} \qquad \qquad \qquad = \frac{\sigma}{\sqrt{2\pi}} - e^{-t^2/2} \Big|_{-\infty}^\infty = 0$$

$$\text{So mean} = E[x] = \mu.$$

Intuitively note that the PDF is symmetric around μ , so the mean can only be μ .

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

$$E[x^2] = \int_{-\infty}^{\infty} x^2 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

$$= \int_{-\infty}^{\infty} (\sigma t + \mu)^2 \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt$$

$$= \mu^2 + \sigma^2 \int_{-\infty}^{\infty} t^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt + 2\mu\sigma \cdot 0$$

$$= \mu^2 + \sigma^2 \left(\left. -t \cdot e^{-\frac{t^2}{2}} \right|_{-\infty}^{\infty} + \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt \right)$$

$$= \mu^2 + \sigma^2$$

$$\therefore \text{Var}(x) = \sigma^2.$$

A normal RV has several special properties.

Theorem. If x is a normal random variable with mean μ and variance σ^2 , and if a, b are scalars then the random variable

$y = ax + b$ is also normal with mean $E[y] = aE[x] + b$ and $\text{var}(y) = a^2\sigma^2$.

Before proving this theorem let us find the PDF of a linear function of any RV.

Let x be a continuous RV with PDF f_x , and let $y = ax + b$.

$$\begin{aligned} P(y \leq y) &= P(ax + b \leq y) \\ &= P(ax \leq y - b) \end{aligned}$$

$$= \begin{cases} P(x \leq \frac{y-b}{a}) & \text{if } a > 0 \\ P(x \geq \frac{y-b}{a}) & \text{if } a < 0 \end{cases}$$

$$f_y(y) = \frac{dF_y(y)}{dy}$$

$$= \begin{cases} \frac{1}{a} f_x\left(\frac{y-b}{a}\right) & \text{if } a > 0 \\ -\frac{1}{a} f_x\left(\frac{y-b}{a}\right) & \text{if } a < 0 \end{cases}$$

$$X \sim f_x \Rightarrow ax+b \sim \frac{1}{|a|} f_x\left(\frac{y-b}{a}\right).$$

Proof of Theorem.

$$\begin{aligned} f_y(y) &= \frac{1}{|a|} f_x\left(\frac{y-b}{a}\right) \\ &= \frac{1}{|a|} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\left(\frac{y-b}{a}-\mu\right)^2/2\sigma^2} \\ &= \frac{1}{\sqrt{2\pi|a|^2\sigma^2}} e^{-\frac{(y-b-\mu a)^2}{2a^2\sigma^2}} \end{aligned}$$

This is a normal PDF with mean $ax+b$ and variance $a^2\sigma^2$. Thus $y=ax+b$ is a normal RV.

A Linear Function of an Exponential RV.

Suppose x is an exponential RV with PDF

$$f_x(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{if } x \geq 0 \\ 0, & \text{otherwise.} \end{cases}$$

Let $y = ax + b$.

$$\begin{aligned} f_y(y) &= \frac{1}{|a|} f_x\left(\frac{y-b}{a}\right) \\ &= \begin{cases} \frac{\lambda}{|a|} e^{-\lambda \frac{y-b}{a}}, & \text{if } \frac{y-b}{a} \geq 0 \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

If $b=0$ and $a>0$ then y is an exponential RV with parameter λ/a . In general, y need not be exponential.

— A normal RV y with zero mean and unit variance is said to be a standard normal RV. Its CDF is denoted by Φ :

$$\Phi(y) = P(Y \leq y) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^y e^{-t^2/2} dt.$$

$$\Phi(-y) = P(y \leq -y)$$

$$= P(y \geq y) = 1 - P(y \leq y)$$

$$= 1 - \Phi(y).$$

- Let x be a normal RV with mean μ and variance σ^2 . We "standardize" x by defining a new RV y given by

$$y = \frac{x - \mu}{\sigma}.$$

$$E[y] = 0, \quad \text{Var}(y) = 1,$$

Example, $X \sim N(\mu = 60, \sigma^2 = 20^2)$.

$$\begin{aligned} P(x \geq 80) &= P\left(\frac{x - 60}{20} \geq \frac{80 - 60}{20} = 1\right) \\ &= P(y \geq 1) = \Phi(1). \end{aligned}$$

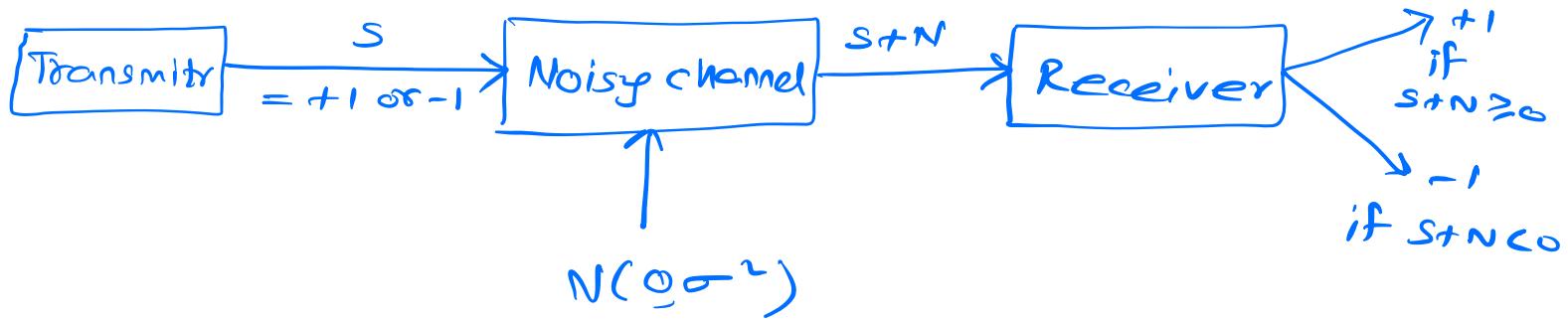
In General

$$P(x \leq x) = \Phi\left(\frac{x - \mu}{\sigma}\right).$$

Normal RVs are often used in signal processing and communications to model noise.

Example (Signal Detection).

A binary message is transmitted as a signal s , which is either $+1$ or -1 .



What is the probability of error?

An error occurs whenever -1 is transmitted and the noise N is at least 1 so that $s+N = -1+N \geq 0$ or whenever $+1$ is transmitted and the noise N is smaller than -1 so that $s+N = 1+N < 0$.

When $s = -1$ $P(N \geq 1)$ is the error probability.

$$\begin{aligned}
 P(N \geq 1) &= 1 - P(N < 1) = 1 - P\left(\frac{N}{\sigma} < \frac{1}{\sigma}\right) \\
 &= 1 - \Phi\left(\frac{1}{\sigma}\right).
 \end{aligned}$$

When $s = +1$ $P(N \leq -1) = P(N \geq 1)$

$$= 1 - \Phi\left(\frac{-1}{\sigma}\right).$$

Joint PDFs of Multiple Random Variables

We say that two random variables associated with the same experiment are jointly continuous if there exists a non-negative function f_{xy} , called as joint PDF such that

$$P((x,y) \in B) = \iint_{(x,y) \in B} f_{xy}(x,y) dx dy,$$

for every subset B of \mathbb{R}^2 .

Exercise. If x and y are jointly continuous prove that they are individually continuous also.

Lecture 15

(5 September 2023)

Joint PDFs of Multiple Random Variables

We say that two random variables associated with the same experiment are jointly continuous if there exists a non-negative function f_{xy} , called as joint PDF such that

$$P((x,y) \in B) = \iint_{(x,y) \in B} f_{xy}(x,y) dx dy,$$

for every subset B of \mathbb{R}^2 .

In particular when B is a rectangle

$$P(a \leq x \leq b, c \leq y \leq d) = \int_{y=c}^d \int_{x=a}^b f_{xy}(x,y) dx dy.$$

$$B = \mathbb{R}^2 \Rightarrow \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{xy}(x,y) dx dy = 1.$$

To interpret the joint PDF, we let δ be a small positive number and consider the probability of a small rectangle.

$$P(a \leq x \leq a+\delta, c \leq y \leq c+\delta)$$

$$= \int_c^{c+\delta} \int_a^{a+\delta} f_{xy}(x,y) dx dy \approx f_{xy}(a,c) \delta^2$$

so we can view $f_{xy}(a,c)$ as the "probability per unit area" in the vicinity of (a,c) .

$$P(x \in A) = P(x \in A \text{ and } y \in (-\infty, \infty))$$

$$= \int_{x \in A} \int_{y=-\infty}^{\infty} f_{xy}(x,y) dy dx.$$

If x and y are jointly continuous, they are individually continuous.

$$\text{Comparing with } P(x \in A) = \int_A f_x(x) dx$$

marginal PDF f_x of x is given by

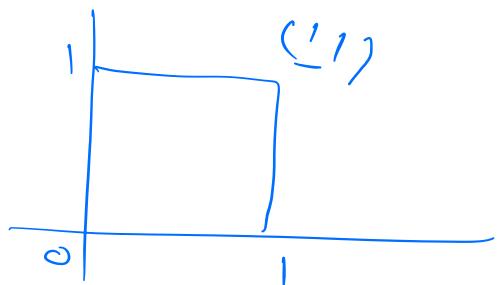
$$f_x(x) = \int_{-\infty}^{\infty} f_{xy}(x,y) dy,$$

$$\text{Similarly } f_y(y) = \int_{-\infty}^{\infty} f_{xy}(x,y) dx.$$

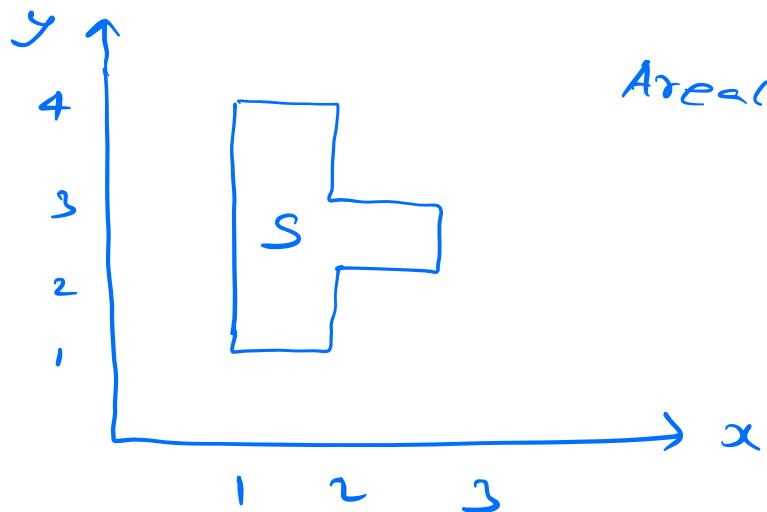
Example (Two-Dimensional Uniform PDF),

$$f_{xy}(x,y) = \begin{cases} c & \text{if } 0 \leq x \leq 1 \text{ and } 0 \leq y \leq 1 \\ 0 & \text{o.w.} \end{cases}$$

where c is a constant.



Example



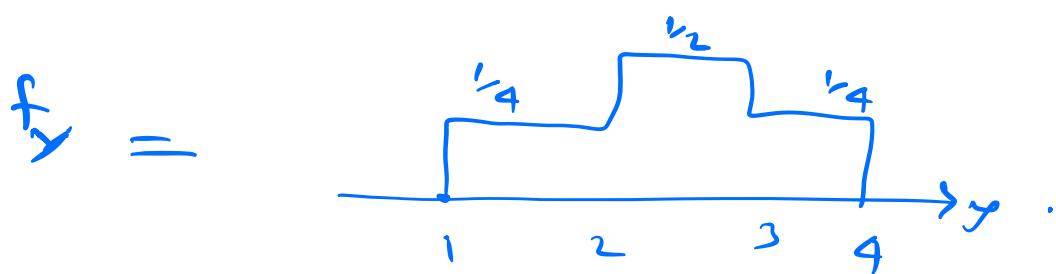
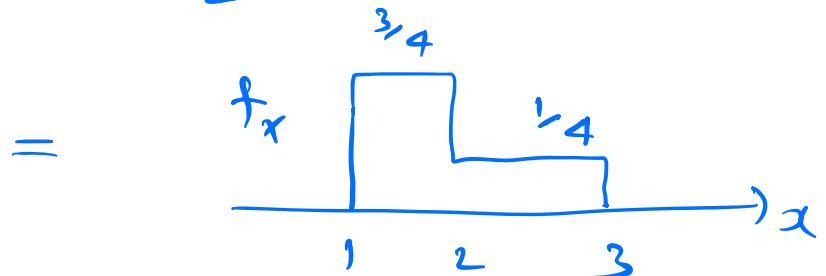
$$\text{Area}(S) = 4$$

$$f_{xy}(x,y) = \begin{cases} \frac{1}{4} & \text{if } (x,y) \in S \\ 0 & \text{o.w.} \end{cases}$$

Find f_x , f_y .

$$f_x(x) = \int_{-\infty}^{\infty} f_{xy}(x,y) dy$$

$$= \begin{cases} \frac{1}{4} dy & - x \in [1, 2] \\ \frac{1}{4} dy & - x \in [2, 3] \end{cases}$$



Joint CDFs. If x and y are two RVs associated with the same experiment we define their joint CDF by

$$F_{xy}(x, y) = P(X \leq x, Y \leq y).$$

For continuous RVs x and y

$$F_{xy}(x, y) = \int_{-\infty}^x \int_{-\infty}^y f_{xy}(s, t) dt ds.$$

$$f_{xy}(x, y) = \frac{\partial F_{xy}(x, y)}{\partial x \partial y}.$$

$$P(x < x \leq x_2, y < y \leq y_2)$$

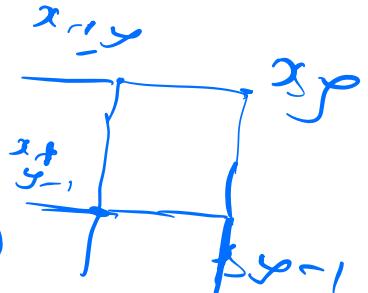
$$= F_{xy}(x_2, y_2) - F_{xy}(x_1, y_2) - F_{xy}(x_2, y_1) + F_{xy}(x_1, y_1)$$

For discrete RVs x and y

$$F_{xy}(x, y) = \sum_{l \leq x} \sum_{k \leq y} P_{xy}(l, k).$$

$$P_{xy}(x, y) = F_{xy}(x, y) - F_{xy}(x-1, y)$$

$$- F_{xy}(x, y-1) + F_{xy}(x-1, y-1)$$



Example. Let x and y be described by a uniform PDF on the unit square,

$$F_{xy}(x, y) = P(x \leq x, y \leq y) = xy \text{ for } 0 \leq x, y \leq 1,$$

$$F_{xy}(x, y) = x \quad \text{if } 0 \leq x \leq 1 \text{ and } y \geq 1.$$

$$f_{xy}(x, y) = \frac{\partial F_{xy}(x, y)}{\partial x \partial y} = 1 \quad \text{for } 0 \leq x, y \leq 1,$$

Properties of Joint CDF

$$(1) \lim_{x \rightarrow \infty} F_{xy}(x, y) = F_y(y),$$

$$\lim_{y \rightarrow \infty} F_{xy}(x, y) = F_x(x),$$

$$\lim_{\substack{x \rightarrow \infty \\ y \rightarrow \infty}} F_{xy}(x, y) = 1.$$

$$(2) \lim_{x \rightarrow -\infty} F_{xy}(x, y) = 0 = \lim_{y \rightarrow -\infty} F_{xy}(x, y).$$

(3) If $x_1 \leq x_2, y_1 \leq y_2$ then

$$F_{xy}(x_1, y_1) \leq F_{xy}(x_2, y_2),$$

$$(4) \lim_{\substack{\varepsilon \rightarrow 0^+ \\ \delta \rightarrow 0^+}} F_{xy}(x + \varepsilon, y + \delta) = F_{xy}(x, y),$$

Expected value Rule

If x and y are jointly continuous RVs and g is some function then $z = g(x, y)$ is

also a RV,

$$E[g(x,y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y) f_{x,y}(x,y) dx dy.$$

$$E[\alpha x + by + c] = \alpha E[x] + bE[y] + c.$$

More than Two Random Variables

x, y and z are jointly continuous if

\exists non-negative f_{xyz} s.t.

$$P((x,y,z) \in B) = \int_{(x,y,z) \in B} f_{xyz}(x,y,z) dx dy dz$$

for any set $B \subseteq \mathbb{R}^3$.

$\Rightarrow x$ and y are jointly continuous with joint PDF

$$f_{xy}(x,y) = \int_{-\infty}^{\infty} f_{xyz}(x,y,z) dz.$$

$$f_x(x) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{xyz}(x,y,z) dy dz.$$

$$E\left[\sum_{i=1}^n a_i x_i\right] = \sum_{i=1}^n a_i E[x_i].$$

Conditioning a RV on an Event

The conditional PDF of a continuous RV x given an event A with $P(A) > 0$, is defined as a non-negative function $f_{x|A}$ that satisfies

$$P(x \in B | A) = \int_B f_{x|A}(x) dx,$$

for any $B \subseteq \mathbb{R}$ and $A \subseteq \Omega$,

$$B = \mathbb{R} \Rightarrow \int_{-\infty}^{\infty} f_{x|A}(x) dx = 1.$$

Suppose A is of the form $A = \{x \in c\}$.

$$\begin{aligned} P(x \in B | x \in c) &= \frac{P(x \in B \cap c)}{P(x \in c)} \\ &= \frac{1}{P(x \in c)} \int_{B \cap c} f_x(x) dx \end{aligned}$$

$$= \int_B \frac{f_x(x) \mathbf{1}_{\{x \in c\}}}{P(x \in c)} dx$$

$$\text{so } f_{x|x \in C}(x) = \begin{cases} \frac{f_x(x)}{P(x \in C)}, & \text{if } x \in C \\ 0, & \text{o.w.} \end{cases}$$

Conditional CDF

$$F_{x|A}(x) = P(x \leq x | A).$$

Theorem. Let A_1, A_2, \dots, A_n be disjoint sets that form a partition of the sample space, and assume that $P(A_i) > 0$ for all i . Then

$$f_x(x) = \sum_{i=1}^n P(A_i) f_{x|A_i}(x).$$

$$\begin{aligned} \text{Proof. } F_x(x) &= P(x \leq x) \\ &= \sum_{i=1}^n P(A_i) F_{x|A_i}(x) \end{aligned}$$

Differentiating w.r.t. x , we get

$$f_x(x) = \sum_{i=1}^n P(A_i) f_{x|A_i}(x).$$

Lecture 16

(9 September 2023)

Conditioning one RV on another

Consider a f_{xy} . For any y with $f_y(y) > 0$, the conditional PDF of x given $y=y_-$ is defined by

$$f_{x|y}(x|y) = \frac{f_{xy}(x,y)}{f_y(y)},$$

$$\int_{-\infty}^{\infty} f_{x|y}(x|y) dx = 1$$

To interpret the conditional probability
consider

$$\begin{aligned} & P(x \leq x \leq x+\delta_1, | y \leq y \leq y+\delta_2) \\ &= P(x \leq x \leq x+\delta_1, y \leq y \leq y+\delta_2) / P(y \leq y \leq y+\delta_2) \\ &\approx \frac{f_{xy}(x,y) \delta_1 \delta_2}{f_y(y) \delta y}. \end{aligned}$$

In view of the above

$$\text{Define } P(X \in A | Y = y) = \int_A f_{X|Y}(x|y) dx.$$

Interpretation:

$$P(X \in A | Y = y) := \lim_{\delta \rightarrow 0} P(X \in A | y \leq Y \leq y + \delta)$$

$$= \lim_{\delta \rightarrow 0} \frac{P(X \in A | y \leq Y \leq y + \delta) / \delta}{P(y \leq Y \leq y + \delta) / \delta}$$

$$= \lim_{\delta \rightarrow 0} \frac{\int \int_{\substack{x \in A \\ y}}^{y+\delta} f_{X|Y}(x|t) dt dx}{\delta} \frac{\int_y^{y+\delta} f_Y(y) dy}{\delta}$$

$$= \int_{x \in A} f_X(x) \left[\lim_{\delta \rightarrow 0} \int_y^{y+\delta} f_{Y|X}(t|x) dt \right] \frac{f_Y(y)}{\delta} dy$$

$$= \int_{x \in A} f_X(x) \frac{d}{dy} F_{Y|X}(y|x) \Big|_{y} \frac{f_Y(y)}{\delta} dy$$

$$= \int_{x \in A} f_X(x) \frac{f_{Y|X}(y|x)}{f_Y(y)} dy$$

Conditional Expectation

$$E[x|A] = \int_{-\infty}^{\infty} x f_{x|A}(x) dx.$$

$$E[x|y=y] = \int_{-\infty}^{\infty} x f_{x|y}(x|y) dx.$$

Expected Value Rule:

$$E[g(x)|A] = \int_{-\infty}^{\infty} g(x) f_{x|A}(x) dx.$$

$$E[g(x)|y=y] = \int_{-\infty}^{\infty} g(x) f_{x|y}(x|y) dx,$$

Total Expectation Theorems

(i) Let A_1, A_2, \dots, A_n be disjoint events that form a partition of the sample space, and let $P(A_i) > 0$ for all i . Then

$$E[x] = \sum_{i=1}^n P(A_i) E[x|A_i].$$

Proof. $f_x(x) = \sum_{i=1}^n P(A_i) f_{x|A_i}(x)$

multiply both sides by x and then integrate to

$$\text{Get } E[x] = \sum_{i=1}^n P(A_i) E[x|A_i].$$

$$(ii) \int_{-\infty}^{\infty} E[x|r=y] f_y(y) dy = E[x].$$

$$\text{Proof. } \int_{-\infty}^{\infty} E[x|r=y] f_y(y) dy$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f_{x|y}(x|y) f_y(y) dy dx$$

$$= \int_{-\infty}^{\infty} x f_x(x) dx = E[x].$$

The total expectation theorem can often be used to calculate mean and variance.

Example, mean and variance of a piecewise constant PDF.

$$f_x(x) = \begin{cases} \frac{1}{3}, & \text{if } 0 \leq x \leq 1 \\ \frac{2}{3}, & \text{if } 1 \leq x \leq 2 \\ 0, & \text{otherwise,} \end{cases}$$

$$A_1 = \{x \text{ lies in } [0, 1]\} \quad A_2 = \{x \text{ lies in } (1, 2]\}.$$

$$P(A_1) = \frac{1}{3}, \quad P(A_2) = \frac{2}{3}.$$

$$E[x] = P(A_1)E[x|A_1] + P(A_2)E[x|A_2]$$

$f_{x|A_1}$, $f_{x|A_2}$ are uniform.

Recall a UNUniform $[a, b] \Rightarrow E[u] = \frac{a+b}{2}$,

$$E[u^2] = \frac{a^2+ab+bc^2}{3}$$

$$E[x|A_1] = \frac{1}{2}, \quad E[x|A_2] = \frac{3}{2}$$

$$E[x^2|A_1] = \frac{1}{3}, \quad E[x^2|A_2] = \frac{2}{3}.$$

$$E[x] = \frac{7}{6}, \quad E[x^2] = \frac{15}{9}.$$

$$\text{Var}(x) = E[x^2] - E[x]^2 = \frac{11}{36}.$$

Independence

Two continuous RVS are independent if

$$f_{xy}(x,y) = f_x(x)f_y(y)$$

This is same as

$$f_{x|y}(x|y) = f_x(x) \text{ for all } y \text{ with } f_y(y) > 0 \text{ and for all } x,$$

Example. Independent Normal RVs.

Let x and y be independent Gaussian RVs with means $\mu_x - \mu_y$, and variances $\sigma_x^2 - \sigma_y^2$ respectively,

$$f_{xy}(x,y) = \frac{1}{2\pi\sigma_x\sigma_y} e^{-\frac{(x-\mu_x)^2}{2\sigma_x^2} - \frac{(y-\mu_y)^2}{2\sigma_y^2}}.$$

Theorem. Two random variables x and y (continuous or discrete) are independent if and only if

$$F_{xy}(x,y) = F_x(x)F_y(y), \text{ for all } x,y.$$

Proof. We prove for the continuous case.

Let x and y are independent.

$$F_{xy}(x,y) = P(x \leq x, y \leq y)$$

$$= \int_{-\infty}^x \int_{-\infty}^y f_{xy}(x,y) dx dy$$

$$= \int_{-\infty}^x f_x(x) dx \int_{-\infty}^y f_y(y) dy = F_x(x)F_y(y).$$

For the other direction let

$$F_{xy}(x,y) = F_x(x)F_y(y).$$

On taking second order mixed partial derivative

$$\frac{\partial^2 F_{xy}(x,y)}{\partial x \partial y} = \frac{\partial F_x(x)}{\partial x} \cdot \frac{\partial F_y(y)}{\partial y}$$

$$\Rightarrow f_{xy}(x,y) = f_x(x)f_y(y).$$

Exercise. Prove the above theorem for the discrete case.

- If x and y are independent, then

$$E[xy] = E[x]E[y].$$

Theorem. If x and y are jointly continuous independent RVS. For any two functions g & h ,

$$E[g(x)h(y)] = E[g(x)]E[h(y)].$$

Proof. It suffices to show that $g(x)$ &

$h(y)$ are independent.

Let $x' = g(x)$, $y' = h(y)$.

$$F_{x', y'}(x', y') = P(g(x) \leq x', h(y) \leq y')$$

$$= \int_{\{x : g(x) \leq x'\}} \int_{\{y : h(y) \leq y'\}} f_{x,y}(x,y) dx dy$$

$$= \int_{x : g(x) \leq x'} f_x(x) dx \cdot \int_{y : h(y) \leq y'} f_y(y) dy$$

$$= P(g(x) \leq x') P(h(y) \leq y')$$

$$= F_{x'}(x') F_{y'}(y').$$

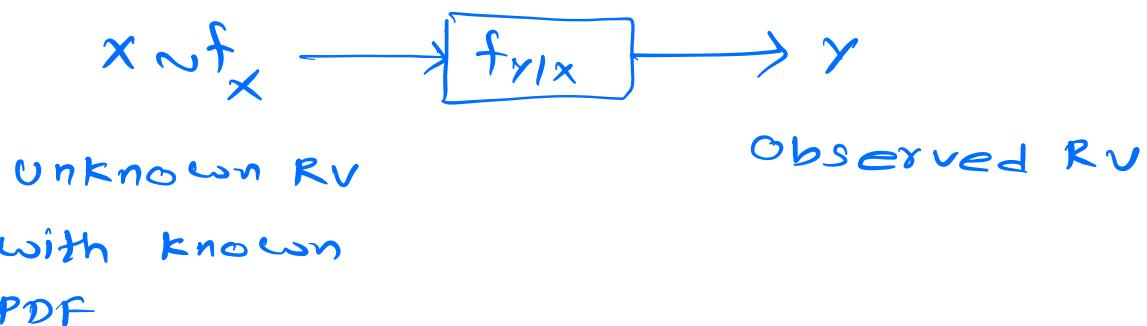
- If x and y are independent, then

$$\text{Var}(x+y) = \text{Var}(x) + \text{Var}(y).$$

Lecture 17

(12 September 2023)

The Continuous Bayes' Rule



Goal: Infer about x ?

The information provided by the event $\{Y=y\}$ is captured by the conditional PDF $f_{x|y}(x|y)$.

$$f_{x,y}(x,y) = f_x(x) f_{y|x}(y|x) = f_y(y) f_{x|y}(x|y)$$

$$\Rightarrow f_{x|y}(x|y) = \frac{f_x(x) f_{y|x}(y|x)}{\int_{-\infty}^{\infty} f_x(t) f_{y|x}(y|t) dt}.$$

Inference about a Discrete RV:

$$P(A|Y=y) = \lim_{\delta \rightarrow 0} P(A | Y \in [y, y+\delta])$$

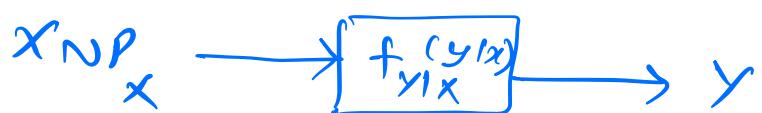
$$= \lim_{\delta \rightarrow 0} \frac{P(A) P(y \in [y, y+\delta] | A) / \delta}{P(y \in [y, y+\delta]) / \delta}$$

$$= P(A) \frac{f_{y|A}(y)}{f_y(y)}.$$

Also $f_y(y) = P(A) f_{y|A}(y) + P(A^c) f_{y|A^c}(y)$,
so

$$P(A | y=y) = \frac{P(A) f_{y|A}(y)}{P(A) f_{y|A}(y) + P(A^c) f_{y|A^c}(y)}$$

Let x be a discrete RV and y be a continuous RV. The above equality gives



$$P_{x|y}(x|y) = \frac{P_x(x) f_{y|x}(y|x)}{f_y(y)}$$

$$= \frac{P_x(x) f_{y|x}(y|x)}{\sum_{x=0}^{M-1} P_x(x) f_{y|x}(y|x)}$$

, where
 $\sum_{x=0}^{M-1} P_x(x) = 1,$

Goal: To estimate the hypothesis x that lead to an observation y .

- A test $\hat{x}(y)$ is a decision rule or
- deterministic function of the observation y ,
- $P_{x|y}(x|y)$ is the probability that hypothesis x is correct, i.e., the probability that $X=x$, conditional on observation y .
 $P_{x|y}(x|y)$ is called 'a posteriori probability'.

Consider the decision rule that maximizes this a posteriori probability.

$$\hat{x}_{MAP}(y) = \arg \max_x P_{x|y}(x|y) \text{ (MAP rule)},$$

where $\arg \max_x$ means the argument $x \in \{0, \dots, M-1\}$ that maximizes the function,
Using the Bayes' law

$$\hat{x}_{\text{MAP}}(y) = \arg \max_x \frac{P_x(x) f_{y|x}(y|x)}{f_y(y)}$$

$$= \arg \max_x P_x(x) f_{y|x}(y|x).$$

when multiple hypotheses achieve the maximum, we arbitrarily choose the largest maximizing x .

- For any test A , $P_{x|y}(\hat{x}_A(y)|y)$ is the probability that $\hat{x}_A(y)$ is the correct decision when test A is used on observation y . Since $\hat{x}_{\text{MAP}}(y)$ maximizes the probability of correct decision we have

$$P_{x|y}(\hat{x}_{\text{MAP}}(y)|y) \geq P_{x|y}(\hat{x}_A(y)|y)$$

for all A and y .

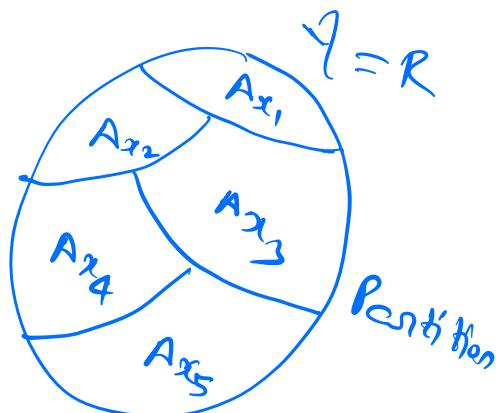
- Probability of correctness for a test A :

$$P(\hat{x}_A(y)=x).$$

Theorem. The MAP rule maximizes the probability of correct decision conditional on each observed sample y . It also maximizes the overall probability of correct decision defined above.

Proof. $A_x = \{y : \hat{x}_A(y) = x\}$ be the set of observations y that test A maps to hypothesis x .

$$\begin{aligned}
 P(\hat{x}_A(y) = x) &= \sum_{x=0}^{M-1} P_x(x) P(\hat{x}_A(y) = x | x=x) \\
 &= \sum_{x=0}^{M-1} P_x(x) P(y \in A_x | x=x) \\
 &= \sum_{x=0}^{M-1} P_x(x) \int_{y \in A_x} f_{y|x}(y|x) dy \\
 &= \int_{y=-\infty}^{\infty} P_x(\hat{x}_A(y)) f_{y|x}(y|\hat{x}_A(y)) dy \\
 &= \int_{y=-\infty}^{\infty} P_{x,y}(\hat{x}_A(y)|y) f_y(y) dy \rightarrow ①
 \end{aligned}$$



Similarly

$$P(\hat{x}_{MAP}(y) = x) = \int_{y=-\infty}^{\infty} P_{x|y}(\hat{x}_{MAP}(y)|y) f_y(y) dy \rightarrow ②$$

Now since

$$P_{x|y}(\hat{x}_{MAP}(y)|y) \geq P_{x|y}(\hat{x}_A(y)|y) \text{ for all } A, y$$

averaging gives

$$② \geq ①, \text{ i.e., } P(\hat{x}_{MAP}(y) = x) \geq P(\hat{x}_A(y) = x).$$

(Revisiting) Binary MAP detection.

Let $P_x(0) = p_0$, $P_x(1) = p_1$, $0 < p_0, p_1$, $p_1 + p_0 = 1$.

Let y be a continuous RV with conditional PDF $f_{y|x}(y|x)$.

$$f_y(y) = p_0 f_{y|x}(y|0) + p_1 f_{y|x}(y|1) > 0.$$

MAP rule for a fixed y ,

$$P_{x|y}(1|y) \geq P_{x|y}(0|y),$$

$\hat{x}(y) = 1$
 $\hat{x}(y) = 0$

Equivalently using Bayes' rule,

$$\frac{f_{Y|X}(y|1)P_X(1)}{f_Y(y)} \stackrel{\hat{x}(y)=1}{\geq} \stackrel{\hat{x}(y)=0}{<} \frac{f_{Y|X}(y|0)P_X(0)}{f_Y(y)},$$

Equivalently,

$$\text{Likelihood ratio} = \frac{f_{Y|X}(y|1)}{f_{Y|X}(y|0)} \stackrel{\hat{x}(y)=1}{\geq} \frac{P_0}{P_1} = \eta.$$

[Threshold rule]

- Note that if the α prior probability P_0 is increased, then the threshold increases and the set of y for which hypotheses o is chosen increases; this corresponds to our intuition - the greater our initial conviction that x is o , the stronger the evidence required to change our minds.

- A special case of the above threshold rule is when $P_0 = P_1$. It is called Maximum Likelihood (ML) test. The ML test is often used when P_0 and P_1 are unknown,

The position of the space of observed sample values is given by

$$A_1 = \{y : \lambda(y) \geq n\} \quad A_0 = \{y : \lambda(y) < n\}$$

\downarrow \downarrow
 $\hat{x} = 1$ $\hat{x} = 0$.

We can compute the probability of error $P(\text{Error} | x=x)$.

$$P(\text{Error} | x=0) = P(\hat{x}(y)=1 | x=0)$$

$$= P(y \in A_1 | x=0)$$

$$= \int_{y \in A_1} f_{x|X}(y|0) dy.$$

Similarly $P(\text{Error} | x=1) = \int_{y \in A_0} f_{x|X}(y|1) dy$.

The overall probability of error

$$= P_0 P(\text{Error} | x=0) + P_1 P(\text{Error} | x=1).$$

Example. (Abstraction of a digital communication system)

$$P_x(b) = P_1, \quad P_x(-b) = P_0. \quad Y = X + Z, \quad Z \sim \mathcal{N}(0, \sigma^2),$$

x and z are independent.

$$P(Y \in B | x = x) = P(x + z \in B | x = x)$$
$$= P(z \in B)$$

$\therefore Y|x=b \sim N(b, \sigma^2) \quad Y|x=-b \sim N(-b, \sigma^2),$

$$f_{Y|X}(y|b) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-b)^2}{2\sigma^2}}$$

$$f_{Y|X}(y|-b) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y+b)^2}{2\sigma^2}}.$$

MAP rule:

$$\Lambda(y) = \frac{f_{Y|X}(y|b)}{f_{Y|X}(y|-b)} = e^{\frac{(y+b)^2 - (y-b)^2}{2\sigma^2}}$$
$$= e^{2yb/\sigma^2}$$
$$\hat{x}(y) = b \geq \frac{p_0}{p_1} = n.$$
$$\hat{x}(y) = -b < \frac{p_0}{p_1} = n.$$

$$\Rightarrow y \geq \frac{\sigma^2}{2b} \log n.$$

$$\hat{x}(y) = -b$$

If $b = 1$ and $\rho_0 = \rho_1$, this recovers the example on signal detection we have seen in an earlier lecture. Then

ML rule maps $y \geq 0 \rightarrow \hat{x} = b$
 $y < 0 \rightarrow \hat{x} = -b.$

Coming back to the case with $b > 0$

$$\begin{aligned}
 P(Er|x=b) &= P(y < 0 | x=b) \\
 &= P\left(\frac{y-b}{\sigma} < \frac{-b}{\sigma} \mid x=b\right) \\
 &= P(N < \frac{-b}{\sigma}) \\
 &= P(N \geq \frac{b}{\sigma}).
 \end{aligned}$$

$$\begin{aligned}
 P(Er|x=-b) &= P(N > \frac{b}{\sigma}) \\
 &= 1 - P(N \leq \frac{b}{\sigma}).
 \end{aligned}$$

$$P(\text{Error}) = P(N \geq \frac{b}{\sigma}).$$

Lecture 18
(16 October 2023)

Functions of Random Variables

$y = g(x)$, i.e., $y(\omega) = g(x(\omega))$.

For the discrete case

$$P_y(y) = \sum_{x: g(x)=y} P_x(x).$$

For the continuous case, we follow the two-step procedure outlined below.

1) Calculate the CDF F_y of y using

$$F_y(y) = P(g(x) \leq y) = \int_{\{x: g(x) \leq y\}} f_x(x) dx.$$

2) Differentiate the CDF to obtain the PDF of y :

$$f_y(y) = \frac{d}{dy} F_y(y).$$

Linear Function of a RV:

$$Y = ax + b, \quad a \neq 0.$$

Discrete case -

$$P_Y(y) = P_X\left(\frac{y-b}{a}\right) \text{ for all } y,$$

continuous case -

$$Y = ax + b, \quad a \neq 0$$

$$\Rightarrow f_Y(y) = \frac{1}{|a|} f_X\left(\frac{y-b}{a}\right).$$

Proof. $a > 0$ (proved earlier)

Let $a > 0$,

$$\begin{aligned} P(ax+b \leq y) &= P(x \geq \frac{y-b}{a}) \\ &= 1 - P(x < \frac{y-b}{a}) \\ &= 1 - F_X\left(\frac{y-b}{a}\right) \\ \Rightarrow f_Y(y) &= \frac{d}{dy} F_Y(y) = \frac{1}{|a|} f_X\left(\frac{y-b}{a}\right). \end{aligned}$$

Example. X is uniform on $[0, 1]$ and let $Y = \sqrt{x}$.

For $0 \leq y \leq 1$,

$$F_Y(y) = P(\sqrt{x} \leq y)$$

$$= P(x \leq y^2)$$

$$= y^2$$

$$\Rightarrow f_Y(y) = 2y, \quad 0 \leq y \leq 1.$$

$$F_Y(y) = 0 \quad y \leq 0 \quad \text{and} \quad F_Y(y) = 1 \quad y \geq 1.$$

$\therefore f_Y(y) = 0$ for y outside $[0, 1]$.

Example. $y = x^3$ - x uniform $[0, 2]$.

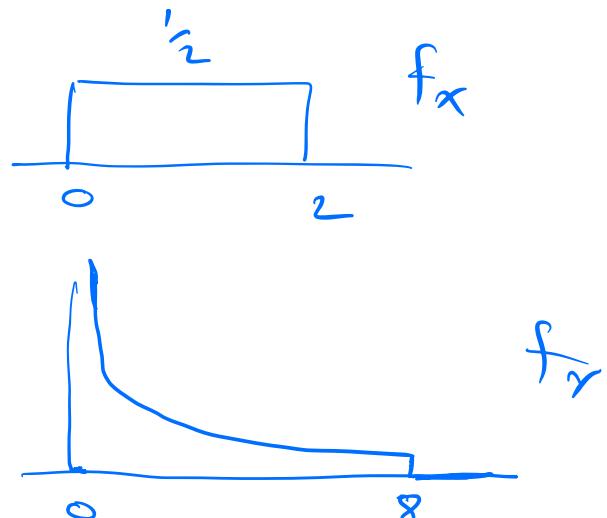
For $0 \leq y \leq 8$

$$F_Y(y) = P(x^3 \leq y)$$

$$= P(x \leq y^{1/3})$$

$$= \frac{y^{1/3}}{2}$$

$$\Rightarrow f_Y(y) = \frac{1}{6} y^{-2/3} \quad 0 \leq y \leq 8.$$



Example. $y = \frac{10}{x}$ - $f_x(x) = \frac{1}{5}$, $5 \leq x \leq 10$.

For $1 \leq y \leq 2$

$$F_Y(y) = P\left(\frac{10}{x} \leq y\right) = P\left(x \geq \frac{10}{y}\right)$$

$$= 1 - P(X \leq \frac{10}{y})$$

$$= 1 - \frac{1}{5}(-5 + \frac{10}{y}) = -\frac{2}{y} + 2$$

$$\Rightarrow f_y(y) = \frac{-2}{y^2}, \quad 1 \leq y \leq 2.$$

Monotonic Function:

Let x be a continuous RV and suppose that its range is contained in a certain interval I , in the sense that $f_x(x)=0$ for $x \notin I$. We consider $y=g(x)$ and assume that g is strictly monotonic over the interval I :

$g(x) < g(x')$ for all $x, x' \in I$ s.t. $x < x'$

(monotonically increasing)
or

$g(x) > g(x')$ for all $x, x' \in I$ s.t. $x < x'$.

For $y \in g^{-1}(I)$, (monotonically decreasing)

$P(g(x) \leq y)$

incr.g

$$= P(X \leq g^{-1}(y)) = F_x(g^{-1}(y))$$

$$\Rightarrow f_y(y) = f_x(g^{-1}(y)) (g^{-1})'(y).$$

For decreasing g ,

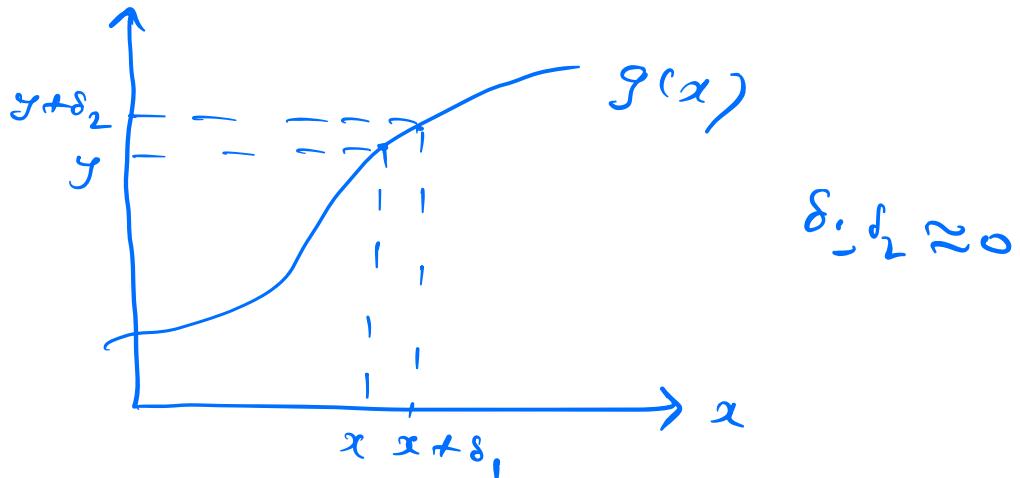
$$\begin{aligned}F_Y(y) &= P(g(x) \leq y) = P(x \geq g^{-1}(y)) \\&= 1 - P(x \leq g^{-1}(y)) = 1 - F_X(g^{-1}(y)) \\&\Rightarrow f_Y(y) = f_X(g^{-1}(y)) |(g^{-1})'(y)|,\end{aligned}$$

$$\therefore f_Y(y) = f_X(g^{-1}(y)) |(g^{-1})'(y)|.$$

This can also be written as

$$f_Y(y) = \begin{cases} f_X(g^{-1}(y)) / |g'(g^{-1}(y))| & \text{if } g(x)=y \text{ for some } x \\ 0 & \text{if } g(x) \neq y \text{ for all } x \text{ with } f_X(x)>0 \end{cases}$$

Illustration:



$$P(y \leq Y \leq y + \delta_2) \approx f_Y(y) \delta_2$$

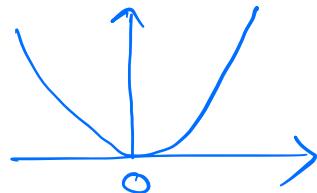
$$= P(x \leq X \leq x + \delta_1) = f_X(x) \delta_1$$

$$f_y(y) = f_x(x) \frac{\delta_1}{\delta_2} = \frac{f_x(x)}{g'(x)} = \frac{f_x(g^{-1}(y))}{g'(g^{-1}(y))}.$$

Non-monotonic Function

$$Y = g(x), \quad g(x) = x^2$$

For $y \geq 0$



$$\begin{aligned} P(g(x) \leq y) &= P(x^2 \leq y) \\ &= P(-\sqrt{y} \leq x \leq \sqrt{y}) \\ &= F_x(\sqrt{y}) - F_x(-\sqrt{y}) \end{aligned}$$

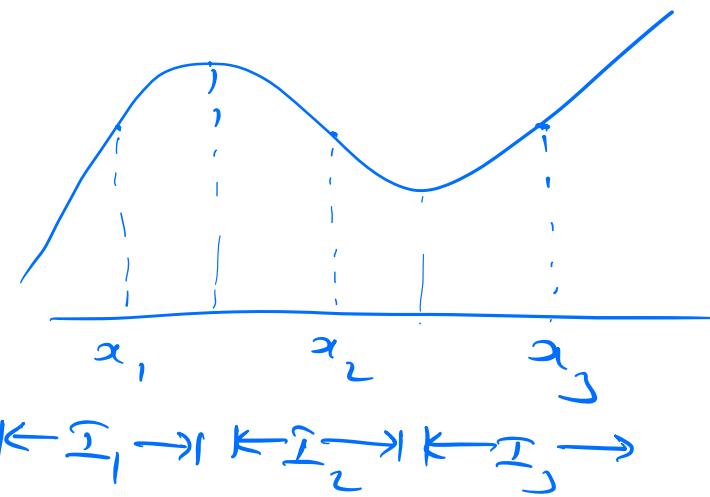
$$\Rightarrow f_y(y) = \frac{f_x(\sqrt{y})}{2\sqrt{y}} + \frac{f_x(-\sqrt{y})}{2\sqrt{y}}.$$

Theorem. Consider a continuous RV x with PDF f_x , and let $y = g(x)$, suppose we partition the domain of f_x into a finite number of intervals such that $g(x)$ is strictly monotone and differentiable on each interval. Then the PDF of y is given by

$$f_y(y) = \sum_{i=1}^n \frac{f_x(x_i)}{|g'(x_i)|} - \text{where}$$

x_1, x_2, \dots, x_n are real solutions to $g(x)=y$ in the domain of f_x .

Proof.



$$\begin{aligned}I_1 &= [l_1, u_1], \\I_2 &= [l_2, u_2], \\I_3 &= [l_3, u_3].\end{aligned}$$

For a given y let $g(x_i) = y$ $i \in \{1, 2, 3\}$.
 $g(x)$ is increasing at x_1 ,

decreasing at x_2 and
increasing at x_3 .

$$\begin{aligned}P(g(x) \leq y) &= P(\{l_1 \leq x \leq x_1\} \cup \{x_2 \leq x \leq u_2\} \\&\quad \cup \{l_3 \leq x \leq x_3\}) \\&= P(l_1 \leq x \leq x_1) + P(x_2 \leq x \leq u_2) + P(l_3 \leq x \leq x_3) \\&\Rightarrow f_y(y) = f_x(x_1) \frac{dx_1}{dy} - f_x(x_2) \frac{dx_2}{dy} + f_x(x_3) \frac{dx_3}{dy}\end{aligned}$$

$$= \frac{f_x(x_1)}{|g'(x_1)|} + \frac{f_x(x_2)}{|g'(x_2)|} + \frac{f_x(x_3)}{|g'(x_3)|}.$$

\rightarrow for $y = x^2$, $x_1 = \sqrt{y}$, $x_2 = -\sqrt{y}$.

$$\text{so } f_y(y) = \frac{f_x(\sqrt{y})}{2\sqrt{y}} + \frac{f_x(-\sqrt{y})}{2\sqrt{y}}.$$

Lecture 19
(26 October 2023)

Functions of Two Random Variables

$$Z = g(X, Y),$$

$$\text{i.e., } Z(\omega) = g(X(\omega), Y(\omega)),$$

Example. $Z = \max(X, Y)$, X & Y are independent.

$$P(Z \leq z) = P(\max(X, Y) \leq z)$$

$$= P(X \leq z, Y \leq z)$$

$$= P(X \leq z) P(Y \leq z)$$

$$= F_X(z) F_Y(z)$$

$$\Rightarrow f_Z(z) = F'_X(z) f_Y(z) + f'_X(z) F_Y(z).$$

As a special case $X, Y \sim \text{Uniform}[0, 1]$,

$$f_Z(z) = \begin{cases} 2z, & 0 \leq z \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

Exercise: $Z = \min\{X, Y\}$ X and Y are indep. Find f_Z .

Example. X and Y are independent RVS that are uniformly distributed on $[0, 1]$. What is the PDF of $Z = Y/X$?

$$P(Z \leq z) = P(Y/X \leq z)$$

For $0 \leq z \leq 1$

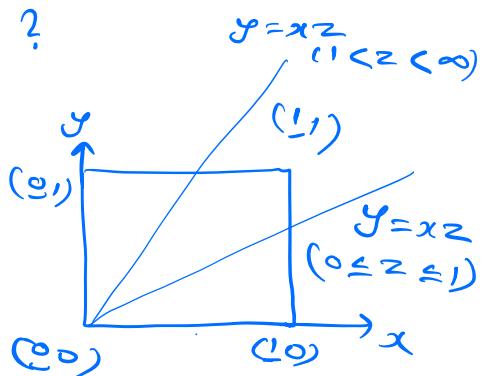
$$\begin{aligned} P(Y \leq xz) &= \int_{x=0}^{zx} \int_{y=0}^1 1 dx dy = z/2. \end{aligned}$$

For $1 < z < \infty$

$$P(Y \leq xz) = 1 - \frac{1}{2z}.$$

$$\Rightarrow F_Z(z) = \begin{cases} z/2 & \text{if } 0 \leq z \leq 1 \\ 1 - \frac{1}{2z} & \text{if } z > 1 \\ 0 & \text{otherwise.} \end{cases}$$

$$f_Z(z) = \begin{cases} \frac{1}{2} & \text{if } 0 \leq z \leq 1 \\ \frac{1}{2z^2} & \text{if } z > 1 \\ 0 & \text{otherwise.} \end{cases}$$



Sum of Independent Random Variables

Let $Z = X+Y$ where X and Y are independent random variables,

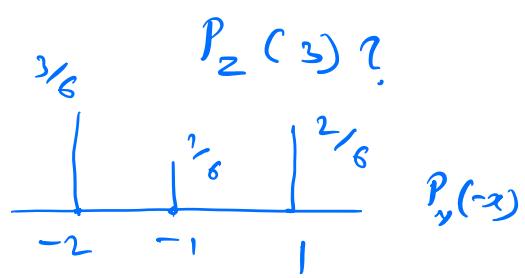
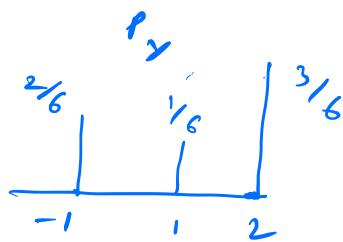
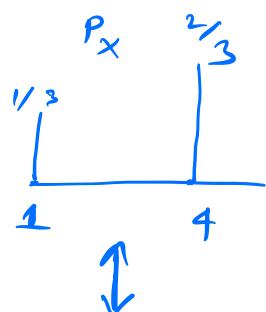
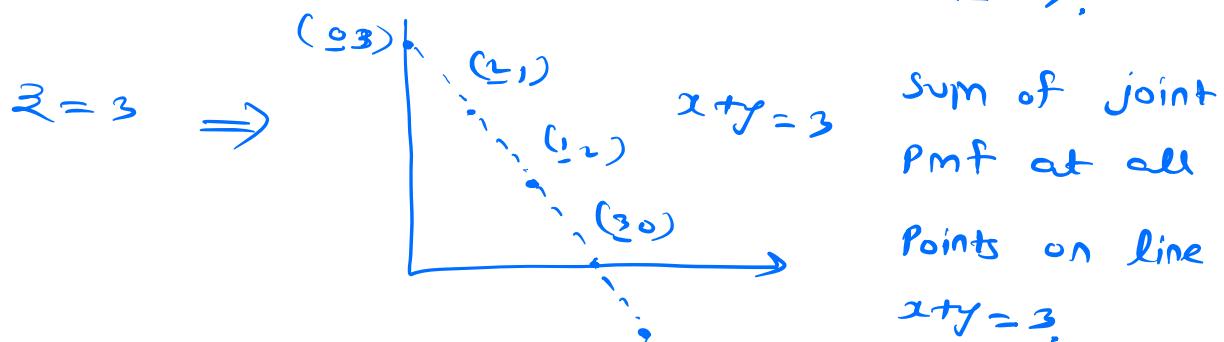
Discrete case:

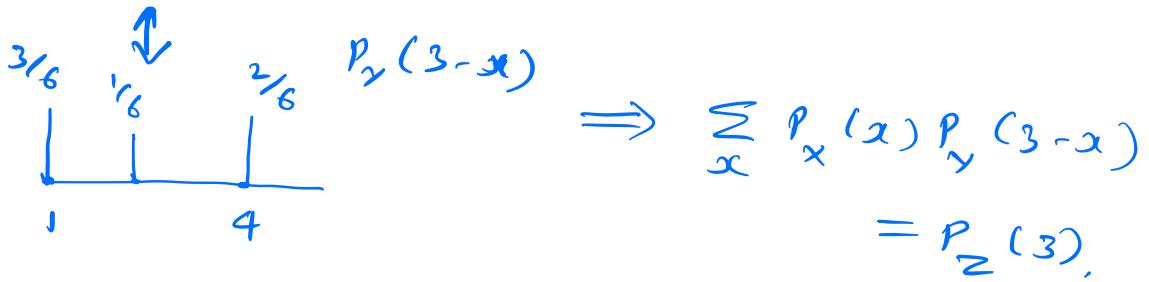
$$P_Z(z) = P(X+Y=z)$$

$$= \sum_{(x,y): x+y=z} P_{XY}(x,y)$$

$$= \sum_x P_{XY}(x, z-x) = \sum_x P_X(x) P_Y(z-x).$$

The resulting PMF P_Z is called the convolution of the PMFs of X and Y .





Continuous case:

$$P(X+Y \leq z) = \int_{x=-\infty}^{\infty} \int_{y=-\infty}^{z-x} f_{xy}(x,y) dy dx$$

$$= \int_{x=-\infty}^{\infty} \int_{y=-\infty}^{z-x} f_x(x) f_y(y) dy dx$$

$$= \int_{x=-\infty}^{\infty} f_x(x) F_y(z-x) dx$$

$$\Rightarrow f_z(z) = \frac{d}{dz} F_z(z)$$

$$= \int_{-\infty}^{\infty} f_x(x) \frac{d}{dz} F_y(z-x) dx$$

$$= \int_{x=-\infty}^{\infty} f_x(x) f_y(z-x) dx.$$

Theorem. If $X \sim N(\mu_1, \sigma_1^2)$, $Y \sim N(\mu_2, \sigma_2^2)$

are independent then

$$X+Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2).$$

Proof. $Z = X+Y$,

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{(z-x)^2}{2}} dx$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{-x^2 - z^2 - x^2 + 2zx}{2}} dx$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2 - z^2 + 2zx}{2}} dx$$

$$= \frac{1}{\sqrt{4\pi}} \cdot e^{-\frac{z^2}{4}} \cdot \int_{-\infty}^{\infty} \frac{1}{\sqrt{\pi}} e^{-\frac{(x-z)^2}{2}} dx$$

$$= \frac{1}{\sqrt{4\pi}} \cdot e^{-\frac{z^2}{4}} \quad \underbrace{\qquad}_{=} = 1$$

Exercise. $X \sim N(\mu_x, \sigma_x^2)$, $Y \sim N(\mu_y, \sigma_y^2)$ indep $\Rightarrow X+Y \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$.

Sum of a discrete RV and a continuous RV:

$$Z = X + Y$$

X-discrete, Y-continuous, X & Y are independent

$$P(Z \leq z) = P(X+Y \leq z)$$

$$= \sum_x P(X+Y \leq z | X=x) P_X(x)$$

$$= \sum_x P(Y \leq z-x | X=x) P_X(x)$$

$$= \sum_x F_Y(z-x) P_X(x)$$

$$= \sum_x \int_{-\infty}^{z-x} f_Y(y) dy P_X(x) \quad y = t-x$$

$$= \sum_x \int_{-\infty}^z f_Y(t-x) dt P_X(x)$$

$$= \int_{-\infty}^z \left(\underbrace{\sum_x f_Y(t-x) P_X(x)}_{f_Z(t)} \right) dt$$

$$\therefore f_Z(z) = \sum_x f_Y(z-x) P_X(x).$$

Functions of Two Random Variables

$(X, Y) \sim f_{XY}$.

$Z = g_1(X, Y)$, $W = g_2(X, Y)$.

$$P(Z \leq z, W \leq w) = P(g_1(X, Y) \leq z, g_2(X, Y) \leq w)$$

$$= \iint_{\substack{xy \\ g_1(x,y) \leq z \\ g_2(x,y) \leq w}} f_{XY}(x, y) dx dy$$

$$g_1(x, y) \leq z$$

$$g_2(x, y) \leq w$$

Sometimes it may be easier to find out the region under the above integral and compute the integral. Then we can find f_{ZW} by taking double derivative:

$$f_{ZW}(z, w) = \frac{\partial^2 F_{ZW}(z, w)}{\partial z \partial w}.$$

Exercise. Suppose x and y are independent uniformly distributed RVs in $[0, 1]$,

$$z = \min\{x, y\} \quad w = \max\{x, y\}, \text{ Find } f_{zw}.$$

- In some cases we may be able to obtain a closed-form expression for f_{zw} in terms of f_{xy} .

Assume that g_1 & g_2 satisfy;

$$(1) \quad z = g_1(x, y), \quad w = g_2(x, y)$$

$$\Leftrightarrow x = h_1(z, w), \quad y = h_2(z, w).$$

$$\text{Example: } z = xy, \quad w = x - y,$$

(2) The functions g_1 & g_2 have continuous partial derivatives at all points (x, y) and are such that

$$J(x, y) = \begin{vmatrix} \frac{\partial g_1}{\partial x} & \frac{\partial g_1}{\partial y} \\ \frac{\partial g_2}{\partial x} & \frac{\partial g_2}{\partial y} \end{vmatrix} = \frac{\partial g_1}{\partial x} \frac{\partial g_2}{\partial y} - \frac{\partial g_1}{\partial y} \frac{\partial g_2}{\partial x} \neq 0$$

at all points (x, y) .

Theorem. In the setting above,

$$f_{\geq \omega}(\geq \omega) = f_{xy}(x,y) | J(x,y)|^{-1}.$$

This is analogous to

$$f_y(y) = \frac{f_x(g^{-1}(y))}{|g'(g^{-1}(y))|},$$

when $y=g(x)$ for a monotonic g .

Lecture 20

(30 October 2023)

Two Functions of Two Random Variables

$(x, y) \sim f_{xy}$.

Let $z = g_1(x, y)$, $w = g_2(x, y)$.

Suppose g_1 & g_2 are such that

(i) \exists functions h_1, h_2 satisfying

$$x = h_1(z, w), \quad y = h_2(z, w).$$

(ii)

$$\begin{aligned} J(x, y) &:= \begin{vmatrix} \frac{\partial g_1(x, y)}{\partial x} & \frac{\partial g_1(x, y)}{\partial y} \\ \frac{\partial g_2(x, y)}{\partial x} & \frac{\partial g_2(x, y)}{\partial y} \end{vmatrix} \\ &= \frac{\partial g_1}{\partial x} \frac{\partial g_2}{\partial y} - \frac{\partial g_2}{\partial x} \frac{\partial g_1}{\partial y} \neq 0 \end{aligned}$$

at all x, y .

Theorem. In the setting above,

$$f_{zw}(z \leq w) = f_{xy}(x \leq y) |J(x,y)|^{-1},$$

where $x = h_1(z \leq w)$, $y = h_2(z \leq w)$.

This is analogous to

$$f_y(y) = \frac{f_x(g^{-1}(y))}{|g'(g^{-1}(y))|},$$

when $y = g(x)$ for a monotonic g .

Proof. Consider

$$P(z \leq Z \leq z + \Delta z, w \leq W \leq w + \Delta w)$$

$$= f_{zw}(z \leq w) \Delta z \Delta w.$$

Alternately,

$$P(z \leq Z \leq z + \Delta z, w \leq W \leq w + \Delta w)$$

$$= P(x \leq X \leq x + \Delta x, y \leq Y \leq y + \Delta y) = f_{xy}(x \leq y) \Delta x \Delta y$$

where $x = h_1(z \leq w)$, $y = h_2(z \leq w)$.

$$\Rightarrow f_{zw}(z, \omega) = f_{xy}(x, y) \frac{\Delta x \Delta y}{\Delta z \Delta w}$$

$$= \frac{f_{xy}(x, y)}{\left(\frac{\Delta z \Delta w}{\Delta x \Delta y} \right)} = \frac{f_{xy}(x, y)}{|J(x, y)|},$$

where

$$J(x, y) = \frac{\partial g_1}{\partial x} \frac{\partial g_2}{\partial y} - \frac{\partial g_1}{\partial y} \frac{\partial g_2}{\partial x}.$$

[See chapter 6
from Papoulis &
Pillai's Book for

Analogy: $g = g(x) \Rightarrow \frac{dg}{dx} = g'(x)$, more details
Similarly $z = g_1(x, y)$ and a formal
 $w = g_2(x, y)$ justification]

$$\Rightarrow \frac{\Delta z \Delta w}{\Delta x \Delta y} = \left| \det \begin{bmatrix} \frac{\partial g_1}{\partial x} & \frac{\partial g_1}{\partial y} \\ \frac{\partial g_2}{\partial x} & \frac{\partial g_2}{\partial y} \end{bmatrix} \right|.$$

- Note that Jacobian is essentially doing the following change of variables:

$$\begin{aligned} P((z, \omega) \in B) &= P((x, y) \in A) = \iint_A f_{xy}(x, y) dx dy \\ &= \iint_B f_{xy}(h_1(z, \omega), h_2(z, \omega)) \cdot \frac{1}{|J(h_1(z, \omega), h_2(z, \omega))|} dz d\omega, \end{aligned}$$

→ more generally, for a given point (z, w) ,
 $f_1(x, y) = z$, $f_2(x, y) = w$ can have many solutions.
 Let $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ represent these multiple solutions such that

$$f_1(x_i, y_i) = z, \quad f_2(x_i, y_i) = w, \quad i \in [1:n].$$

Then

$$f_{zw}(z, w) = \sum_{i=1}^n f_{xy}(x_i, y_i) |J(x_i, y_i)|^{-1}.$$

Example. $\begin{bmatrix} z \\ w \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$.

$$J(x, y) = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{12}a_{21} \neq 0,$$

$$f_{zw}(z, w) = \frac{f_{xy}(x, y)}{|a_{11}a_{22} - a_{12}a_{21}|}, \quad \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}^{-1} \begin{bmatrix} z \\ w \end{bmatrix}.$$

Moment Generating Functions

- A transform associated with a random variable and provides an alternative representation of a probability law.
- Particularly useful for certain types of mathematical manipulations
- Moment generating function (MGF) of a RV X is a function $M_X(s)$ defined as
$$M_X(s) = E[e^{sx}], \text{ for a scalar } s.$$

Discrete case:

$$M_X(s) = \sum_x e^{sx} P_X(x).$$

Continuous case:

$$M_X(s) = \int_{-\infty}^{\infty} e^{sx} f_X(x) dx, \quad \rightarrow M_X(0) = 1.$$

Example. Let $P_X(1) = \frac{1}{2}$, $P_X(2) = \frac{1}{4}$, $P_X(3) = \frac{1}{4}$,

$$M_X(s) = \frac{1}{2} e^s + \frac{1}{4} e^{2s} + \frac{1}{4} e^{3s}.$$

Example. Exponential RV with parameter λ .

$$f_X(x) = \lambda e^{-\lambda x}, x \geq 0,$$

$$M_X(s) = \int_0^\infty e^{sx} \lambda e^{-\lambda x} dx$$

$$= \lambda \int_0^\infty e^{(s-\lambda)x} dx$$

$$= \lambda \left[\frac{e^{(s-\lambda)x}}{(s-\lambda)} \right]_0^\infty$$

$$= \frac{\lambda}{\lambda-s} \quad \text{if } s < \lambda,$$

otherwise the integral is infinite.

— $M_X(s)$ is only defined for those values of s for which $E[e^{sx}]$ is finite.

Exercise.

1) If x takes only non-negative integer values then show that

$$\lim_{s \rightarrow -\infty} m_x(s) = P(x=0).$$

2) Let x be a Poisson RV with parameter λ . Find $m_x(s)$.

3) Prove that $m_{ax+b}(s) = e^{sb} m_x(as)$.

Example. $x \sim N(0,1)$.

$$m_x(s) = \int_{-\infty}^{\infty} e^{sx} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2} + sx} dx$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-s)^2}{2}} \cdot e^{s^2/2} dx$$

$$= e^{s^2/2} \cdot \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-s)^2}{2}} dx$$

$$= e^{s^2/2}.$$

$$Y \sim N(\mu - \sigma^2) \Rightarrow M_Y(s) = e^{sb} m_X(s)$$

$\left[Y = \sigma X + \mu \right]$

$$= e^{sb} \cdot e^{\frac{-\sigma^2 s^2}{2}}$$

$$= e^{\mu s + \frac{-\sigma^2 s^2}{2}}.$$

- MGFs are useful mainly for two reasons

(1) MGFs of X gives us all moments of X .

$$M_X(s) = E[e^{sx}]$$

$$= \int_{-\infty}^{\infty} e^{sx} f_X(x) dx$$

$$\frac{d}{ds} M_X(s) = \int_{-\infty}^{\infty} x e^{sx} f_X(x) dx$$

$$\frac{d}{ds} M_X(s) \Big|_{s=0} = \int_{-\infty}^{\infty} x f_X(x) dx$$

Similarly $\frac{d^n}{ds^n} M_X(s) \Big|_{s=0} = \int_{-\infty}^{\infty} x^n f_X(x) dx.$

(2) MGF (if it exists) uniquely determines

the CDF of random variable.

Consider two random variables x and y , suppose that there exists a positive constant c such that MGFs of x and y are finite and identical for all values of s in $[-c, c]$. Then,

Uniqueness

$$F_x(t) = F_y(t) \text{ for all } t \in \mathbb{R}. \quad \text{Property}$$

Proof follows from connection with Laplace transform.

This can also be seen through characteristic functions.

$$\Phi_x(t) = E[e^{itx}] \quad i = \sqrt{-1}.$$

- characteristic functions always exist.

- Related to Fourier transform.

Fact:

$\int_{-\infty}^{\infty} g(x) dx$ exists if $\int_{-\infty}^{\infty} |g(x)| dx$ exists.

Since $|e^{itx}| = |\Phi_x|$ always exists,

- If x and y are independent RVs

$$M_{x+y}(s) = M_x(s)M_y(s).$$

Exercise . (i) $X \sim N(\mu_x, \sigma_x^2)$, $Y \sim N(\mu_y, \sigma_y^2)$

$$\Rightarrow X+Y \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2).$$

(ii) Let x and y be independent Poisson RVs with parameters λ_1 & λ_2 , respectively.

Then $X+Y$ is Poisson with parameter $\lambda_1 + \lambda_2$

Lecture 21

(2 November 2023)

Module 5 : Probability Bounds & Limit Theorems

— Suppose we want to compute $P(x \geq a)$. In some scenarios it may be sufficient to have bounds on this probability instead of its exact value, e.g., when the distribution of x is unavailable or hard to compute. In such scenarios if we have exact values or bounds for the mean and variance of x , we can obtain meaningful bounds on the quantity of interest.

Markov's Inequality

If x is a non-negative random variable

then $P(x \geq a) \leq \frac{E[x]}{a}$, for all $a > 0$.

Interpretation:

"If $x \geq 0$ and $E[x]$ is small, then the probability that x takes a large value must be small".



$$E[x] = \sum_x x P_x(x), \quad P_x(1000) \rightarrow \text{small}$$

Proof. we first prove that if $x \geq y$, then $E[x] \geq E[y]$.

$z = x - y \geq 0$, i.e., $x(\omega) - y(\omega) \geq 0$ whenever,

$$E[z] \geq 0 \Rightarrow E[x] \geq E[y].$$

Let $y = a \mathbf{1}\{x \geq a\}$, we have

$$x \geq y.$$

$$\Rightarrow E[x] \geq E[y]$$

$$\Rightarrow E[x] \geq E[a \mathbf{1}\{x \geq a\}]$$

$$= a P(x \geq a)$$

$$\therefore P(x \geq a) \leq \frac{E[x]}{a}.$$

Example. Let $X \sim \text{Binomial}(n, p)$. Using Markov's inequality find an upper bound on $P(X \geq \alpha n)$, where $0 < \alpha < 1$. Evaluate the bound for $p = \frac{1}{2}$ and $\alpha = \frac{3}{4}$.

$$P(X \geq \alpha n) \leq \frac{E[X]}{\alpha n} = \frac{np}{n\alpha} = \frac{p}{\alpha}$$

$$= \frac{1}{2} \cdot \frac{4}{3} = \frac{2}{3}.$$

$$\therefore P(X \geq \alpha n) \leq \frac{2}{3}.$$

Example. $X \sim \text{Uniform}[-4, 4]$.

$$P(X \geq 3) \leq P(|X| \geq 3)$$

$$\leq \frac{E[|X|]}{3}$$

$$= \frac{2}{3}.$$

Chebychev's Inequality

If x is a random variable with mean μ and variance σ^2 then

$$P(|x-\mu| \geq c) \leq \frac{\sigma^2}{c^2} \text{ - for all } c > 0,$$

Proof. Let $y = |x-\mu|$.

$$P(y \geq c) = P(y^2 \geq c^2)$$

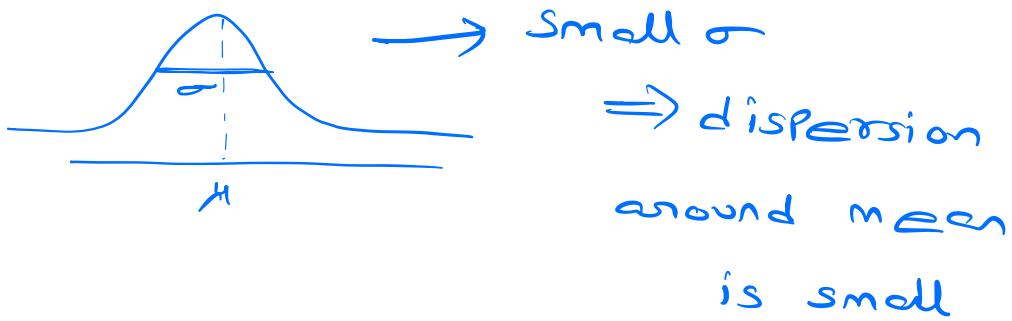
$$\begin{aligned} &\leq \frac{E[y^2]}{c^2} = \frac{E[|x-\mu|^2]}{c^2} \\ &= \frac{\sigma_x^2}{c^2}. \end{aligned}$$

Interpretation:

"If a random variable has small variance then the probability that it takes a value far from its mean is also small."

- Recall that variance measures the spread

of Rv X around its mean.



- An alternative form of chebyshев's inequality:

$$P(|X-\mu| \geq k\sigma) \leq \frac{1}{k^2}.$$

Example. $X \sim \text{Binomial}(n, p)$. Using chebyshev's inequality find an upper bound on $P(X \geq \alpha n)$ where $p < \alpha < 1$. Evaluate for $p = \frac{1}{2}$, $\alpha = \frac{3}{4}$.

$$P(X \geq \alpha n) = P(X - np \geq n(\alpha - p))$$

$$\leq P(|X - np| \geq n(\alpha - p))$$

$$\begin{aligned} &\leq \frac{\text{Var}(X)}{n^2(\alpha - p)^2} = \frac{np(1-p)}{n^2(\alpha - p)^2} \\ &= \frac{4}{n}. \end{aligned}$$

Chernoff Bounds

If X is a random variable then for any $a \in \mathbb{R}$, we can write

$$P(X \geq a) = P(e^{sx} \geq e^{sa}), \text{ for } s > 0$$

$$P(X \leq a) = P(e^{sx} \geq e^{sa}), \text{ for } s < 0.$$

Note that e^{sx} is always a positive random variable for all $s \in \mathbb{R}$. Thus we can apply Markov's inequality.

$$\begin{aligned} \text{For } s > 0, \quad P(X \geq a) &= P(e^{sx} \geq e^{sa}) \\ &\leq \frac{E[e^{sx}]}{e^{sa}} \\ &\leq \min_{s > 0} E[e^{sx}] \cdot e^{-sa}. \end{aligned}$$

Similarly for $s < 0$ -

$$P(X \leq a) \leq \min_{s < 0} E[e^{sx}] e^{-sa}.$$

Example, $X \sim \text{Binomial}(n, p)$. Using Chernoff bounds bound $P(X \geq \alpha n)$ where $p < \alpha < 1$. Evaluate the bound for $p = \frac{1}{2}$ and $\alpha = \frac{3}{4}$.

$$X = \sum_{i=1}^n X_i, \quad P_{X_i}(1) = p = 1 - P_{X_i}(0),$$

X_i are i.i.d.

$$\begin{aligned} E[e^{sx}] &= m_X(s) \\ &= \prod_{i=1}^n m_{X_i}(s) \\ &= (pe^s + 1-p)^n. \end{aligned}$$

$$P(X \geq \alpha n) \leq \min_{s > 0} e^{-\alpha ns} (pe^s + 1-p)^n$$

$$\frac{d}{ds} (e^{-\alpha ns} (pe^s + 1-p)^n) = 0$$

$$\Rightarrow e^s = \frac{\alpha(1-p)}{p(1-\alpha)}$$

$$\Rightarrow s = \log \frac{\alpha(1-p)}{p(1-\alpha)} > 0$$

Since

$$\frac{1-p}{p} \cdot \frac{\alpha}{1-\alpha} > 1,$$

Also check double derivative ≥ 0 .

We get

$$\begin{aligned} P(X \geq \alpha n) &\leq \left(\frac{\alpha(1-p)}{p(1-\alpha)} \right)^{-\alpha n} \left(\frac{p\alpha(1-p)}{p(1-\alpha)} + 1-p \right)^n \\ &= \left(\frac{\alpha(1-p)}{p(1-\alpha)} \right)^{-\alpha n} \left(\frac{1-p}{1-\alpha} \right)^n \\ &= \left(\frac{1-p}{1-\alpha} \right)^{(1-\alpha)n} \left(\frac{p}{1-\alpha} \right)^{\alpha n} \\ &= \left(\frac{1-\frac{1}{2}}{1-\frac{3}{4}} \right)^{\frac{n}{4}} \cdot \left(\frac{2}{3} \right)^{\frac{3n}{4}} \\ &= \frac{2^{\frac{n}{4}} \cdot 2^{\frac{3n}{4}}}{27^{\frac{n}{4}}} = \left(\frac{16}{27} \right)^{\frac{n}{4}}. \end{aligned}$$

Comparison between Markov, Chebyshev
and Chernoff bounds:

$$P(X \geq \alpha n) \leq \frac{2}{3} \quad [\text{Markov}]$$

$$P(X \geq \alpha n) \leq 4/n \quad [\text{Chebyshev}]$$

$$P(X \geq \alpha n) \leq \left(\frac{16}{27} \right)^{\frac{n}{4}} \quad [\text{Chernoff}]$$

- The bound given by Markov is the weakest bound. It is a constant (does not depend on n).
- Chebyshev's bound is stronger than Markov's. In particular note that $4/n \rightarrow 0$ as $n \rightarrow \infty$.
- Chernoff bound is the strongest bound. It goes to zero exponentially fast.

Example. Suppose X is a RV taking values in $[a, b]$. Obtain a bound on $P(|X - \mu| \geq c)$ using Chebyshev's inequality.

$$\begin{aligned}
 \text{claim: } X \in [a, b] \Rightarrow \sigma_x^2 &\leq \frac{(b-a)^2}{4}, \\
 E[(X-\mu)^2] &= E[X^2] - 2E[X]\mu + \mu^2 \\
 \sigma_x^2 &\leq E[(X - \frac{a+b}{2})^2] = E[((X-a) + (X-b))^2]/4 \\
 &= \frac{1}{4} E \left[((X-a) - (X-b))^2 + 4(X-a)(X-b) \right] \\
 &\leq \frac{(b-a)^2}{4}.
 \end{aligned}$$

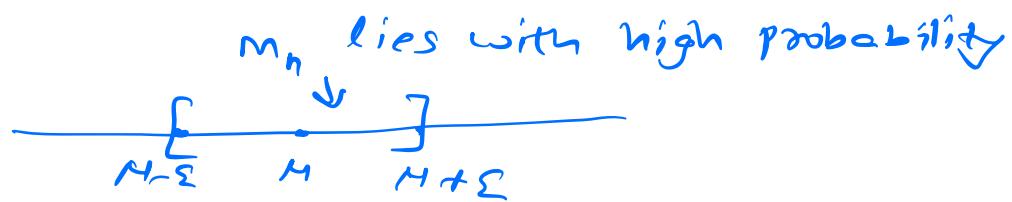
$$P(|X - \mu| \geq c) \leq \frac{\sigma_x^2}{c^2} \leq \frac{(b-a)^2}{4c^2}.$$

Weak Law of Large Numbers (WLLN)

Let x_1, x_2, \dots be a sequence of independent and identically distributed (i.i.d.) random variables with mean μ . For every $\varepsilon > 0$, we have

$$P\left(\left|\frac{\sum_{i=1}^n x_i}{n} - \mu\right| \geq \varepsilon\right) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Interpretation: Let $M_n = \frac{\sum_{i=1}^n x_i}{n}$. WLLN asserts that the sample mean of a large number of i.i.d. RVs is very close to the true mean.



Proof. $E[M_n] = \mu$, $\text{Var}(M_n) = \frac{\sigma^2}{n}$.

$$P(|M_n - E[M_n]| \geq \epsilon)$$

$$= P(|M_n - M| \geq \epsilon) \leq \frac{\text{Var}(M_n)}{\epsilon^2}$$

$$= \frac{\sigma^2}{n\epsilon^2} \rightarrow 0 \quad \text{as} \quad n \rightarrow \infty,$$

Lecture 22

(6 November 2023)

WLLN. Let x_1, x_2, \dots be a sequence of independent and identically distributed RVs with mean μ . For every $\varepsilon > 0$ we have

$$P\left(\left|\frac{\sum_{i=1}^n x_i}{n} - \mu\right| \geq \varepsilon\right) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Remark. Note that the proof we have seen assumes finite variance. It turns out that this law remains true even if the x_i have infinite variance, but a much more elaborate argument is needed.

Another Interpretation:

Let a_1, a_2, \dots, a_n be the realizations of a sequence of i.i.d. random variables $\sim p_x$.

$$a_i \in X = \{x_1, x_2, \dots, x_k\}.$$

$$\sum_{i=1}^n a_i/n = \frac{\sum_{i=1}^K x_i k_i}{n} \rightarrow \sum_{i=1}^K k_i = n$$

k_i = no. of times x_i occurs

If we interpret probabilities as relative frequency

$$\sum_{i=1}^K \left(\frac{k_i}{n}\right) x_i \rightarrow \sum_{i=1}^K P_x(x_i) x_i = E[x].$$

Example (Polling). Each voter in a population selects a particular candidate A with probability p and chooses another with probability 1-p.

$$x_i = \begin{cases} 1 & \text{ith voter chooses A} \\ 0 & \text{o.w.} \end{cases}$$

$$P_{x_i}(1) = p = 1 - P_{x_i}(0).$$

x_i are independent and identically distributed. We are interested in knowing the value of p.

$$M_n = \sum_{i=1}^n x_i / n,$$

$$P(|M_n - p| \geq 0.1) \leq \frac{p(1-p)}{n(0.1)^2} \leq \frac{1}{400(0.1)^2} \quad \text{for } n=100 \\ = 0.25$$

With a sample size of 100 the probability that our estimate is incorrect by more than 0.1 is smaller than 0.25.

$$P(|M_n - p| \geq 0.01) \leq \frac{1}{4n(0.01)^2}$$

$$n \geq 50000 \Rightarrow P(|M_n - p| \geq 0.01) \leq 0.05.$$

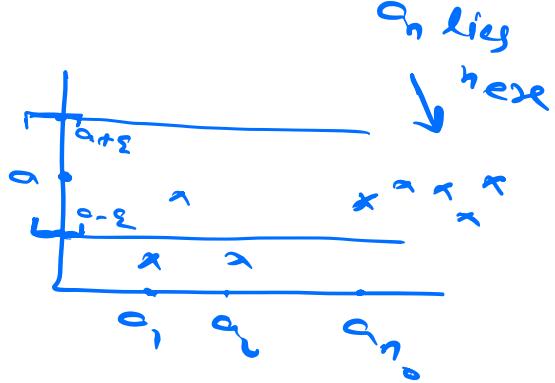
Convergence in Probability

We can interpret WLLN as stating that " $M_n = \frac{1}{n} \sum_{i=1}^n X_i$ converges to mean μ " in the sense of "convergence in probability".

Recall the convergence of a deterministic sequence.

Let a_1, a_2, \dots be a sequence of real numbers and let a be another real. We say that a_n converges to a , or $\lim_{n \rightarrow \infty} a_n = a$, if for every $\epsilon > 0$ there exists some n_0 s.t.

$|Y_n - a| \leq \varepsilon$ for all $n \geq n_0$.



Convergence in Probability

Let Y_1, Y_2, \dots be a sequence of random variables, and let a be a real number. We say that Y_n converges to a in probability if for every $\varepsilon > 0$ we have

$$\lim_{n \rightarrow \infty} P(|Y_n - a| \geq \varepsilon) = 0.$$

- If RVs Y_1, Y_2, \dots have a PMF or a PDF and converge in probability to a , then "almost all" of the PMF or PDF of Y_n is concentrated within ε of a for large values of n .
- The above condition of convergence in probability can be equivalently written as follows: for every $\varepsilon > 0$ and for every $\delta > 0$ there exists some n_0 such that

$$P(|Y_n - a| \geq \varepsilon) \leq \delta \text{ for all } n \geq n_0.$$

$\varepsilon \rightarrow$ accuracy level, $\delta \rightarrow$ confidence level.

Example. Let x_1, x_2, \dots be a sequence of independent and uniformly distributed rvs in $[0, 1]$, and let $y_n = \min\{x_1, \dots, x_n\}$.

Note that $y_1 \geq y_2 \geq \dots$. In fact,

$y_n \rightarrow 0$ in probability.

$$\begin{aligned} P(|y_n - 0| \geq \varepsilon) &= P(x_1 \geq \varepsilon, \dots, x_n \geq \varepsilon) \\ &= (1 - \varepsilon)^n \rightarrow 0 \text{ as } n \rightarrow \infty \end{aligned}$$

for every $\varepsilon > 0$.

- The following example shows that if $y_n \rightarrow a$, then $E[y_n]$ may not converge to a .

Example. $P(y_n = y) = \begin{cases} 1/n, & \text{for } y = 0 \\ 1/n, & \text{for } y = n^2 \\ 0, & \text{elsewhere} \end{cases}$

For every $\varepsilon \geq 0$ $P(|y_n - 0| \geq \varepsilon) = \frac{1}{n} \rightarrow 0$ as $n \rightarrow \infty$.

$E[y_n] = n \rightarrow \infty$ as $n \rightarrow \infty$,

i.e., $E[y_n] \not\rightarrow a$.

Central Limit Theorem (CLT)

Let x_1, x_2, \dots be a sequence of independent and identically distributed RVs with common mean μ and variance σ^2 .

Define $Z_n = \frac{\sum_{i=1}^n x_i - n\mu}{\sqrt{n}\sigma}$. Then

$$\begin{aligned}\lim_{n \rightarrow \infty} P(Z_n \leq z) &= \Phi(z) \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-x^2/2} dx.\end{aligned}$$

$$- S_n = \sum_{i=1}^n x_i \quad \text{---} \quad \text{Var}(S_n) = n\sigma^2 \quad \rightarrow \infty \text{ as } n \rightarrow \infty$$

$$M_n = \sum_{i=1}^n x_i / n \quad \text{---} \quad \text{Var}(M_n) = \frac{\sigma^2}{n} \quad \rightarrow 0 \text{ as } n \rightarrow \infty$$

$$Z_n = \sum_{i=1}^n x_i / \sqrt{n} \quad \text{---} \quad \text{Var}(Z_n) = 1 \text{ independent of } n \quad (\text{constant})$$

Proof of CLT:

Fact, If $M_{Z_n}(s) \rightarrow M_Z(s)$ then

$$F_{Z_n}(z) \rightarrow F_Z(z) \text{ for all } z.$$

Proof of the fact is related to the continuity of inverse Fourier transform,

We assume that $M_X(s)$ is finite when $-d < s < d$. Let x_i has zero mean & variance 1.

$$Z_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i / \sqrt{n}.$$

$$M_{Z_n}(s) = E \left[e^{s \sum_{i=1}^n x_i / \sqrt{n}} \right]$$

$$= \prod_{i=1}^n E \left[e^{s x_i / \sqrt{n}} \right] = M_X(s/\sqrt{n})^n.$$

$$\text{Let } \angle(t) = \log M_X(t).$$

$$\angle(0) = 0, \quad \angle'(0) = \frac{m_x'(0)}{m_x(0)} = E[x] = 0,$$

$$\angle''(0) = \frac{m_x(0)m_x''(0) - m_x'(0)^2}{m_x(0)^2} = E[x^2] = 1,$$

we first show that

$$m_x(s/\sqrt{n}) \xrightarrow{\text{MGF of } N(0,1)} e^{t^2/2}$$

or equivalently that

$$n\angle(s/\sqrt{n}) \rightarrow t^2/2 \text{ as } n \rightarrow \infty.$$

$$\lim_{n \rightarrow \infty} \frac{\angle(s/\sqrt{n})}{n^{-1}} = \lim_{n \rightarrow \infty} \frac{-\angle'(s/\sqrt{n}) \cdot s}{-\frac{2}{n^2}} \quad (\text{by L'Hopital Rule})$$

$$= \lim_{n \rightarrow \infty} \frac{\angle'(s/\sqrt{n}) \cdot s}{2n^{-3/2}}$$

$$= \lim_{n \rightarrow \infty} \frac{-\angle''(s/\sqrt{n}) s^2}{-2n^{-3/2}} \quad (\text{by L'Hopital Rule})$$

$$= \lim_{n \rightarrow \infty} \frac{\angle''(s/\sqrt{n}) s^2}{2} = \frac{s^2}{2}, \quad (\because \angle''(0) = 1)$$

If x_i 's are of mean μ & variance σ^2

$$\sum_{i=1}^n \frac{x_i - \mu}{\sqrt{n}\sigma} = \frac{\sum_{i=1}^n \frac{x_i - \mu}{\sigma}}{\sqrt{n}} \rightarrow \begin{matrix} \text{mean } 0 \\ \text{var} = 1 \end{matrix}$$

Fact. If $m_{Z_n}(t) \rightarrow m_Z(t)$ then

$$F_{Z_n}(z) \rightarrow F_Z(z) \text{ for all } z.$$

Proof is related inverse Fourier transform
and its continuity.

This completes the proof of the central limit theorem.

Lecture 23

(9 November 2023)

Central Limit Theorem.

Let x_1, x_2, \dots be a sequence of i.i.d. random variables with common mean μ and variance σ^2 . Define

$$Z_n = \frac{\sum_{i=1}^n x_i - n\mu}{\sigma\sqrt{n}}. \quad \text{Then, the CDF of}$$

Z_n converges to the standard normal CDF in the sense that

$$\begin{aligned} \lim_{n \rightarrow \infty} P(Z_n \leq z) &= \Phi(z) \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-x^2/2} dx, \end{aligned}$$

for every z .

Proof. we prove this for $\mu=0$ & $\sigma^2=1$,
i.e., $Z_n = \frac{\sum_{i=1}^n x_i}{\sqrt{n}}$.

$$\begin{aligned}
 M_{Z_n}(s) &= E[e^{sZ_n}] \\
 &= E\left[e^{s\sum_{i=1}^n X_i/\sqrt{n}}\right] \\
 &= \prod_{i=1}^n M_X(s/\sqrt{n}) = (M_X(s/\sqrt{n}))^n.
 \end{aligned}$$

Let $\angle(s) = \log M_X(s)$.

$$\log M_{Z_n}(s) = n \log M_X(s/\sqrt{n}) = n \angle(s/\sqrt{n}).$$

we have $\angle(0) = \angle'(0) = 0$ & $\angle''(0) = 1$,

$$\text{Consider } \lim_{n \rightarrow \infty} n \log M_X(s/\sqrt{n})$$

$$= \lim_{n \rightarrow \infty} \frac{\log M_X(s/\sqrt{n})}{n^{-1}}$$

$$= \lim_{n \rightarrow \infty} \frac{M_X'(s/\sqrt{n}) s n^{-3/2}}{2 M_X(s/\sqrt{n}) n^{-2}}$$

$$= \lim_{n \rightarrow \infty} \frac{\angle'(s/\sqrt{n}) s}{2 n^{-1/2}}$$

$$= \lim_{n \rightarrow \infty} \frac{\angle''(s/\sqrt{n}) s^2 n^{-3/2}}{2 n^{-3/2}} = s^2/2.$$

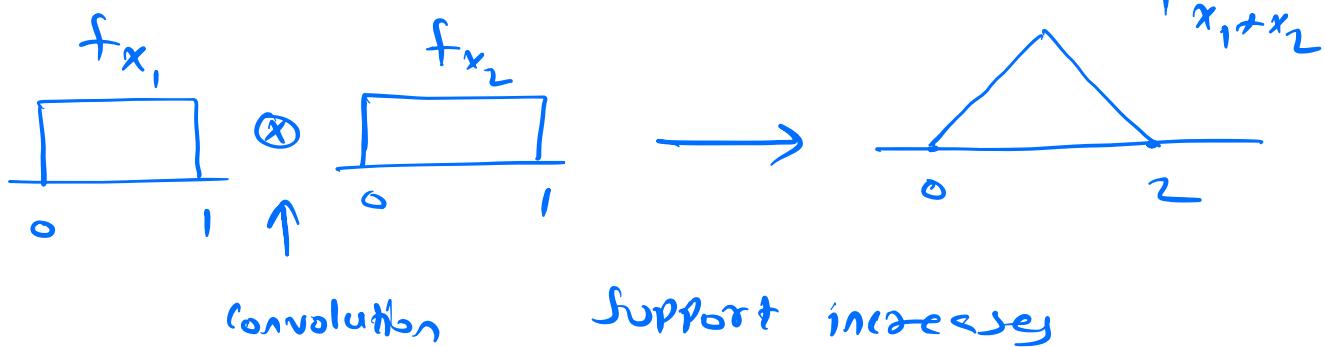
Examples. (i) $X_i \sim N(\mu, \sigma^2)$

$$\frac{\sum_{i=1}^n X_i - n\mu}{\sqrt{n}\sigma} \sim N(0, 1) \text{ for all } n \in \mathbb{N},$$

$$\frac{X_1 - \mu}{\sigma}, \frac{X_1 + X_2 - 2\mu}{\sqrt{2}\sigma}, \dots \sim N(0, 1),$$

$$\lim_{n \rightarrow \infty} P(Z_n \leq z) = P(Z \leq z), \quad Z \sim N(0, 1).$$

(ii) $X_i \sim \text{uniform}[0, 1]$.



Intuitively as $n \rightarrow \infty$, support will be real line
and $\sum_{i=1}^n X_i \rightarrow \text{Gaussian RV.}$

Normal Approximation Based on CLT:

Let $S_n = \sum_{i=1}^n x_i$, where x_i are i.i.d.

with mean μ and variance σ^2 . We are interested in computing $P(S_n \leq c)$. If n is large, it can be approximated in the following way.

$$P(S_n \leq c) = P\left(\frac{S_n - n\mu}{\sqrt{n}\sigma} \leq \frac{c - n\mu}{\sqrt{n}\sigma}\right)$$

$$\approx \Phi(z) \text{ where } z = \frac{c - n\mu}{\sqrt{n}\sigma}.$$

Example, x_i , $i \in [1:100]$ are i.i.d. and uniformly distributed in the interval $[5, 50]$. Find $P(S_{100} > 3000)$, $S_{100} = \sum_{i=1}^{100} x_i$. It is not easy to calculate CDF and the desired probability exactly, but an approximate answer can be quickly obtained using CLT.

$$\mu = \frac{5+50}{2} = 27.5, \quad \sigma^2 = \frac{(50-5)^2}{12} = 168.75.$$

$$P(S_{100} > 3000) = 1 - P(S_{100} \leq 3000)$$

$$= 1 - P\left(\frac{S_{100} - n\mu}{\sqrt{n}\sigma} \leq \frac{3000 - 2750}{10(168.75)}\right)$$

$$= 1 - P\left(\frac{S_{100} - n\mu}{\sqrt{n}\sigma} \leq 1.92\right)$$

$$= 1 - \Phi(1.92)$$

$$\approx 1 - 0.9726$$

(using standard Normal table)

$$= 0.0274.$$

Example (Polling). $x_i = \begin{cases} 1 & \text{if } i\text{th voter votes for A} \\ 0 & \text{otherwise} \end{cases}$

$$E[x_i] = p \quad M_n = \sum_{i=1}^n x_i / n .$$

WLLN (or Chebyshev's inequality):

For $n=100$, we got that

$$P(|M_n - p| \geq 0.1) \leq 0.25 .$$

We will see now if CLT can improve this bound.

$$P\left(\left|\frac{\sum_{i=1}^n x_i}{n} - p\right| \geq 0.1\right)$$

$$= P\left(\left|\frac{\sum_{i=1}^n x_i - np}{n}\right| \geq 0.1n\right)$$

$$= P\left(\left|\frac{\sum_{i=1}^n x_i - np}{\sqrt{n}\sigma}\right| \geq \frac{0.1\sqrt{n}}{\sigma}\right)$$

$$= P\left(\frac{\sum_{i=1}^n x_i - np}{\sqrt{n}\sigma} \geq \frac{0.1\sqrt{n}}{\sigma}\right) + P\left(\frac{\sum_{i=1}^n x_i - np}{\sqrt{n}\sigma} \leq -\frac{0.1\sqrt{n}}{\sigma}\right)$$

$$\approx 2P\left(Z \geq \frac{0.1\sqrt{n}}{\sigma}\right) \sim N(0,1)$$

$$\leq 2P\left(Z \geq 0.2\sqrt{n}\right) \text{ as } \sigma \leq \frac{1}{2}$$

$$= 2 - 2P(Z \leq 0.2\sqrt{n})$$

$$= 2 - 2\Phi(-0.2\sqrt{100}) = 2 - 2\Phi(-2) = 2 - 2 \times 0.977 = 0.046.$$

This is a much better estimate than 0.25 from WLLN.

Strong Law of Large Numbers

- SLLN also deals with the convergence of the sample mean to the true mean. However, SLLN refers to a different type of convergence.

WLLN

$$M_n = \frac{1}{n} \sum_{i=1}^n x_i \rightarrow E[x] \text{ in Probability}$$

SLLN

$$M_n = \frac{1}{n} \sum_{i=1}^n x_i \rightarrow E[x] \text{ with probability 1,}$$

(or almost sure
convergence)

Almost Sure convergence:

Let x_1, x_2, \dots be a sequence of RVS and c be a real number. Then we say $x_n \xrightarrow{a.s.} c$ if $P(\{\omega \in \Omega : \lim_{n \rightarrow \infty} x_n(\omega) = c\}) = 1$.

Example. Let $\Omega = [0, 1]$, consider a probability law defined by $P([a, b]) = b - a$ for all $0 \leq a \leq b \leq 1$.

Define $X_n(\omega) = \omega^n$ for $n \in \mathbb{N}$.

Note that

$$\lim_{n \rightarrow \infty} X_n(\omega) = \begin{cases} 0 & \text{if } 0 \leq \omega < 1 \\ 1 & \text{if } \omega = 1. \end{cases}$$

$$P(\{\omega : \lim_{n \rightarrow \infty} X_n(\omega) = 0\}) = P([0, 1)) = 1$$

Since the singleton set $\{1\}$ has zero probability,

$$\therefore X_n \xrightarrow{\text{a.s.}} 0.$$

SLLN. Let x_1, x_2, \dots be a sequence of i.i.d. RVs with mean μ . Then,

$$M_n = \frac{\sum_{i=1}^n x_i}{n} \xrightarrow{\text{a.s.}} \mu, \quad \text{i.e.,}$$

$$P(\{\omega : \lim_{n \rightarrow \infty} M_n(\omega) = \mu\}) = 1.$$

Proof. Assume that $E[x_i^4] = k < \infty$.

$$S_n = \sum_{i=1}^n x_i. \quad E[S_n^4] = E\left[\left(\sum_{i=1}^n x_i\right)^4\right].$$

This will have terms of the form

$$x_i^4, x_i^3 x_j, x_i^2 x_j^2, x_i^2 x_j x_k, x_i x_j x_k x_l$$

where i, j, k, l are different.

Assume $n=0$. Then because of independence it follows that

$$E[x_i^3 x_j] = E[x_i^3] E[x_j] = 0$$

$$E[x_i^2 x_j x_k] = E[x_i^2] E[x_j] E[x_k] = 0$$

$$E[x_i x_j x_k x_l] = E[x_i] E[x_j] E[x_k] E[x_l] = 0.$$

$$\begin{aligned} \text{so } E[S_n^4] &= n E[x_i^4] + 6 \binom{n}{2} E[x_i^2 x_j^2] \\ &= n E[x_i^4] + 3n(n-1) E[x_i^2] E[x_j^2] \\ &\qquad\qquad\qquad \underbrace{\qquad\qquad\qquad}_{= (E[x_i^2])^2 \leq E[x_i^4]} \\ &\leq nk + 3n(n-1)k \stackrel{\text{as } \text{var}(x_i^2) \geq 0}{\leq} 3n^2 k \\ \Rightarrow E[S_n^4/n^4] &\leq 3k/n^2 \end{aligned}$$

$$\Rightarrow E\left[\sum_{n=1}^{\infty} \frac{s_n^4}{n^4}\right] = \sum_{n=1}^{\infty} E[s_n^4/n^4]$$

$$= \sum_{n=1}^{\infty} \frac{1}{n^2} \cdot (3K) < \infty.$$

This implies that

$$P\left(\sum_{n=1}^{\infty} \frac{s_n^4}{n^4} < \infty\right) = 1$$

$$\sum_{n=1}^{\infty} \frac{s_n^4}{n^4} < \infty \Rightarrow \lim_{n \rightarrow \infty} \frac{s_n^4}{n^4} = 0$$

$$\Rightarrow 1 = P\left(\sum_{n=1}^{\infty} \frac{s_n^4}{n^4} < \infty\right) \leq P\left(\lim_{n \rightarrow \infty} \frac{s_n^4}{n^4} = 0\right)$$

$$\Rightarrow P\left(\lim_{n \rightarrow \infty} \frac{s_n}{n} = 0\right) = 1.$$

In the above proof we have used

$$E\left[\sum_{i=1}^{\infty} z_i\right] = \sum_{i=1}^{\infty} E[z_i], \text{ This is not}$$

necessarily true always for any z_i 's.

However this holds true when all z_i are non-negative random variables. This is because of monotone convergence theorem (not covered in this course).

Lecture 24

(16 November 2023)

Random Processes

- A random or stochastic process is a mathematical model of a probabilistic experiment that evolves in time and generates a sequence of numerical values.
- Each numerical value in the sequence is modeled by a random variable so a random process is simply a (finite or infinite) sequence of random variables.

Recall that a RV $X : \Omega \rightarrow \mathbb{R}$,

Formally, a random process is a family of random variables $(X_t : t \in T)$, all on the same probabilistic model (Ω, P) .

In many applications, T is a set of times. If $T = \mathbb{N}$ then $(X_t : t \in T)$ is called a discrete-time process. If $T = \mathbb{R}$, then continuous-time process.

- For each $t \in T$, x_t is a RV.

$$\omega \mapsto x_1(\omega) \ x_2(\omega) \ \dots$$

Discrete-time

$$\omega \mapsto x_t(\omega) \quad t \in T.$$

Continuous-time

For a fixed $\omega \in \Omega$, $(x_t(\omega), t \in T)$ is called the sample path at ω .

Mean Function:

$$\mu_x(t) = E[x_t] \text{ for } t \in T,$$

Covariance Function:

$$\begin{aligned} C_x(t_1, t_2) &= \text{Cov}(x_{t_1}, x_{t_2}) \\ &= E[x_{t_1} x_{t_2}] - E[x_{t_1}] E[x_{t_2}] \end{aligned}$$

Correlation function:

$$R_x(t_1, t_2) = E[x_{t_1} x_{t_2}].$$

- A random process is statistically specified by its complete set of n th order distribution function for all $n \in \mathbb{N}$,

$$F_x(x_1, \dots, x_n; t_1, t_2, \dots, t_n) = P(X_{t_1} \leq x_1, \dots, X_{t_n} \leq x_n).$$

For the continuous case we can obtain the PDF as

$$f(x_1, x_2; t_1, t_2) = \frac{\partial^2 F(x_1, x_2; t_1, t_2)}{\partial x_1 \partial x_2},$$

Arrival - Type Processes we are interested in occurrences that have the character of an arrival such as message receptions at a receiver, customer purchases at a store etc. we will focus on models in which the interarrival times (the times between successive arrivals) are independent random variables.

Bernoulli process - case where arrivals occur in discrete time and the interarrival times are geometrically distributed

Poisson process - case where arrivals occur in continuous time and the interarrival times are exponentially distributed.

The Bernoulli process

Consider a sequence of independent coin tosses with the probability of heady p in the range $0 < p < 1$.

- A Bernoulli process is a sequence x_1, x_2, \dots of independent Bernoulli Rvs x_i with
 - $P(x_i = 1) = P(\text{success at the } i\text{th trial})$
 - $P(x_i = 0) = P(\text{failure at the } i\text{th trial})$
- No. of arrivals within a certain time period or no. of successes in n independent trials:

$$P_s(k) = \binom{n}{k} p^k (1-p)^{n-k} \quad k=0, 1, \dots, n$$

$$E[S] = np, \quad \text{var}(S) = np(1-p).$$

This is Binomial with parameters n & p .

- The time until the first arrival or no. of trials up to (and including) the first success :

$$P_T(t) = (1-p)^{t-1} p, \quad t=1, 2, \dots$$

$$E[T] = \frac{1}{p}, \quad \text{var}(T) = \frac{1-p}{p^2}.$$

Independence & Memorylessness

Due to independence property in the Bernoulli process — whatever has happened in the past trials provides no information on the outcome of future trials.

Fresh-start property :

For any given n , consider

$x_1, x_2, \dots, x_n, x_{n+1}, x_{n+2}, \dots$, we notice that

$x_i = x_{n+i}$ are independent Bernoulli trials and therefore form a Bernoulli process,

Memolessness property: T = time until k^{th} success

$$P(T-n=t | T>n)$$

$$= \frac{P(T-n=t - T>n)}{P(T>n)}$$

$$= \frac{P(T=n+t)}{P(T>n)} = \frac{(1-p)^{n+t-1} \cdot p}{(1-p)^n}$$

$$= (1-p)^{t-1} p = P(T=t),$$

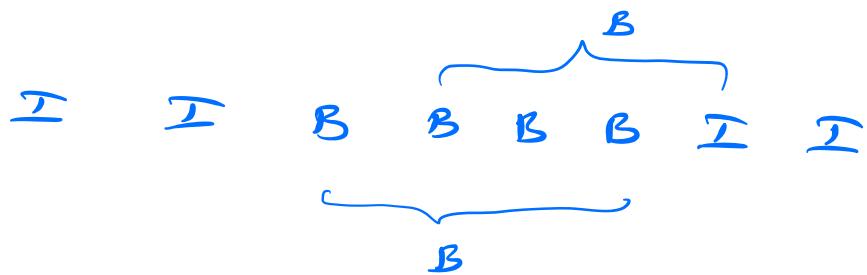
- we call a slot *busy* (B) if there is an arrival, and otherwise let us call it *idle*, we call a string of idle slots flanked by busy slots, an 'idle period'. Similarly, we can define 'busy period'.

Let us derive the PMF, mean, and variance of the following random variables.

T = the time index of the first idle slot

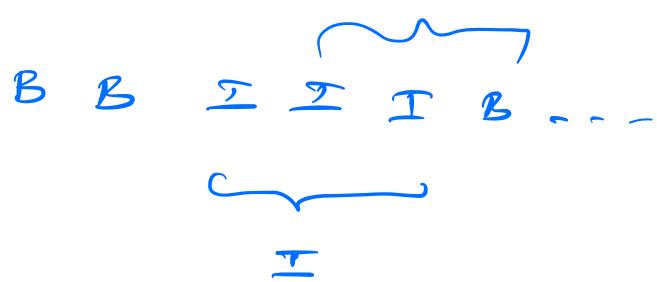
$$P_T(k) = p^{k-1} (1-p) \quad k=1, 2, \dots$$

\angle = the length (no. of slots) of the first busy period



$$P_B(k) = p^{k-1} (1-p) \quad k=1, 2, \dots$$

Σ = the length of the first idle period



$$P_I(k) = (1-p)^{k-1} p \quad k=1, 2, \dots$$

Example. Let N be the first time that we have a success immediately following a previous success. That is N is the first time i for which $x_{i-1} = x_i = 1$. What is the probability that there are no successes in the two trials that follow, i.e., $P(x_{N+1} = x_{N+2} = 0)$?

$$P(X_{N+1} = x_{N+2} = 0) \\ = \sum_{n=1}^{\infty} P(X_{N+1} = x_{N+2} = 0 \mid N=n) P(N=n)$$

$$P(X_{N+1} = x_{N+2} = 0 \mid N=n) \\ = P(X_{n+1} = x_{n+2} = 0) = (1-p)^2$$

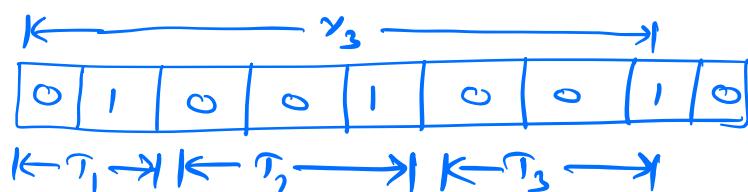
Equality because $N=n$ is determined by some function of (x_1, x_2, \dots, x_n) .

Interarrival Times

Let y_k denote the time of the k^{th} success (or arrival).

$$T_1 = y_1 \quad T_k = y_k - y_{k-1}, \quad k=2, 3, \dots$$

T_k represents the k^{th} interarrival time i.e., no. of trials following the $(k-1)^{\text{st}}$ success until the next success.



We have $y_k = \sum_{i=1}^k T_i$.

Interarrival times T_1, T_2, T_3, \dots are independent and all have the same geometric distribution.

Properties of the k^{th} Arrival Time

$$y_k = T_1 + T_2 + \dots + T_k$$

Each T_i is a geometric RV with parameter p and all T_i 's are independent.

$$E[y_k] = \sum_{i=1}^k E[T_i] = k/p,$$

$$\text{Var}(y_k) = \sum_{i=1}^k \text{Var}(T_i) = \frac{k(1-p)}{p^2},$$

To find the pmf of y_k

$$P(y_k = t) = P(A \cap B), \text{ where}$$

$$A = \left\{ \text{trial } t \text{ is a success} \right\}$$

$$B = \left\{ \text{exactly } k-1 \text{ successes in } t-1 \text{ trials} \right\}$$

$$P_{Y_k}(t) = P(A), P(B)$$

$$= p \cdot \binom{t-1}{k-1} p^{k-1} (1-p)^{t-k}$$

$$= \binom{t-1}{k-1} p^k (1-p)^{t-k},$$



Known as Pascal PMF of order k.

Poisson Process

- A counting process counts the number of arrivals, e.g., no. of customers who arrive at a restaurant — $X(t) = \text{no. of arrivals in time } t$.
- A counting process has independent and stationary increments. $[X(0)=0]$

Independent increments:

For all $0 \leq t_1 < t_2 < t_3 < \dots < t_n$ the RVs

$X(t_2) - X(t_1)$, $X(t_3) - X(t_2)$, ... are independent.

Note that $x(t_i) - x(t_{i-1})$ represents the no. of arrivals in the interval $(t_{i-1}, t_i]$.

Stationary increments:

For all $t_2 > t_1 \geq 0$ & $\tau > 0$, the RVs

$x(t_2) - x(t_1)$ and $x(t_2 + \tau) - x(t_1 + \tau)$ have the same distribution.

- Thus a counting process has independent increments if the no. of arrivals in non-overlapping intervals are independent.

Also it has stationary increments if, for all $t_2 > t_1 \geq 0$, $x(t_2) - x(t_1)$ has the same distribution as $x(t_2 - t_1)$.

- Before we define Poisson process we recall Poisson RV.

$X \sim \text{Poisson}(\lambda)$ means $P_x(k) = \frac{\lambda^k e^{-\lambda}}{k!}$ - $k=0, 1, 2, \dots$

- A counting process $(N(t), t \in [0, \infty))$ is called a Poisson process with rate λ if:

- (1) $N(0) = 0$,
- (2) $N(t)$ has independent increments,
- (3) the no. of arrivals in any interval of length $T > 0$ has Poisson (λT) distribution.

$$E[N(t)] = \lambda t, \quad \text{var}(N(t)) = \lambda t.$$

$$\begin{aligned} R_X(t_1, t_2) &= E[N(t_1)N(t_2)] \quad (t_1 < t_2) \\ &= E[N(t_1)(N(t_2) - N(t_1) + N(t_1))] \\ &= E[N(t_1)(N(t_2) - N(t_1))] + E[N^2(t_1)] \\ &= E[N(t_1)] E[N(t_2 - t_1)] + E[N^2(t_1)] \\ &= \lambda t_1 \cdot \lambda(t_2 - t_1) + \lambda t_1 + \lambda^2 t_1^2 \\ &= \lambda^2 t_1 t_2 - \cancel{\lambda^2 t_1^2} + \lambda t_1 - \cancel{\lambda^2 t_1^2} \\ &= \lambda t_1 + \lambda^2 t_1 t_2. \end{aligned}$$

$$C(t_1, t_2) = \lambda \min\{t_1, t_2\},$$

Lecture 25

(18 November 2023)

- A random process $X = (X_t; t \in T)$ is specified by n^{th} -order distribution function,

$$F_{X_{t_1}, \dots, X_{t_n}}(x_1, x_2, \dots, x_n)$$

$$= P(X_{t_1} \leq x_1, \dots, X_{t_n} \leq x_n)$$

for all $n \in \mathbb{N}$,

Example. Consider the random process

$$X(t) = A \sin(\omega_0 t + \phi)$$

where A and ϕ are independent and ϕ is uniformly distributed over $[-\pi, \pi]$.

$$M_x(t) = E[A \sin(\omega_0 t + \phi)]$$

$$= M_A \cdot E[\sin(\omega_0 t + \phi)]$$

$$= M_A \cdot \int_{-\pi}^{\pi} \frac{1}{2\pi} \cdot \sin(\omega_0 t + \phi) d\phi$$

$$= 0.$$

$$\begin{aligned}
 R_x(t_1, t_2) &= E[x(t_1)x(t_2)] \\
 &= E[A^2 \sin(\omega_0 t_1 + \phi) \sin(\omega_0 t_2 + \phi)] \\
 &= E[A^2] E[\sin(\omega_0 t_1 + \phi) \sin(\omega_0 t_2 + \phi)] \\
 &= E[A^2] \frac{1}{2} \left(E[\cos(\omega_0(t_1 - t_2))] - E[\cos(\omega_0(t_1 + t_2) + 2\phi)] \right) \\
 &\quad \qquad \qquad \qquad \underbrace{\qquad\qquad\qquad}_{=0} \\
 &= \frac{1}{2} E[A^2] [\cos(\omega_0(t_1 - t_2))].
 \end{aligned}$$

Stationary Process

A random process $(x(t), t \in \mathbb{R})$ is stationary if, for all $t_1, t_2, \dots, t_n, T \in \mathbb{R}$,

$$F_{x(t_1), x(t_2), \dots, x(t_n)}(x_1, x_2, \dots, x_n) = F_{x(t_1+T), x(t_2+T), \dots, x(t_n+T)}(x_1, x_2, \dots, x_n)$$

$\forall n \in \mathbb{N},$

Exercise. Write down a similar definition for discrete-time random process

Example. Consider i.i.d. discrete-time random process x_1, x_2, \dots , Is this stationary?
Yes!

$(x_{n_1}, x_{n_2}, \dots, x_{n_8})$ has same distribution as
 $(x_{n_1+T}, \dots, x_{n_8+T})$.

— If a process is stationary, the analysis is usually simpler as the probabilistic properties do not change by time.

However, it turns out that not many real-life processes are stationary.
Even if a process is stationary, it might be difficult to prove it. Fortunately, often a weaker notion of stationarity suffices.

Wide-Sense stationary (WSS)

A random process $x(t)$ is called WSS if for all $t, t_1, t_2 \in \mathbb{R}, T \in \mathbb{R}$,

$$M_x(t) = M_x(t+T)$$

$$R_x(t_1, t_2) = R_x(t_1+T, t_2+T).$$

so $M_x(t)$ is a constant and $R_x(t_1, t_2)$ is a function of $t_2 - t_1$.

Example. $x(t) = A \sin(\omega_0 t + \theta)$ — A, θ are independent & $\theta \sim \text{uniform } [-\pi, \pi]$.

$x(t)$ is WSS because

$$M_x(t) = 0 \quad (\text{a constant})$$

$$\begin{aligned} R_x(t_1, t_2) &= \frac{1}{2} E[A^2] \cos(\omega_0(t_2 - t_1)) \\ &= R_x(t_2 - t_1). \end{aligned}$$

- For WSS random process it suffices to denote autocorrelation by $R_x(T)$.

Exercise. Prove that a stationary process is WSS.

Properties of $R_x(\tau)$

- $R_x(0) = E[x(t)x(t)] = E[x^2(t)] \geq 0$,
 $E[x^2(t)]$ is called expected (or average)
power in $x(t)$ at time t . For wss
this is not a function of time.
- $R_x(-\tau) = E[x(t+\tau)x(t)]$
= $E[x(t)x(t+\tau)]$
= $R_x(\tau)$.
- $|R_x(\tau)| \leq R_x(0)$ for all $\tau \in \mathbb{R}$.
This follows from Cauchy-Schwarz's
inequality.

$$|E[xy]| \leq \sqrt{E[x^2]E[y^2]}.$$

Consider $E[(x-\alpha y)^2] \geq 0$

$$\Rightarrow E[x^2 + \alpha^2 y^2 - 2\alpha xy] \geq 0$$
$$\Rightarrow \alpha^2 E[y^2] - 2\alpha E[xy] + E[x^2] \geq 0$$

So discriminant ≤ 0

$$\Rightarrow 4E^2[xy] \leq 4E[x^2]E[y^2]$$

$$\Rightarrow |E[xx^T]| \leq \sqrt{E[x^2]E[g^2]}.$$

$$|E[x(t)x(t+\tau)]| \leq \sqrt{E[\tilde{x}(t)]E[\tilde{x}(t+\tau)]}$$

$$\Rightarrow |R_x(\tau)| \leq \sqrt{R_x(0) \cdot R_x(0)} = R_x(0).$$

Power Spectral Density (PSD)

Power spectral density is the Fourier transform of $R_x(\tau)$, mathematically

$$S_x(f) = \int_{-\infty}^{\infty} R_x(\tau) e^{-j2\pi f \tau} d\tau, \quad j = \sqrt{-1}$$

is PSD of $x(t)$,

It can be shown that the inverse Fourier transform of $S(f)$ is $R_x(\tau)$.

$$\int_{-\infty}^{\infty} S_x(f) e^{j2\pi f \tau} df = R_x(\tau)$$

Properties of PSD

- $S_x(-f) = S_x(f)$.

This is because $R_x(\tau)$ is even.

- $E[x(t)^2] = R_x(0) = \int_{-\infty}^{\infty} S_x(f) df$.

Inverse Fourier Transform

so we get expected or average power in $x(t)$ by integrating the PSD of $x(t)$.
This is why $S_x(f)$ is called power spectral density.

- $S_x(f) \geq 0$.

(Proof omitted)

Lecture 26

(20 November 2023)

Q) (a) Is an event A independent of itself? No in general.

$$P(A \cap A) = P(A)^2 \Rightarrow P(A) = 0 \text{ or } 1.$$

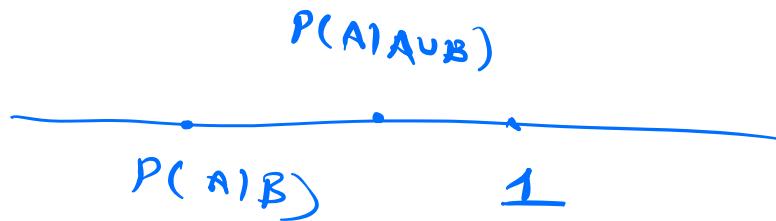
(b) Find the probability that exactly one of the events A or B occurs in terms of $P(A)$, $P(B)$, $P(A \cap B)$.

$$\begin{aligned} P((A \cap B^c) \cup (A^c \cap B)) &= P(A \cap B^c) + P(A^c \cap B) \\ &= P(A) + P(B) - 2P(A \cap B). \end{aligned}$$

(c) Is it true that $P(A|A \cup B) \geq P(A|B)$?

$$\begin{aligned} P(A|A \cup B) &= P(A|(A \cup B) \cap B) P(B|A \cup B) + \\ &\quad P(A|(A \cup B) \cap B^c) P(B^c|A \cup B) \\ &= P(A|B) P(B|A \cup B) + P(A|A \cap B^c) P(B^c|A \cup B) \\ &\quad \underbrace{\qquad\qquad}_{=1} \end{aligned}$$

$$= \neg P(A|B) + (1-\gamma),$$



$$\therefore P(A|A \cup B) \geq P(A|B),$$

Q) If X is a positive integer valued RV that satisfies memorylessness property

$$P(X > m+n | X > m) = P(X > n), \text{ for any } m, n \in \mathbb{N}.$$

Then prove that X is a geometric random variable.

Proof. Let $P(X > n) = a_n$ for $n \in \mathbb{N}$,

The memorylessness property gives

$$\frac{a_{m+n}}{a_m} = a_n$$

$$\Rightarrow a_{m+n} = a_m a_n, \forall m, n \in \mathbb{N}$$

$$\Rightarrow a_{m+1} = a_m a_1 = a_1^{m+1}$$

where $a_1 = P(X > 1) = 1 - P(X = 1) \triangleq 1 - p$,
 $P(X = n) = a_{n+1} - a_n = (1-p)^{n+1} - (1-p)^n = (1-p)^{n-1}p$.

Q) Let X be a discrete RV, Y be a continuous RV, and I is a binary RV s.t. X is independent of I , Y is independent of I .

Define

$$Z = \begin{cases} X & \text{if } I=1 \\ Y & \text{if } I=0 \end{cases}$$

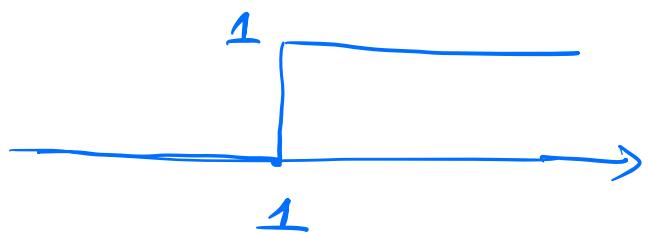
$$Z = \begin{cases} 1 & \text{if } I=1 \\ \text{Uniform } [0,2] & \text{if } I=0 \end{cases}$$

Z is neither a discrete nor a continuous RV.

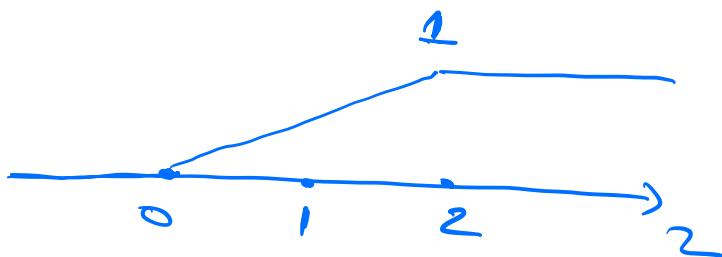
$$\begin{aligned} F_Z(z) &= P(Z \leq z) \\ &= P(Z \leq z | I=1)P(I=1) + P(Z \leq z | I=0) \\ &\quad P(I=0) \\ &= P F_X(z) + (1-p) F_Y(z) \\ &\quad (P = P(I=1)) \end{aligned}$$

Z is a mixed random variable,

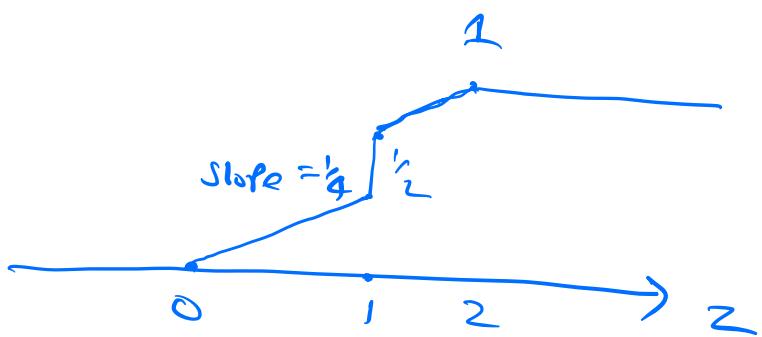
$$F_x(z)$$



$$F_y(z)$$



$$\text{For } p = \frac{1}{2} \quad f_z(2) =$$



Application of Probability in Information

Leakage,

The "information content" of a message depends on the degree to which the content of the message is surprising. If a highly likely event occurs, the message carries very little information. If a highly

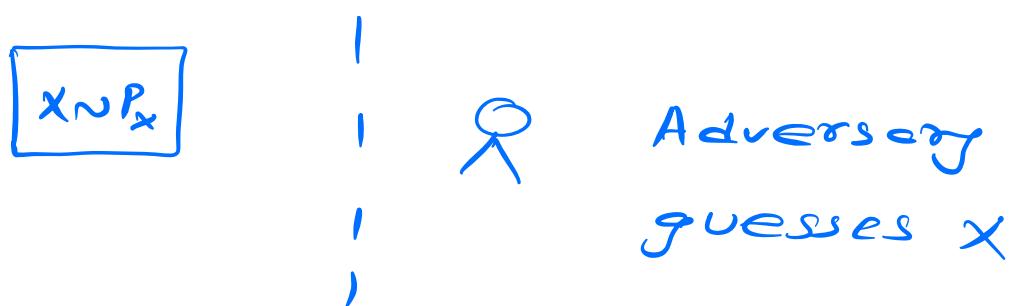
unlikely event occurs, the message is more informative.

The information content in an event with probability $p = \log \frac{1}{p}$.

[all logarithms are to the base 2]

Vulnerability: Suppose $x \sim p_x$. The vulnerability of x is given by

$$v(x) = \max_{x \in X} p_x(x).$$



$v(x)$ is the worst-case probability that an adversary could guess the value of x correctly.

$\log \frac{1}{v(x)}$ can be viewed as information measure

$$\text{min-entropy of } x - H_\infty(x) = \log \frac{1}{v(x)}$$

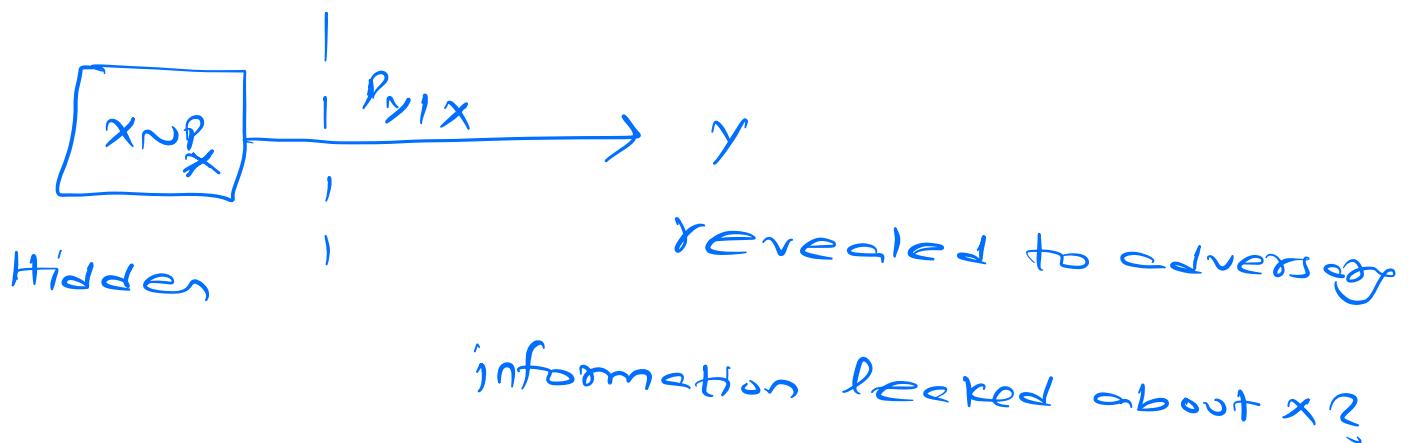
$$= \log \frac{1}{\max_x p_x(x)}.$$

- Given two RVS, $(x, y) \sim P_{xy}$, the conditional vulnerability

$$v(x|y) = \sum_{y \in Y} p_y(y) v(x|y=y)$$

$$= \sum_{y \in Y} p_y(y) \max_x p_{x|y}(x|y).$$

$$H_\infty(x|y) = \log \frac{1}{v(x|y)}.$$



Information leaked about x from y

= initial uncertainty in x

-

remaining uncertainty in x after observing y

$$\text{Leakage } \angle_{P_X}(x \rightarrow y) = H_\infty(x) - H_\infty(x|y)$$

$$= \log \frac{v(x|y)}{v(x)}.$$

If P_X is unknown, the worst-case leakage is of interest.

Theorem. $\max_{P_X} \angle_{P_X}(x \rightarrow y) = \angle_{P_U}(x \rightarrow y)$

$$= \sum_y \max_x P_{Y|X}(y|x).$$

(sum of maximum of the columns of $P_{Y|X}$ matrix)

Proof. $\angle_{P_X}(x \rightarrow y)$

$$= \frac{v(x|y)}{v(x)}$$

$$= \frac{\sum_y P_Y(y) \max_x P_{X|Y}(x|y)}{\max_{x'} P_X(x')}$$

$$= \sum_y P_Y(y) \max_x \frac{P_{Y|X}(y|x) P_X(x)}{P_Y(y)} / \max_{x'} P_X(x')$$

$$\begin{aligned}
 &= \frac{\sum_{y} \max_x P_{Y|X}(y|x) P_X(x)}{\max_{x'} P_X(x')} \\
 &\leq \frac{\sum_y \max_x P_{Y|X}(y|x) \cdot \cancel{\max_{x'} P_X(x')}}{\cancel{\max_{x'} P_X(x')}} \\
 &= \sum_y L_{P_0}(x \rightarrow y),
 \end{aligned}$$

where $P_0(x) = \frac{1}{|x|}$, $\forall x \in X$.

Remark. This is an operationally motivated leakage measure and satisfies all the axioms of a leakage measure. Interestingly, mutual information $I(X; Y) = H(X) - H(X|Y)$ does not satisfy all.

Q) Let (X, Y) have the joint PDF

$$f_{XY}(x, y) = \begin{cases} xy & (x, y) \in [0, 1]^2 \\ 0 & \text{otherwise,} \end{cases}$$

Let $Z = x(1+y)$, $W = x^2$.

$$g_1(x,y) = x(1+y) \quad g_2(x,y) = x^2.$$

$$x = h_1(z \geq \omega) = \sqrt{\omega}, \quad y = h_2(z \geq \omega) = \frac{z}{\sqrt{\omega}} - 1.$$

$$\frac{\partial g_1}{\partial x} = 1+y \quad - \frac{\partial g_1}{\partial y} = x$$

$$\frac{\partial g_2}{\partial x} = 2x \quad - \quad \frac{\partial g_2}{\partial y} = 0.$$

$$J = \begin{vmatrix} 1+y & x \\ 2x & 0 \end{vmatrix} = 2x^2 = 2\omega$$

$$f_{Z \geq \omega}(z \geq \omega) = \frac{f_{XY}(\sqrt{\omega}, \frac{z}{\sqrt{\omega}} - 1)}{2\omega}$$

$$= \frac{\sqrt{\omega} + \frac{z}{\sqrt{\omega}} - 1}{2\omega},$$

for $(z \geq \omega)$ s.t. $0 \leq \omega \leq 1$, $\sqrt{\omega} \leq z \leq 2\sqrt{\omega}$.

0 otherwise,