

Statistical Inference (2nd Edition)

Chapter 1: Probability Theory

Statistics builds upon the foundation of prob theory.

outline some of the basic ideas of prob theory fundamental to statistics

- Probability theory builds upon set theory:

1.1: Set Theory

Def (sample space): the set, S , of all possible outcomes of a particular experiment

- can be countable or uncountable; dictates the way in which probs can be assigned

countable: elements of S can be put into a 1-1 correspondence w/ subset of integers

Recall: a set can be finite or infinite

→ an infinite set is either countable or uncountable

→ a countable set is either finite or countably infinite

Def (count): any collection of possible outcomes of an experiment, i.e. any subset of S (including S itself)

constraint: $A \subset B \Leftrightarrow A \in \sigma_B$

equality: $A \subset B \Leftrightarrow A \in B \text{ and } B \in A$

- Given any two events (or sets) A and B , we have the following

<u>elementary set operations:</u> <ul style="list-style-type: none"> Union: $A \cup B = \{x : x \in A \text{ or } x \in B\}$ Intersection: $A \cap B = \{x : x \in A \text{ and } x \in B\}$ Complementation: $A^c = \{x : x \notin A\}$
--

Remark: \emptyset denotes the empty set

(+) can be contained; being careful, we can treat sets as numbers

Theorem (Properties of set operations)

For any three events, A , B , and C , defined on a sample space, S ,

a. Commutativity

$$A \cup B = B \cup A$$

$$A \cap B = B \cap A$$

b. Associativity

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

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

c. Distributive laws

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

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

d. De Morgan's Laws

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

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

the operations of union and intersection can be extended to infinite collections of sets as well:

Thm: If A_1, A_2, A_3, \dots is a collection of sets, all defined on a sample space, S , then

countable unions

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

$$= \{x \in S : \exists i \in \mathbb{N} \text{ s.t. } x \in A_i\}$$

$$= \lim_{n \rightarrow \infty} \bigcup_{i=1}^n A_i$$

$$\bigcap_{i=1}^{\infty} A_i = \{x : \forall i \in \mathbb{N} \text{ s.t. } x \in A_i\}$$

$$= \{x \in S : \forall i \in \mathbb{N} \text{ s.t. } x \in A_i\}$$

$$= \lim_{n \rightarrow \infty} \bigcap_{i=1}^n A_i$$

Remark: it is also possible to define unions and intersections over uncountable collections of sets:

corollary: if I is an index set (set of elements to be used as indices), then

$$\bigcup_{i \in I} A_i = \{x : \exists i \in I \text{ s.t. } x \in A_i\}$$

$$= \{x \in S : \exists i \in I \text{ s.t. } x \in A_i\}$$

$$= \lim_{n \rightarrow \infty} \bigcup_{i=1}^n A_i$$

$$\bigcap_{i \in I} A_i = \{x : \forall i \in I \text{ s.t. } x \in A_i\}$$

$$= \{x \in S : \forall i \in I \text{ s.t. } x \in A_i\}$$

uncountable unions & intersections

not being covered here

note: don't play major role in statistics; sometimes provide a useful mechanism as seen later (Ch-B).

. It can be verified that $\bigcup_{i=1}^{\infty} A_i$ and $\bigcap_{i=1}^{\infty} A_i$ are infinite collections of subsets too, i.e.

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

$$\text{Proof: } w \notin \bigcup_{i=1}^{\infty} A_i \Leftrightarrow \forall i, w \notin A_i \Leftrightarrow w \notin A_i^c \Leftrightarrow w \in \bigcap_{i=1}^{\infty} A_i^c.$$

Def: two events A and B are disjoint (mutually exclusive) if $A \cap B = \emptyset$

the events A_1, A_2, \dots are pairwise disjoint (mutually exclusive) if $A_i \cap A_j = \emptyset \quad \forall i \neq j$

. Disjoint sets are sets w/ no points in common

Def: if A_1, A_2, \dots are pairwise disjoint and $\bigcup_{i=1}^{\infty} A_i = S$, then the collection A_1, A_2, \dots forms a partition of S

. Generally, partitions are useful, enabling us to divide the sample space into small, non-overlapping pieces

1.2: Basics of Probability Theory

. the 'frequency of occurrence' of an outcome can be thought of as a probability

. we describe outcomes of an experiment probabilistically to analyze the experiment statistically

describe basics of prob. here

→ define w.r.t. to a mathematically simpler axiomatic approach (instead of in terms of frequency)

Recall:

Def (real number):

any number that can be placed on the number line

Numbers: Real = Rational \cup Irrational

Def (rational Number): Any number that can be written as a fraction

$$\frac{p}{q} \text{ where } \begin{cases} p \text{ and } q \text{ are integers} \\ q \neq 0 \end{cases}$$

e.g. integers, fractions, terminating & repeating decimals

Def (irrational Number): real number that violates (1) and/or (2)

s.t. its decimal expansion goes on forever,

e.g. π , e, square roots of non-perfect squares ($\sqrt{2}$)

Irrational (\mathbb{Q}') Rational (\mathbb{Q})

integers (\mathbb{Z}) $-1, -1, 0, 1, 2, \dots$

rationals (\mathbb{Q}) $0, 1, 2, \dots$

irrationals (\mathbb{Q}') $\pi, e, \sqrt{2}, \dots$

$$\mathbb{R} = \bigcup_{n=1}^{\infty} (-n, n) \cup \mathbb{Z}, \quad n \in \mathbb{N}$$

$$\mathbb{R} = \bigcup_{n \in \mathbb{Z}} (-2, 2)$$

Notes (Countable Unions):

Conceptual Layer	What It Means	How It Practically Used
Set Theory	$x \in \bigcup_{i=1}^{\infty} A_i \Leftrightarrow \exists i, x \in A_i$	logical membership
Prob. Theory	union = "at least one event occurs"	aggregating likelihood

Notes (Countable Intersections):

Conceptual Layer	What It Means	Prob. Theory
Conceptual	common elements across sets	Joint occurrence of events
Abstract	$x \in \bigcap_{i=1}^{\infty} A_i \Leftrightarrow \forall i, x \in A_i$	$P(A_1 \cap A_2 \cap \dots) = \prod_i P(A_i)$: prob. all occur

+ concerned not w/ prob. interpretations but instead define prob's by a function satisfying the axioms
- discuss subjective prob. interpretation as a belief in the chance of an event occurring

(2.1) Axiomatic Foundations

- For each event A in Ω , we want to associate w/ it a number between 0 and 1 that we call $P(A)$

+ not as simple to define as follows:

for each $A \in \Omega$, we define $P(A)$ = prob. that A occurs

→ technical difficulties arise often more relevant to probabilists than statisticians

→ however, a firm understanding of statistics requires us at least a passing familiarity with the following:

Def:	A collection of subsets of S is called a <u>sigma algebra</u> (or <u>Borel field</u>), denoted by \mathcal{B} , if it satisfies the following properties:
a.	$\emptyset \in \mathcal{B}$
b.	If $A \in \mathcal{B}$, then $A^c \in \mathcal{B}$ (closed under complementation)
c.	If $A_1, A_2, \dots \in \mathcal{B}$, then $\bigcup_{i=1}^{\infty} A_i \in \mathcal{B}$ (closed under countable unions)

Remarks:

\mathcal{B} is also closed under countable intersections (by deMorgan's)

$$i.e. \text{ if } A_1, A_2, \dots \in \mathcal{B} \text{ then } \bigcap_{i=1}^{\infty} A_i \in \mathcal{B}$$

The collection of two sets \emptyset, S is a sigma algebra, usually called the trivial sigma algebra.

Ex(C): If S is a finite or countable, then these technicalities rarely do

not arise, for we define, for a given sample space S ,

$$\mathcal{B} = \text{all subsets of } S, \text{ including } S \text{ & } \emptyset$$

Remarks: If S has n elements, there are 2^n sets in \mathcal{B} :

$$i.e. S = \{1, 2, \dots, n\}, \text{ then } |\mathcal{B}| = 2^n = 8: \begin{cases} \{1\}, \{2\}, \{3\}, \{4\} \\ \{1, 2\}, \{1, 3\}, \{1, 4\} \\ \{2, 3\}, \{2, 4\} \\ \{3, 4\} \\ \emptyset \end{cases}$$

Ex(C): Let $S = (-\infty, \infty)$ on the real line, then choose \mathcal{B} to contain all sets of the form

$$[a_1, b_1], [a_2, b_2], (a_3, b_3) \text{ and } [a, b] \quad \forall a, b \in \mathbb{R}.$$

Remark: From the properties of \mathcal{B} , it follows that \mathcal{B} contains all sets that can be formed by taking (possibly countably infinite) unions and intersections of sets of the above varieties.

We are now in a position to define a probability function.

Def:	Given a sample space S and an associated sigma algebra \mathcal{B} , a probability function is a function P w/ domain \mathcal{B} that
Satisfies:	$\left\{ \begin{array}{l} P(\emptyset) = 0 \text{ for all } \emptyset \in \mathcal{B} \\ \text{if } A_1, A_2, \dots \text{ are pairwise disjoint,} \\ \text{then } P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i) \end{array} \right.$

Remarks: These three properties are called (Kolmogorov's) Axioms of Probability (A)

- Any function P that satisfies (A) is called a prob. function

. we assume finite additivity holds, i.e. if $A \in \mathcal{B}$ and $B \in \mathcal{B}$ are disjoint, then

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

For any sample space, many different prob. functions can be defined

→ which ones reflect what is likely to be observed in a particular experiment
is still to be discussed

Ex(Defining Probabilities-I): Tossing a fair coin: $S = \{\text{H, T}\}$.

$$P(H) = P(T) = \frac{1}{2}$$

Tossing an unfair coin: $\begin{cases} P(H) = q, \\ P(T) = p, \end{cases}$ for some $q, p \in \mathbb{R}$

Tossing two coins in a row:
 $S = \{\text{HH, HT, TH, TT}\}$

$$P(HH) = p_1, \dots, P(TT) = q_1$$

$$\text{where } \begin{cases} p_1 = p_2 = p \\ q_1 = q_2 = q \end{cases}$$

We need general methods of defining prob. functions that we know will always satisfy Kolmogorov's Axioms

→ the following gives a common method of defining a legitimate prob. function

Thm: Let $S = S_1, \dots, S_n$ be a finite set. Let \mathcal{B} be any sigma algebra of subsets of S . Let P_1, \dots, P_n be nonnegative numbers that sum to 1.

For any $A \in \mathcal{B}$, define $P(A) = \sum_{S_i \in A} P_i$

Corollary: this remains true if $S = S_1, S_2, \dots$ is a countable set.

The physical reality of the experiment might dictate the prob. assignment, e.g.

Ex(Defining Probabilities-II): In a dart throwing example,

$$P(\text{Scoring } i \text{ points}) = \frac{\text{area on ring } i}{\text{area of dartboard}} = \frac{(6-i)^2 - (5-i)^2}{5^2}, i \in \{1, \dots, 5\}$$

→ the sum of the disjoint regions equals the area of the dart board.

REMARKS:

Q: How to show that sets are elements of \mathcal{B} ?

Q: Verify that the set in question can be obtained from known elements of \mathcal{B} via countably many set operations

→ in particular, sometimes useful to use fact that any real number can be constructed as the limit of a sequence of rational numbers

Let $S = (-\infty, a)$ where $a \in \mathbb{Q} \Leftrightarrow \exists x \in \mathbb{Q} \text{ s.t. } a = x$.

$$\begin{array}{ccccccc} -\infty & \xrightarrow{\text{+1}} & x & \xrightarrow{\text{+1}} & x+1 & \xrightarrow{\text{+1}} & \infty \\ \downarrow & & \downarrow & & \downarrow & & \downarrow \\ -\infty & & a & & a+1 & & \infty \end{array}$$

Q: Let $a \in \mathbb{Q}$. I will show that $\{a\} \in \mathcal{B}$.

For $n \in \mathbb{N}$, $R = \mathbb{Q} \cap [a, a+1]$.

$$\text{Then, } (-\infty, a] = \lim_{n \rightarrow \infty} \bigcup_{i=1}^n (-\infty, a - \frac{1}{i}] \in \mathcal{B}.$$

$$\text{Since } (-\infty, a] \subset \{a\}, \quad \{a\} = (-\infty, a] \setminus (-\infty, a) \in \mathcal{B}.$$

(ii) Let $a \in \mathbb{R}$, $n \in \mathbb{N}$

Pick rationals p_1, p_2, \dots (decreasing sequence of rationals converging to a form above), i.e. $p_i > a \forall i$.

$$\text{Then, } (-\infty, a] = \lim_{n \rightarrow \infty} \bigcup_{i=1}^n (-\infty, p_i] \in \mathcal{B}.$$

Pick rationals r_1, r_2, \dots (increasing sequence of rationals converging to a form below), i.e. $r_i < a \forall i$.

$$\text{Then, } (-\infty, a) = \lim_{n \rightarrow \infty} \bigcup_{i=1}^n (-\infty, r_i) \in \mathcal{B}.$$

$$\text{Thus, } \{a\} = (-\infty, a] \setminus (-\infty, a) \in \mathcal{B}.$$

1.2.2: THE CALCULUS OF PROBABILITIES

* From the Axioms of Probability (AOP), we can build up many properties of the probability function, some of which are quite helpful in the calculation of more complicated probabilities.

THEM: IF P IS A prob. function and A is any set in \mathcal{B} , then

- a. $P(\emptyset) = 0$
- b. $P(A) \leq 1$
- c. $P(A^c) = 1 - P(A)$

REMARK: Fairly self-evident; applied to a single event

THEM: IF P IS A prob. function and A and B are any sets in \mathcal{B} , then

- a. $P(A \cap A^c) = P(\emptyset) = 0$
- b. $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- c. (If $A \subset B$, then $P(A) \leq P(B)$)

Corollary: $P(A \cap B) = P(A) + P(B) - I$ (Bonferroni's Inequality) (a)

REMARKS: i) allows us to bound the prob. of a simultaneous event (intersection) in terms of the prob.s of the individual events
+ unless the probs. of the individual events are sufficiently large, the (counterc)bound is a useful (but conservative) number

THEM: IF P IS A probability function, then

a. $P(C_i) = \sum_{i=1}^n P(A_i C_i)$ for any partition C_1, C_2, \dots

b. $P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$ for any sets A_1, A_2, \dots (Boole's Inequality)

$$\text{PROOF: } P(B) = P(A) + P(B \cap A^c) \\ \Rightarrow P(B) \geq P(A).$$



$$\text{PROOF: } \text{By (b), } P(A \cap B) \\ = P(A) + P(B) - P(A \cup B) \\ = P(A \cap B) \geq P(A) + P(B) - I$$

REMARKS: - A partition (C_1, C_2, \dots) of S means C_1, C_2, \dots are pairwise disjoint, and $\bigcup_{i=1}^n C_i = S$

- Boole's inequality is a more general version of the Bonferroni inequality, i.e.

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i) - (n-1)I$$



$$\text{PROOF: } P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i) \quad \text{Applying Boole's Inequality} \\ 1 - P\left(\bigcap_{i=1}^n A_i^c\right) \leq 1 - \sum_{i=1}^n P(A_i^c), \quad \bigcap_{i=1}^n A_i^c = (A_1 \cap A_2 \cap \dots \cap A_n)^c = 1 - P(A_1 \cap A_2 \cap \dots \cap A_n) \\ P\left(\bigcap_{i=1}^n A_i\right) \geq \sum_{i=1}^n P(A_i) - (n-1)I$$

1.2.3: COUNTING

- elementary task of counting can become sophisticated for a statistician

- most often, methods of counting are used to construct prob assignments on finite sample spaces

• Break down counting problems into a series of simple tasks that are easy to count, and employ known rules of combining tasks

Fundamental Rule of Counting: If a job consists of k separate tasks, the i^{th} of which can be done in n_i ways, $i=1, \dots, k$, then the entire job can be done in $n_1 n_2 \dots n_k$ ways

REMARK: although the FRC is a reasonable place to start, in applications there are usually more aspects of a problem to consider, i.e.

counting
with/without replacement
uniform/not order matters

AS AN aside, we first discuss helpful notation below:

DEF: For a positive integer n , $n!$ (n factorial) is the product of all of the positive integers less than or equal to n , i.e.

$$n! = n(n-1)(n-2)\dots(3)(2)(1) \quad \text{where } 0! = 1.$$

REMARK: # of distinct ways the original n objects can be shuffled

DEF: For nonnegative integers n and r , where $n \geq r$, we define n permutations: P_r^n

$$P_r^n = n(n-1)\dots(n-r+1) = \frac{n!}{(n-r)!}$$

REMARK: an ordered arrangement of r distinct objects out of n total

DEF: For nonnegative integers n and r , where $n \geq r$, we define n choose r : $\binom{n}{r}$ as

$$\binom{n}{r} = \frac{n!}{r!(n-r)!}$$

REMARK: the # of (unordered) combinations of n objects chosen r at a time

EX (lottery ticket): From the numbers 1, 2, ..., 44 a person may pick any six for their ticket

a) ordered w/o replacement: $44 \times 43 \times 42 \times 41 \times 40 \times 39 = \frac{44!}{6!} = 3,028,512,440$

b) ordered w/r replacement: $44 \times 44 \times 44 \dots 44 = 44^6 = 7,562,313,856$

c) unordered w/o replacement

- must divide out the redundant orderings $\frac{44!}{6!}$ $\Rightarrow \frac{44!}{6! \cdot 38!} = \frac{44!}{6! \cdot (44-6)!} = \binom{44}{6} = 7,059,052$
e.g. how six numbers arranged in 6! ways

d) unordered w/r replacement

- common mistake (but too small) is $\frac{44^6}{6!}$
- consider placing 6 markers in the 44 numbers, i.e.
count all arrangements of 44 values (44 bins \rightarrow 45 slots, disregard two ends)
 \Rightarrow there are 43! ways that can be arranged $43!$ ways,
enumerating the redundant orderings, i.e.

$$\frac{43!}{6!} = \frac{(44+6-1)!}{6!} = 13,983,816$$

PROOF: Let $k=2$. The first task can be done in n_1 ways, and for each of these ways, we have n_2 choices for the second task.

thus, we can do the job in

$$(n_1) + (n_2) + \dots + (n_k) = n_1 n_2 \dots n_k$$

ways, establishing the theorem for $k=2$.

Visually: Sampling ordered lottery tick w/o replacement:

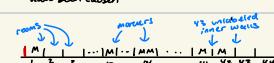
$$\begin{array}{ccccccccc} \text{First \#} & \text{Second \#} & \text{Third \#} & & & & & & \\ \text{contests: } & 44 & \times & 43 & \times & 42 & \times & \dots & \end{array}$$

REMARKING

- keep all the 44 numbers in a fixed order and create $\binom{44}{6}$ "rows", one for each number

+ rows are separated by bars (w/ 43 total)

+ place the 6 markers in 43 rows where the numbers have been chosen



GIVEN: $\binom{44}{6}$ rows
43 inner walls
44 positions

ABOVE: $\binom{44}{6}$ (10, 11, 12, 13, 14, 15) were chosen

* key obs: # tickets \equiv

unique ways to place 6 markers to 44 bins

= # ways to arrange 6 markers w/ 43 inner walls

$$= \binom{6+43}{6}$$

possible methods of counting (summary)

Number of possible arrangements of size r from n objects

	with replacement	without replacement
ordered	$P_r^n = \frac{n^r}{(n-r)!}$	n^r
unordered	$C_r^n = \frac{n!}{r!(n-r)!} = \binom{n}{r}$	$\binom{n-r}{r-1}$

Remark: • $\binom{n}{r}$ are referred to as binomial coefficients

$$\bullet C_r^n = \frac{P_r^n}{r!} = \frac{n^r}{r!(n-r)!}; \text{ dividing by redundant orderings}$$

Prob. since $r! \geq 1$, it follows that $C_r^n \leq P_r^n$

1.2.4: Enumerating options

* the counting techniques of the previous section are useful when the sample space S is a finite set and all of the outcomes are equally likely

→ then prob's of events can be calculated by simply counting the number of outcomes in the event

Prob. Suppose $S = \{1, 2, \dots, n\}$ is a finite sample space.

Assume that all outcomes are equally likely s.t.

$$P(E \in S) = \frac{1}{n} \quad \forall E \subseteq S$$

Then, using AOP(3),

$$P(A) = \sum_{E \subseteq A} P(E \in S) = \sum_{E \subseteq A} \frac{1}{n} = \frac{\# \text{ of elements of } A}{n} = \frac{\# \text{ of elements in } S}{n}$$

Remark: For large sample spaces, the counting techniques might be used to calculate both the numerator and denominator of this expression

Ex (Poker): suppose the events do not depend on order, so we use the unordered outcomes

Q: How many 5-card hands can be chosen from a 52-card deck

$$\bullet \binom{52}{5} \Rightarrow P(\{\text{random hand}\}) = 1/\binom{52}{5} = 2,598,960 \text{ (#)}$$

$$\bullet P(\{\text{4 aces}\}) = \frac{48}{\binom{52}{5}} \text{ where } 48 \text{-# different ways of specifying fifth card}$$

Q: $P(\{\text{4 of a kind}\}) = \frac{13 \cdot 48}{\binom{52}{5}}$ where 48-as above

13 - which denomination there will be four of

30 2H 2S 2H 2S

$$\bullet P(\{\text{exactly one pair}\}) = \frac{1}{\binom{52}{5}} \cdot 13 \cdot \binom{4}{2} \cdot \binom{48}{3} \cdot 4^3$$

where 13-as above $\binom{4}{2} = 6$

$\binom{48}{3}$ - # ways to specify the two cards from that denomination

4^3 - # of ways of specifying the other three denominations

4^3 - # of ways of specifying the other three cards from those denominations

- When Sampling w/o replacement, if we want to calculate the prob. of an event that does not depend on the order, we can either use the ordered or unordered sample space

→ Each outcome in the unordered sample space corresponds to $r!$ outcomes in the ordered sample space

. When Sampling w/ replacement, each outcome in the unordered sample space corresponds to some outcomes in the ordered sample space, but the number of outcomes differs

→ the formula for the number of outcomes in the unordered sample space is useful for enumerating the outcomes, but ordered outcomes must be counted to correctly calculate prob's, e.g.

Ex (Sampling w/o replacement)

"uniform sampling": $r=2$ items from $n=3$ items w/o replacement

ordered	(1,1)	(2,1)	(3,1)	(1,2)(2,1)	(1,3)(2,1)	(2,3)(2,2)
ordered	$\frac{1}{3} \cdot \frac{1}{2}$	$\frac{2}{3} \cdot \frac{1}{2}$	$\frac{3}{3} \cdot \frac{1}{2}$	$\frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{1}$	$\frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{1}$	$\frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{1}$
Prob.	1/6	1/6	1/6	1/6	1/6	1/6

1.3: Conditional Probability and Independence

In many instances, we are in a position to update the sample space based on new info

→ In such cases, we want to be able to update prob. calculations, i.e. to calculate conditional prob's

Def: if A and B are events in S , and $P(B) > 0$, then the conditional probability (if A given B) is

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



Remarks: Equivalently, $P(A \text{ and } B \text{ true}) = \frac{P(\text{both } A \text{ and } B \text{ true})}{P(B \text{ true})} \cdot P(B)$; relative prob. of overlap inside B

where $(*)$ is "renormalizing" under condition into B

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

Joint prob. = marginal + conditional

Prob. Generally, $P(A \cap B \cap C) = P(A) \times P(B \cap C|A) = P(A) \times P(B|A) \times P(C|B, A)$

→ Hence, Chain formula: $= P(A_1 \cap \dots \cap A_n) = \prod_{i=1}^n P(A_i | A_1 \cap \dots \cap A_{i-1})$

(*)

$$\text{Ex (Four Aces): } P(\{\text{4 aces}\}) = \frac{1}{\binom{52}{4}} = \frac{(\text{first ace})(\text{second ace})(\text{third ace})(\text{fourth ace})}{\binom{52}{4}} = \frac{P(\text{1st ace}) \cdot P(\text{2nd ace}|\text{1st ace}) \cdot P(\text{3rd ace}|\text{1st, 2nd ace}) \cdot P(\text{4th ace}|\text{1st, 2nd, 3rd ace})}{\binom{52}{4}}$$

s.t. (*) updates the sample space after each draw of a card

Prop. If $P(A) > 0$, then $P(B|A) = P(A)P(B|A)$
 $= P(B)P(A|B)$ joint prob. of both happening

Conseq.: hence, $P(A|B) = \frac{P(A)P(B|A)}{P(B)}$ for $P(B) > 0$, $A, B \in \mathcal{B}$.

Remarks:

- * conditional probabilities provide the means to quantify how we may have different assessments of the same event (due to different info we use or have access to)
- * what event/info we choose to condition on can affect the answer
- * conditional prob's provide the mathematical machinery of Bayesian statistics which bases the inference on the posterior distribution, i.e. cond. prob. of quantities of interest conditionally given observed data
- * cond. Prob's are obtained usually via the Bayes' formula

Ex (Three Prisoners): three prisoners A, B, and C one of them chosen (random) to be pardoned.

A asks the warden, who is supposed to keep secret:
 - which among B and C will be executed?
 to which the warden responds
 - B is to be executed.

Warden's thinking: the prob. that either A or B or C gets pardoned is 1/2.
 Between B and C, at least one is executed.
 So, "2nd" guess A no new information on his life.

A's thinking: since either C or I gets the pardon, my chance of being alive has gone up to 1/2

Let A, B, C be pardon events, i.e.
 $P(A) = P(B) = P(C) = 1/2$.

IF Prisoner is pardoned	Then, warden tells A	Prob
A	$\rightarrow B$ dies	1/1
	$\rightarrow C$ dies	1/2
B	$\rightarrow C$ dies	1
C	$\rightarrow B$ dies	1

Let us think what warden tells A that B dies

The warden's reasoning is as follows:

$$P(A|W) = \frac{P(A \cap W)}{P(W)} = \frac{P(A)P(W|A)}{P(W)} = \frac{(1/2)(1/2)}{(\frac{1}{2})(\frac{1}{2}) + (\frac{1}{2})(1)} = \frac{1/4}{1/2} = \frac{1}{2}.$$

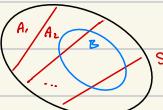
A's reasoning comes from

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B|A)}{1 - P(B)} = \frac{1/2}{1 - 1/2} = \frac{1}{2}.$$

Prop. Let A_1, A_2, \dots be a partition of S , i.e. $S = \bigcup_{i=1}^{\infty} A_i$

$$\text{then } P(B) = \sum_i P(B \cap A_i) = \sum_i P(B|A_i)P(A_i)$$

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



Remarks: (a) is the law of total probability
 - turning around conditional prob's

Corollary: $\forall i: P(A_i|B) \propto P(B|A_i)P(A_i)$

Remarks: \propto : "is proportional to" (up to multiplying constant)

in Bayesian statistics

$P(A_1|B) \propto P(B|A_1)P(A_1)$
 || same likelihood prior dist.

Posterior dist.

Proof: We want to be able to convey that an event B has no effect on A by insisting that

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

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

Moreover,

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

Def: two events are statistically independent if $P(A \cap B) = P(A)P(B)$

Remarks: does not require $P(A) > 0$ or $P(B) > 0$
 A & B independent does not mean $A \cap B = \emptyset$ (disjoint def.)

Exercise Rolling: $S = \{1, \dots, 6\}$. Here, $\bar{S} = \{1, \dots, 6\}^4 \in \mathbb{I}^{12}$

$$Q: P(\text{at least one 6 in 4 rolls})$$

Assuming outcome of rolls are independent

$$= 1 - P(\text{no 6 in 4 rolls})$$

$$= 1 - \left(\frac{5}{6}\right)^4$$

$$= 1 - (1 - 5/6)^4 = 1 - (1/6)^4 = 0.58$$

Independence of A and B implies independence of complements, i.e.

Thm: If $A \perp\!\!\!\perp B$, then
 i) $A \perp\!\!\!\perp B^c$
 ii) $A^c \perp\!\!\!\perp B$
 iii) $A^c \perp\!\!\!\perp B^c$

Proof (i): $P(A \cap B^c) = P(A) - P(A \cap B)$, and
 $= P(A) - P(A)P(B)$, and
 $= P(A)(1 - P(B))$
 $= P(A)P(B^c)$.

Def: A collection of events A_1, \dots, A_n are mutually independent, if from any subcollection A_{i_1}, \dots, A_{i_k} , we have

$$P(A_{i_1} \cap \dots \cap A_{i_k}) = \prod_{j=1}^k P(A_{i_j})$$

Remarks: mutual independence much stronger than pairwise ind.

Corollary: pairwise ind. does not imply mutual ind., i.e.
given 3 events A, B, C s.t.

$$\begin{aligned} P(A)P(B) &= P(A \cap B) \\ P(A)P(C) &= P(A \cap C) \\ P(B)P(C) &= P(B \cap C) \end{aligned}$$

conversely, in general,

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

Ex (Letters): Let $S = \{abc, acb, bac, bca, cab, cba\}$ s.t. each element in S is equally ($1/6$) probable.

Let $A_i = \{\text{3rd place in the triple is occupied by } a_i\}$.

Then, $P(A_1) = P(A_2) = P(A_3) = 1/3$.

$$\begin{aligned} P(A_1 \cap A_2) &= P(A_1 \cap A_2) = P(A_1 \cap A_2) \\ &= P(A_1 \cap A_2 \cap A_3) = P(\{\text{all } a\}) = 1/6. \end{aligned}$$

Thus, A_1, A_2, A_3 are pairwise independent, but

$$P(A_1 \cap A_2 \cap A_3) \neq P(A_1)P(A_2)P(A_3)$$

and therefore, not mutually independent.

Ex (Three coin tosses): Assume three coin tosses, each of which produces H or T.

$$S = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\} \text{ s.t. all 8 outcomes are equally probable.}$$

Let $A_i = \{\text{ith toss is } H\}$, $i=1, 2, 3$.

$$\text{Then, } P(A_i) = P(\{\text{HHH, HHT, HTH, HTT}\}) = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = \frac{1}{2}.$$

$$P(A_1 \cap A_2) = P(\{\text{HHHH, HHTH, HTTH}\}) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}.$$

$$P(A_1 \cap A_2 \cap A_3) = P(\{\text{HHHH}\}) = \frac{1}{8}.$$

Likewise, $P(A_i) = 1/2 \quad \forall i=1, 2, 3$

$$P(A_1 \cap A_2) = 1/4 \quad \forall i, j \in \{1, 2, 3\}$$

$$\begin{cases} P(A_1 \cap A_2 \cap A_3) = P(A_1)P(A_2)P(A_3) \\ P(A_1 \cap A_2 \cap A_3) = P(A_1)P(A_2) \quad \forall i, j, k \end{cases}$$

Hence, A_1, A_2, A_3 are pairwise and mutually independent.

1.4 Random Variables

- In many experiments, it is easier to deal w/ a summary variable than w/ the original prob. structure
- # In defining the quantity X , we have a mapping (a function) from the original sample space to a new sample space, usually a set of real numbers.

Def: A random variable is a function from a sample space S into the real numbers.

Ex (RV): Experiments | Random Variables

i) Toss 2 dice
 $S = \{1, \dots, 6\} \times \{1, \dots, 6\}$

Suppose $S = \{(i_1, i_2)\}$, then $X = X(i) = i_1 + i_2$

ii) Toss a coin 2 times
 $S = \{H, T\}^2$

Suppose $S = \{c_1, \dots, c_{10}\}$, then $X = X(c) = \sum_{i=1}^2 c_i + 1$

note: indicator function $I(A) = \begin{cases} 1, & \text{if } A \\ 0, & \text{otherwise} \end{cases}$

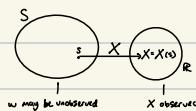
iii) Apply different amt of fertilizer to flower plants

$X = \text{yield/acre}$

- In defining a RV, we have also defined a new sample space (the range of the RV)
 → we now must formally check that our prob. function, which is derived on the original sample space, can be used for the RV.

Prop: suppose we have a sample space $S = \{s_1, \dots, s_n\}$ w/ prob. function P
 and we define a RV X w/ range $\{x_1, \dots, x_m\}$. we can define a prob. function P_X on \mathbb{R} in the following way:

$$P_X(X=x) = P(\{s_i \in S : X(s_i) = x\})$$



Remarks: - we observe $X=x$, i.e. the outcome of the random experiment is on

$s \in S : X(s) = x$ (equivalence in words)

- RV's will always be denoted w/ uppercase letters and the realized values of the variable (its range) will be denoted by corresponding lowercase letters

Ex (Three coin tosses II): Let X : number of heads out of tossing coin 3 times

$$\begin{aligned} S &= \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\} \\ X(s) &= \begin{cases} 3 & \text{if } s = HHH \\ 2 & \text{if } s = HHT, HTH, THH \\ 1 & \text{if } s = HTT, THT, TTH \\ 0 & \text{if } s = TTT \end{cases} \end{aligned}$$

so, the range of (RV) X is:

	0	1	2	3
P_X	$1/8$	$3/8$	$3/8$	$1/8$

Ex (discr. of a RV): Let $X = \# \text{ heads after tossing coin } n \text{ times}$

$$\text{Then, } X = \# \text{ heads after tossing coin } n \text{ times} \Rightarrow \text{Range of } X = \{0, 1, 2, \dots, n\}. \text{ So } P(X=i) = \frac{\binom{n}{i}}{2^n} \quad \forall i \in \mathbb{N}.$$

- Previous definition for $P_A(\cdot)$ applies to settings w/ a finite S and \mathcal{X}

→ more straightforward; also the case when \mathcal{X} is countable

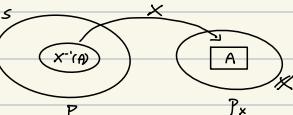
- If \mathcal{X} is uncountable, we define the induced prob. function P_A as follows:

Def: Suppose the range of random variable X is \mathcal{X} . For any subset $A \subset \mathcal{X}$, the probability distribution of X is a probability function, denoted P_A , on the sigma-algebra of subsets of \mathcal{X} s.t.

\forall such subset $A \subset \mathcal{X}$

$$P_A(X \in A) = P(X^{-1}(A)) = P(\{\omega \in S : X(\omega) \in A\})$$

Remarks: For complicated space (e.g. uncountable), we need to specify prob's on "simple sets" which generate the associated sigma-algebra



[5]: Distribution Functions

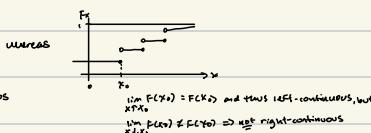
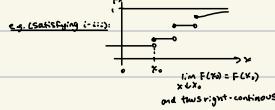
- With every RV X , we associate a function called the cumulative distribution function of X :

Def: The cumulative distribution function (cdf) of a RV X is defined by $F_X(x) = P(X \leq x) \text{ for all } x \in \mathcal{X}$

- Every cdf satisfies certain properties, some of which are obvious when we think of the definition of $F_X(x)$ in terms of prob.:

Thm: The function $F(x)$ is a cdf i.f.f. the following three conditions hold:

- $\lim_{x \rightarrow -\infty} F(x) = 0$ and $\lim_{x \rightarrow \infty} F(x) = 1$
- $F(x)$ is a nondecreasing function of x
- $F(x)$ is right-continuous, i.e. $\lim_{x \rightarrow x_0^+} F(x) = F(x_0)$



Ex (rolling for a head): Let $p := \text{prob. that a coin turns head}$

$X := \# \text{ independent coin tosses needed to get a head}$
 $\Leftrightarrow \# \text{ trials until first head}$
 $\Leftrightarrow \# \text{ trials on which first head occurs}$

Note: $X \sim \text{Geo}(p)$; Aside: Some authors define $\begin{cases} X^* := \# \text{ failures before first success} \\ X := \# \text{ trials on which first success occurs} \end{cases}$
 $\Rightarrow X^* = X - 1 \Rightarrow P(X=x) = P(X^*=x-1)$.

$$\text{Then, } X \in \{1, 2, \dots, \infty\}, \text{ so, } P(X=x) = \begin{cases} (1-p)^{x-1} p, & \forall x \in \{1, 2, \dots, \infty\} \\ 0, & \text{otherwise} \end{cases}$$

$$\Rightarrow F_X(x) = P(X \leq x), \quad \forall x \geq 0$$

$$= \sum_{k=1}^{x+1} (1-p)^{k-1} p$$

$$= p \sum_{k=0}^{x+1} (1-p)^k$$

$$= p \frac{1-(1-p)^{x+1}}{1-(1-p)}$$

$$= 1 - (1-p)^x, \quad x \geq 1$$

recall: (geometric series)

$$\text{Prop: } \forall x \neq 1 \in \mathbb{R}, \quad \sum_{k=0}^{x-1} r^k = \frac{1-r^x}{1-r}$$

Prop: $\forall |r| < 1$, i.e. $r \in (-1, 1)$, $\sum_{k=0}^{\infty} r^k = \frac{1}{1-r}$; otherwise, series diverges.

Proof (cfd): $\begin{cases} \lim_{x \rightarrow 0} F_X(x) = 0 \text{ since } F_X(0) = 0 \quad \forall x \geq 0 \\ \lim_{x \rightarrow \infty} F_X(x) = \lim_{x \rightarrow \infty} (1-p)^{x-1} + 1 = 1 \text{ for } x \in \mathbb{Z}. \end{cases}$

$\therefore \sum_{k=1}^{\infty} (1-p)^{k-1} p$ contains more positive terms as k increases

$\therefore \lim_{x \rightarrow \infty} F_X(x) = F_X(\infty)$.

\therefore $\lim_{x \rightarrow \infty} F_X(x) = F_X(\infty)$.

recall: (geometric series)

$$\text{Prop: } \forall x \neq 1 \in \mathbb{R}, \quad \sum_{k=0}^{x-1} r^k = \frac{1-r^x}{1-r}$$

Prop: $\forall |r| < 1$, i.e. $r \in (-1, 1)$, $\sum_{k=0}^{\infty} r^k = \frac{1}{1-r}$; otherwise, series diverges.

Ex (continuous cdf): Let $F_X(x) = \frac{1}{1+p} e^{-px}$ which satisfies

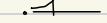
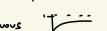
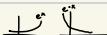
$$\text{1) } \lim_{x \rightarrow -\infty} F_X(x) = 0 \text{ since } \lim_{x \rightarrow -\infty} e^{-px} = 0.$$

$$\text{2) } \lim_{x \rightarrow \infty} F_X(x) = 1 \text{ since } \lim_{x \rightarrow \infty} e^{-px} = 0.$$

$$\text{3) } \frac{d}{dx} F_X(x) = \frac{d}{dx} \left(\frac{1}{1+p} e^{-px} \right) = \frac{-pe^{-px}}{(1+p)^2} > 0$$

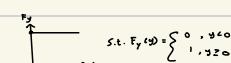
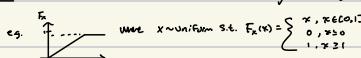
\Rightarrow increasing \Rightarrow nondecreasing

\therefore $F_X(x)$ is (right/left) continuous



Def: A RV X is continuous if $F_X(x)$ is a continuous function of x
A RV X is discrete if $F_X(x)$ is a step function of x .

Remark: cdf can also be a "mixture" of continuous segments and jumps



$$\text{s.t. } F_Y(y) = \begin{cases} 0, & y < 0 \\ 1, & y \geq 0 \end{cases}$$

$$\text{NOTICE: } F_2(z) = \frac{1}{2} F_X(z) + \frac{1}{2} F_Y(z)$$

still satisfies (i)-(iii) \Rightarrow valid cdf.

and in context,

$$\text{flip a fair coin } C = \{H, T\} \text{ s.t. } \begin{cases} \text{if } C=H, Z=X \\ \text{if } C=T, Z=Y \end{cases}$$

REMARK: motivates RV model construction from constituent RV mixtures

- we close w/ a theorem formally stating that F_A completely determines the prob. distn of a RV X
 \rightarrow this is true if $P(X \in A)$ is defined only for events $A \in \mathcal{B}$, the smallest sigma algebra containing all of the intervals of real numbers of the form (a, b) , $[a, b]$, $(a, b]$, and $[a, b)$
 \rightarrow we don't consider a larger class of events than above; avoids such pathological cases

DEF: the RVs X and Y are identically distributed, i.e.,
 $P(X \in A) = P(Y \in A) \quad \forall A \in \mathcal{B}$

REMARK: $X \stackrel{d}{=} Y$ does not necessarily imply $X=Y$

EX (identically distribute RV's): toss a fair coin 3 times.

then $X = \# \text{heads}$ and $Y = \# \text{tails}$

$$\Rightarrow X \stackrel{d}{=} Y \text{ since } P(X=i) = P(Y=i) \quad \forall i \in S.$$

Note that, $X \neq Y$ since for no sample points do we have $X(\omega) = Y(\omega)$.

ITEM: the following two statements are equivalent:

- the RVs X and Y are identically distributed
- $F_X(x) = F_Y(x) \quad \forall x \in \mathbb{R}$

1.6: Density and Mass Functions

- Associated w/ a RV X and its cdf F_X is another function, called either the probability density function (pdf) or probability mass function (pmf).
 \rightarrow respectively refer to the continuous & discrete case
 \rightarrow concerned w/ "point probabilities" of RV's

DEF: the probability mass function (pmf) of a discrete RV X is given by $f_X(x) = P(X=x) \text{ for all } x$

EX (discrete prob's): Let $X \sim Geom(p)$, $0 < p < 1$

$$\text{then } f_X(x) := P(X=x) = \begin{cases} (1-p)^{x-1} p, & x=1, 2, \dots \\ 0, & \text{otherwise} \end{cases}$$

$$\begin{aligned} \text{s.t. } P(A \subseteq X \subseteq B) &= \sum_{k=a}^b f_X(k) = \sum_{k=a}^b (1-p)^{k-1} p \\ P(X \in A) &= \sum_{k \in A} f_X(k) = f_X(A) \end{aligned} \quad (*)$$

Prop: the probability density function (pdf), $f_X(x)$, of a continuous RV X is the function that satisfies

$$f_X(x) = \int_x^{\infty} f_X(t) dt \text{ for all } x$$

REMARKS: if F_X is differentiable, then the pdf always exists
and it follows: $f_X(x) = \frac{d}{dx} F_X(x)$

\cdot F_X or f_X contain all info there is not the distn of the RV X , we denote

$$X \sim F_X \text{ or } X \sim f_X \text{ (equivalently)}$$

"is distributed as"

* can use either F_X or f_X to solve problems; try to choose simpler one

* we can instead understand the purpose of the pdf, i.e.
instead of summing (as seen in the discrete case \rightarrow),
substitute integrals for sums:

$$P(X \leq x) = F_X(x) = \int_{-\infty}^x f_X(t) dt$$

PROOF: $P(X=x) = 0$ whenever X -continuous RV

COROLLARY: $P(a \leq X \leq b) = P(a < X \leq b) = P(a \leq X < b) = P(a < X \leq b)$

PROOF (sketch): we must be careful in our definition of a pdf in the continuous case

if we naively try to calculate $P(X=x)$ for a continuous RV:

since $\{x\} = \{x\} \subset \{x - \epsilon < X \leq x + \epsilon\}$ for any $\epsilon > 0$,

$$P(X=x) \leq P(x - \epsilon < X \leq x) = F_X(x) - F_X(x - \epsilon) \text{ for any } \epsilon > 0.$$

Therefore, $0 \leq P(X=x) \leq \lim_{\epsilon \rightarrow 0} [F_X(x) - F_X(x - \epsilon)] = 0$ by continuity of F_X . \square



EX (logistic probabilities): recall the logistic cdf: $F_X(x) = \frac{e^x}{1+e^{-x}}$

$$f_X(x) = \frac{1}{1+e^{-x}}$$

which gives

$$f_X(x) = \frac{d}{dx} F_X(x) = \frac{e^{-x}}{(1+e^{-x})^2}$$

symmetric since $f_X(x) = f_X(-x)$, i.e.

$$f(x) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{e^{-x}e^{-x}}{(1+e^{-x})e^{-x}} = \frac{e^{-x}}{e^{-x}(1+e^{-x})} = \frac{e^{-x}}{(1+e^{-x})^2} = f(x).$$

then, $P(X \in [a, b]) = P(a < X < b) = F_X(b) - F_X(a)$

$$\begin{aligned} &= \int_{-\infty}^b f_X(x) dx - \int_{-\infty}^a f_X(x) dx \\ &= \int_a^b f_X(x) dx = \int_a^b \frac{e^{-x}}{(1+e^{-x})^2} dx \end{aligned}$$

- There are really only two requirements for a pdf (pmf), both of which are immediate consequences of the definition:

Theorem: A function $f_X(x)$ is a pdf (pmf) of a RV X .

i.f.f.

- a) $f_X(x) \geq 0$ for all x
- b) $\sum_x f_X(x) = 1$ (pmf)
or
 $\int_{-\infty}^{\infty} f_X(x) dx = 1$ (pdf)

Chapter 2: Transformations and Expectations

- often, if we are able to model a phenomenon in terms of a RV X w/ cdf $F_X(x)$, we also want to study the behavior of functions of X .

* ch 2 covers techniques that allow us to gain info abt functions of X
→ can range from very complete (dist. of these functions) to more vague (avg. behavior)

2.1: Distributions of Functions of RVs

Def: Let X - RV with cdf $F_X(x)$. Then, any function of X , denoted $g(X)$ is also a RV.

Remarks: often write $Y = g(X)$ to denote new RV of interest

Def: Since Y is a function of X , we can describe the probabilistic behavior, i.e.

$$P(Y \in A) = P(g(X) \in A)$$

- depending on the choice of g , it is sometimes possible to obtain a tractable expression for this prob.

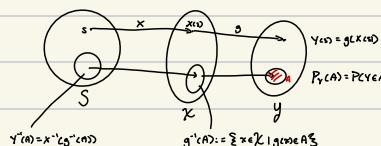
Prop: Formally, let $g: S \rightarrow Y$ for S : sample space. Then,

$$Y(s) = g(X(s)) = g \circ X(s)$$

and the following properties hold for the prob. distn of $Y(s)$:

By def., $\forall A \subset Y$ (prob. distn of $Y \Rightarrow$)

$$\begin{aligned} P(Y \in A) &= P(g(X) \in A) \\ &:= P(g^{-1}(A)) = P(X \in g^{-1}(A)) \\ &= P(g^{-1}(y)) \\ &= P(g^{-1}(y)) \\ &= P(Y \in A) \end{aligned}$$



Prop: Show that prob. distn for Y satisfies Kolmogorov's Axioms.

Let X - discrete RV $\Rightarrow S$: countable.
Then, the sample space for $Y = g(X)$:

$$Y = \{y : y = g(x), x \in S\}, \text{ also a countable set.}$$

Thus, Y is also a discrete RV S.t.

$$P_Y(y) = P(Y=y) = \sum_{x \in g^{-1}(y)} P(X=x) = \sum_{x \in g^{-1}(y)} f_X(x) \text{ for } y \in Y,$$

$$\text{and } f_Y(y) = 0 \text{ for } y \notin Y.$$

⇒ finding $f_Y(y)$ involves simply identifying $g^{-1}(y)$, for each $y \in Y$, and summing the appropriate probs. ■

Remark 3: • notation change from X to X .

• if there is only one x for which $g(x)=y$, then $g^{-1}(y)$ is the point set $\{x\}$ and we write $g^{-1}(y) = x$, i.e. $g^{-1}(2y) = \{x \in S : g(x)=y\}$

Continuous RVs

If X is a cont. RV, g a nice (cont.) function,

then $Y=g(X)$ is a cont. RV.

so,

$$\begin{aligned} P_Y(y) &= P(Y \leq y) \\ &= P(g(X) \leq y) \\ &= P(g^{-1}(\{y\}) \leq X) \\ &= \int_{g^{-1}(\{y\})} f_X(x) dx \\ &= \int_{g^{-1}(\{y\})} f_X(x) dx \end{aligned}$$

The set $\{x : g(x) \leq y\}$ may be difficult to identify
so P_Y may be hard to derive in general!

Prop: If g is (strictly) monotone, this gets simpler
s.t. either

- [] g is increasing, i.e. $g(x_1) > g(x_2)$ if $x_1 > x_2$
- [] g is decreasing, i.e. $g(x_1) < g(x_2)$ if $x_1 > x_2$

note: $g^{-1}(y) = \{x : g(x) = y\} \Rightarrow$ a singleton set.

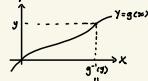
recall: monotonic function values only move in one direction (increasing/decreasing) as the input increases

Thm. Let X have the cdf $F_X(x)$,
 $K = \{x : F_X(x) > 0\} \subseteq \text{support of } f_X$
 $y = g(x) = \frac{1}{g'(x)} \cdot x \text{ for some } x \in K$

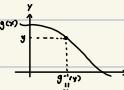
a) If $g: K \rightarrow Y$ is increasing, then $Y = g(X)$ is a RV taking values in Y with the cdf
 $F_Y(y) = F_X(g^{-1}(y)) \quad \forall y \in Y$

b) If $g: K \rightarrow Y$ is decreasing, then $Y = g(X)$ is a RV in Y with cdf
 $F_Y(y) = 1 - F_X(g^{-1}(y)) \quad \forall y \in Y$

Proof (a): If g is strictly increasing, where $Y = g(X)$, then
 $\{x : g(x) \leq y\} = \{x : x \leq g^{-1}(y)\}$
 $\Rightarrow F_Y(y) = \int_{-\infty}^{g^{-1}(y)} f_X(x) dx = F_X(g^{-1}(y))$



Proof (b): If g is strictly decreasing, where $Y = g(X)$, then
 $\{x : g(x) \leq y\} = \{x : x \geq g^{-1}(y)\}$
 $\Rightarrow F_Y(y) = \int_{g^{-1}(y)}^{\infty} f_X(x) dx = 1 - F_X(g^{-1}(y))$



Corollary: we can deduce the pdf from cdf:

$$\begin{aligned} \text{from (a), } f_Y(y) &= \frac{d}{dy} F_Y(y) \\ &= f_X(g^{-1}(y)) \frac{1}{|g'(y)|} \cdot \underbrace{g'(y)}_{> 0 \text{ since } g \text{ is inc}} \\ \text{and from (b), } f_Y(y) &= -f_X(g^{-1}(y)) \frac{1}{|g'(y)|} \cdot \underbrace{g'(y)}_{< 0 \text{ since } g \text{ is dec}} \end{aligned}$$

combining the two cases to obtain

$$\text{change-of-var-formula: } f_Y(y) = f_X(g^{-1}(y)) \left| \frac{1}{|g'(y)|} \right|$$

Thm (Probability Integral Transform). Let X -continuous RV w/ cdf $F_X(x)$.
(i) Let $Y = F_X(x) = g(x)$
then
 $Y \sim \text{uniform}(0,1)$

Remarks: • doesn't hold for discrete RV's (F_X would be a step function)
• F_X might not be strictly increasing from 0 to 1

Recall (minimum): largest number that is less than or equal to every element in a set

Let $A \subseteq \mathbb{R}$,
then $\inf A = \begin{cases} \text{largest } a \in A \text{ s.t. } a \leq x \ \forall x \in A \\ +\infty \text{ if } A \text{ doesn't exist} \end{cases}$
e.g. $\inf(\{a, b\}) = a$
 $\inf(\{a, +\infty\}) = a \neq \inf(\{a, b\})$ which DNE

Inverse Function: $F_X^{-1}(y) = \inf \{x : F_X(x) \geq y\} \quad \forall y \in (0,1)$

where $\cdot F_X^{-1}(y) = +\infty$ only if $F_X(x) < y \ \forall x$.
 $\cdot F_X^{-1}(0) = -\infty$

Thm. Let $F_X(x)$ be a cdf of a RV.

(i) Let $Y \sim \text{uniform}(0,1)$,
and set $Z = F_X^{-1}(y) = \inf \{x : F_X(x) \geq y\} \quad \forall y \in (0,1)$
then Z has cdf F_X , i.e. $F_Z(z) = F_X(z)$

Remarks: • useful for generating RV's from Uniform RV's
• no restriction to continuity; works generally (for real-valued RV's)

2.2: Expectation

Def: the expectation of a RV $g(X)$ is

$$E[g(X)] = \begin{cases} \int_{-\infty}^{\infty} g(x) f_X(x) dx, & X \text{-continuous} \\ \sum_{x \in K} g(x) f_X(x), & X \text{-discrete} \end{cases}$$

Remarks: • also known as 'expected values', 'average' of a RV or the prob. dist. of the RV $g(X)$.

• expectation is associated w/ a distribution

Q: what is the expectation of X ?

Prop: If $X \in \mathbb{R}$, then by letting $g(x) = x$:

$$E(X) = \begin{cases} \int_{-\infty}^{\infty} x f_X(x) dx, & X \text{-cont.} \\ \sum_{x \in K} x f_X(x), & X \text{-discrete} \end{cases}$$

Note: If the domain of X is not a subset of \mathbb{R} (Euclidean space), then $E(X)$ may be invalid but a version of $E(X)$ may still be designed via the expectations of [big(x)] in #.

Remark: • A proper expectation exists only when the positive and negative part of integral are finite ($\int_{\mathbb{R}} f_X(x) dx < \infty$)
• $E(X)$ is undefined if $\int_{\mathbb{R}} |f_X(x)| dx$ is not absolutely convergent

Prop:

Let X -RV with dist. P_X .

Let X_1, \dots, X_n be n mutually independent RV's that have identical distributions as X , i.e. X_1, \dots, X_n is an n -iid sample of P_X

Def (Empirical Distribution): If \hat{P}_n is a prob. dist., denoted by P_n , s.t.

$$\text{if } Y \sim P_n \text{ then } y \in Y = \{X_1, \dots, X_n\} \text{ and } P(Y=y) = \frac{1}{n} \text{ for } y \in Y$$

Remark: • Thus, Y is discrete (regardless of X) and

$$E(Y) = \sum_y P(Y=y) \cdot y = \frac{1}{n} (X_1 + \dots + X_n) \neq E(X) = \int x f_X(x) dx$$

• Consequently, this $E(Y)$ is also called the "average" of the (data) sample X_1, \dots, X_n .

Linearity of Expectation

Thm: Let X, Y : real-valued RV's for which the expectations exist.

$$(i) E(aX+bY+c) = aE(X)+bE(Y)+c \quad \text{for } a, b, c \in \mathbb{R}$$

(ii) If $X \geq 0$ almost surely, i.e. $P(X=0) \geq 1$, then $E(X) \geq 0$

(iii) If $X \geq Y$ almost surely, i.e. $P(X \geq Y) \geq 1$, then $E(X) \geq E(Y)$

(iv) If $P(X \in \{a, b\}) = 1$, then $E(X) = a + b$

2.3: MOMENTS

Def: For each $n \in \mathbb{N}$, the n -th moment of X , or $F_X^{(n)}$
is $m_n := E(X^n)$

The n -th central moment of X is

$$\mu_n := E(X - E(X))^n$$

Remarks: • Let $Y = X - E(X)$, then $E(Y) = E(X - E(X)) = E(X) - [E(X)]^2$

• Let $n=2$. The second central moment is the variance, i.e.

$$\begin{aligned} \text{var}(X) &:= (E(X) - E(X))^2 \\ \text{sd}(X) &:= \sqrt{\text{var}(X)} \end{aligned}$$

• $E(X)$ captures the location (center) of X dist.

• $\text{var}(X)$ and/or $\text{sd}(X)$ captures the spread

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

Thm: If X is a RV w/ finite variance, then for any $a, b \in \mathbb{R}$:

$$\text{var}(aX+b) = a^2 \text{var}(X)$$

$$\text{Recall: } E(aX+b) = aE(X)+b$$

Remark: Let $X = X_1 + \dots + X_n$ where X_1, \dots, X_n are mutually independent.

Suppose X_1, \dots, X_n are identically distributed.

Then, we will later show (more formally) that

$$\text{var}(X) = \text{var}(X_1) + \dots + \text{var}(X_n) = n \cdot \text{var}(X_1)$$

whereas $E(aX+b) = aE(X) + b$ by (previous) defined L.O.E.

Def: Let X -RV w/ cdf F_X . The moment-generating function (mgf)

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

Provided that the expectation exists in some neighborhood of $t=0$ (i.e. $t \in \text{some } (-h, h) \text{ for } h > 0$).

Remarks: • $M_X(t) = \begin{cases} \int_{-\infty}^{+\infty} e^{tx} f_X(x) dx, & X-\text{cont. RV} \\ \sum_x e^{tx} f_X(x), & X-\text{discrete RV} \end{cases}$

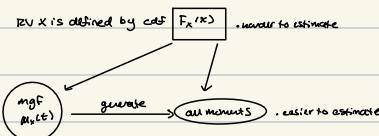
• the mgf can be used to "characterize" the dist. (of the RV)

• mgf uniquely determines all moments

Theorem: $\forall n \in \mathbb{N}, E(X^n) = \left. \frac{d^n}{dt^n} M_X(t) \right|_{t=0}$

- Remarks: - we also write RHS as $M_X^{(n)}(0)$
 - the derivative of the mgf evaluated at $t=0$ gives the n -th moment
 - gives meaning of "moment generating"

2.3: Moments (cont.)



Q1: Does the mgf uniquely determine the cdf?

Q2: Does the set of moments $\{E(X^k)\}_{k=1}^\infty$ uniquely determine the cdf?

Theorem: Let $F_X(x)$ and $F_Y(y)$ be two cdf's, all of whose moments exist.

i) If X and Y have bounded support, i.e. $\mathbb{R} = \{x : F_X(x) > 0\} \subset \mathbb{R}$, then

$$F_X(u) = F_Y(u) \Leftrightarrow E(X^k) = E(Y^k) \quad \forall k \geq 1, u \in \mathbb{R}$$

ii) If $M_X(t)$ and $M_Y(t)$ exist, and

$$M_X(t) = M_Y(t) \text{ for all } t \text{ in some neighborhood of } 0,$$

$$\text{then } F_X(u) = F_Y(u) \quad \forall u$$

Remark: the moments do not capture all information abt the dist. in the unbounded support scenarios

- motivation lies in the necessity to approximate pdf's

Def: $X_i \rightarrow X$ in distribution $\Leftrightarrow F_{X_i}(x) \rightarrow F_X(x)$
 If $F_{X_i}(x) \rightarrow F_X(x)$ at all continuity points of F_X

- implicit in the below thm.

Theorem (Convergence of mgf's leads to convergence of CDF):

Suppose X_1, X_2, \dots is a sequence of random variables,
 each w/ mgf $M_{X_i}(t)$.

Suppose $\lim_{i \rightarrow \infty} M_{X_i}(t) \rightarrow M(t)$

for all t in a neighborhood of 0, i.e.

$$\exists \delta > 0 \text{ s.t. } \forall t \in (-\delta, \delta), M_{X_i}(t) = M(t)$$

and $M(t)$ is a valid mgf.

Then $X_i \rightarrow X$ in distribution

where X is a RV w/ mgf $M_X(t)$

i.e. $F_{X_i}(x) \rightarrow F_X(x)$ at all points x where cdf F_X is continuous

Remark (Poisson Approx): By the convergence thm. of mgf,

If $X_n \sim \text{Binomial}(n, p)$, $x \in \mathbb{N}$, then,

$F_{X_n}(x) \rightarrow F_Y(x)$ for $Y \sim \text{Poisson}(\lambda)$

as $n \rightarrow \infty$ & λ where F_Y is continuous.

$$\Leftrightarrow M_{X_n}(t) \rightarrow M_Y(t) \quad \forall t$$

For Poisson, F_Y is a step function w/ continuity at $n \in \mathbb{N}$:

for $n \in \mathbb{N}$, $F_Y(y) = F_Y(y+1)$ continuous at $y+1/2$.

$$\Rightarrow P(X_n = x) \rightarrow P(Y = x) \quad \forall x \in \mathbb{N} \text{ (in fact, here } \lambda = np \in \mathbb{R})$$

Still holds if $X_n \sim \text{Binomial}(n, p_n)$ where $n \rightarrow \infty$, $p_n \rightarrow p$, $np_n \rightarrow \lambda$

Theorem: $\forall a, b \in \mathbb{R}$,

$$M_{aX+b}(t) = e^{bt} M_X(at)$$

2.4: Tools

Theorem (Leibniz's Rule): If

$f(x, \theta), a(x), b(x)$ are differentiable wrt θ

$\frac{\partial f}{\partial \theta}(x, \theta)$ is continuous on $\overline{\mathcal{X}} \times \mathcal{A}(\theta, \mathbb{R})$

can be measured

Then for $a(\theta), b(\theta) \in \mathcal{C}_b(\mathcal{X}, \mathbb{R})$:

$$\begin{aligned} \frac{d}{d\theta} \int_{\mathcal{X}} f(x, \theta) dx &= f(b(\theta), \theta) \frac{d}{d\theta} b(\theta) + f(a(\theta), \theta) \frac{d}{d\theta} a(\theta) \\ &\quad + \int_{\mathcal{X}} \frac{\partial f}{\partial \theta}(x, \theta) dx \end{aligned}$$

Corollary: If $a < b$, $b > 0$, then

$$\frac{d}{d\theta} \int_a^b f(x, \theta) dx = \int_a^b \frac{\partial f}{\partial \theta}(x, \theta) dx$$

To strengthen the above results for infinite domain of the integral, we need a new tool.

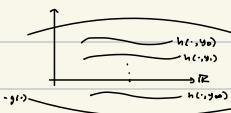
Theorem: Suppose

- $h(x, y)$ is continuous at $y=y_0$ for each fixed x .
- there's a function $g(x)$ s.t.

$$\text{envelope condition } \left\{ \begin{array}{l} h(x, y) \leq g(x) \forall y \\ g(x) \text{ is } \text{abs. cont.} \end{array} \right. \quad \text{on } \Omega$$

Then,

$$\lim_{n \rightarrow \infty} \int_{-\infty}^{\infty} h(x, y) dy = \int_{-\infty}^{\infty} \lim_{n \rightarrow \infty} h(x, y) dy$$



- all $h(x, y)$ must be contained in $(-g, g)$
- all g must be integrable

Remarks: g is called the dominating or envelope function of f .

- known as a version/consequence of Lebesgue's dominated convergence theorem (DCT).

- think of h as a collection of functions in $\mathcal{L}^1(\Omega \times \mathbb{R})$ etc.

Theorem: Suppose

- for each $x \in \mathbb{R}$, $f(x, \theta)$ is differentiable wrt θ , at $\theta \in \mathbb{R}$
- for each $\theta \in \mathbb{R}$, there's a function $g(x, \theta)$ and $\delta_0 > 0$ s.t.

$$\left\{ \begin{array}{l} |\frac{\partial}{\partial \theta} f(x, \theta)| \leq g(x, \theta) \quad \forall \theta \in (\theta - \delta_0, \theta + \delta_0) \\ \int_{-\infty}^{\infty} g(x, \theta) d\theta \text{ is cor. for each } \theta \in \mathbb{R} \end{array} \right.$$

Then,

$$\frac{d}{d\theta} \int_{-\infty}^{\infty} f(x, \theta) dx = \int_{-\infty}^{\infty} \frac{d}{d\theta} f(x, \theta) dx$$

holds for each $\theta \in \mathbb{R}$.

Remarks: - Leibniz's rule holds under this weaker envelope condition for $f(x, \theta)$.

- If we need only to differentiate at $\theta = \theta_0$, then it is sufficient that the envelope condition be satisfied for a neighborhood of θ_0 i.e. $\exists \delta \in (\theta_0 - \delta_0, \theta_0 + \delta_0)$ for some $\delta_0 > 0$, and only a dominating function $g(x, \theta)$ (for fixed $\theta = \theta_0$) is required.

- This theorem is a direct consequence of Lebesgue's dominated convergence theorem (DCT).

Properties: Let $X \sim \text{Exp}(\lambda)$; $f(x) = \frac{1}{\lambda} e^{-x/\lambda}$, $\lambda > 0, \infty$.

$$\text{Moments: } E(X^n) = \int_0^\infty \frac{1}{\lambda} x^n e^{-x/\lambda} dx, n=1, 2, \dots$$

This gives a recursive relation, i.e.

$$E(X^{n+1}) = \lambda E(X^n) + \lambda^2 \frac{d}{dx} E(X^n); \text{ (similar identity holds for broad family of dist's)}$$

Def (Pointwise convergence): $\lim_{n \rightarrow \infty} \sum_{x=0}^{\infty} h(x, \theta_n) = \sum_{x=0}^{\infty} h(x, \theta)$

$$S_n(\theta) \xrightarrow{\theta_n \rightarrow \theta} S(\theta)$$

Recall: Let $A \subset \mathbb{R}$. Then

- $S(A) = \text{sup}_{\theta \in A} S(\theta) \text{ s.t. } \forall n \in A, \epsilon > 0$,
- $\exists \delta > 0$ s.t. $\forall \theta \in A$, $|S(\theta) - S(A)| < \epsilon$.
 - $\text{if } A = (a, b)$, then $S(A) = S(a, b)$.
 - $\text{if } A = (a, \infty)$, then $S(A) = \infty$.

Def (Uniform convergence): $\sup_{\theta \in [a, b]} |S_n(\theta) - S(\theta)| \rightarrow 0$ as $n \rightarrow \infty$

Suppose $\sum_{x=0}^{\infty} h(x, \theta)$ exists (i.e. converges pointwise) $\forall \theta \in \mathbb{C}(a, b)$

Moreover, assume

- $\sum_{x=0}^{\infty} h(x, \theta)$ is continuous in θ for each x
- $\sum_{x=0}^{\infty} h(x, \theta)$ converges uniformly for all θ in a closed subinterval of (a, b)

Then,

$$\frac{d}{d\theta} \sum_{x=0}^{\infty} h(x, \theta) = \sum_{x=0}^{\infty} \frac{d}{d\theta} h(x, \theta)$$

Remark: This is a consequence of Lebesgue's DCT

Theorem: Suppose $\sum_{x=0}^{\infty} h(x, \theta)$ exists (i.e. converges pointwise) $\forall \theta \in [a, b]$

Moreover, assume

- $h(x, \theta)$ is continuous in θ for each fixed x
- $\sum_{x=0}^{\infty} h(x, \theta)$ converges uniformly on $[a, b]$

Then,

$$\int_a^b \sum_{x=0}^{\infty} h(x, \theta) d\theta = \sum_{x=0}^{\infty} \int_a^b h(x, \theta) d\theta$$

3.1: Discrete Distributions

(Recall: RV X is discrete if the range of X ($\subset \mathbb{R}$) is countable)

$$\text{Def (Uniform RV)} := X \sim \text{Uniform}(1, N) \text{ if } P(X=x) = \frac{1}{N}, \quad x=1, \dots, N.$$

we sometimes write
 $P(X=x|N) = \frac{1}{N}$

Remarks: $E(X) = \frac{N+1}{2}$, $V(X) = \frac{(N+1)(N-1)}{12} = \frac{(N-1)^2}{12}$

Def (Hypergeometric Dist):

Experiment: N balls, M red, $N-M$ green
 Pick K balls uniformly at random (w/o replacement)

Let X : number of reds.

then, $P(X=x|N, M, K) = \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}}, \quad x=0, 1, \dots, K$

Remarks: counting argument implies $\sum_{x=0}^K \binom{M}{x} \binom{N-M}{K-x} = \binom{N}{K}$

$E(X) = \frac{KM}{N}, V(X) = \frac{KM}{N} \frac{(N-M)(N-K)}{N(N-1)}$

Def (Bernoulli Dist): Let $X \sim \text{Bernoulli}(p)$, $0 < p \leq 1$

$$P(x) := \begin{cases} P(X=1|p) = p \\ P(X=0|p) = 1-p \end{cases}$$

Remarks: $E(X)=p, V(X)=p(1-p)$.

Def (Binomial Dist): Perform n -independent-and-identically-distributed (iid) Bernoulli trials.

Let Y : number of successes ($\sim \text{Bin}(n, p)$).

Then $Y \sim \text{Binomial}(n, p)$ and

$$P(Y=y|n, p) = \binom{n}{y} p^y (1-p)^{n-y}$$

Remarks: $E(Y)=np, V(Y)=np(1-p), M_Y(t) = (pe^t + (1-p))^n$

Def (Poisson Dist): Let $X \sim \text{Poisson}(\lambda), \lambda > 0$.

$$P(X=x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x=0, 1, \dots$$

Remarks: $e^\lambda = 1 + \lambda + \frac{\lambda^2}{2!} + \frac{\lambda^3}{3!} + \dots$

- arises in modelling the following (count) experiments:
 - number of buses arriving within a period of time
 - number of stars observed in a telescope
 - number of incidents found in a segment of a network
- $E(X) = \lambda = V(X), M_X(t) = \exp[\lambda(e^t - 1)]$.

Prop (Poisson Approx): (of the Binomial Dist):

If $X_n \sim \text{Binomial}(n, p_n)$ and $np_n \rightarrow \lambda$ as $n \rightarrow \infty$

then $X_n \xrightarrow{d} Y$ where $Y \sim \text{Poisson}(\lambda)$.

Def (Negative Binomial Dist):

Experiment: count the number of independent Bernoulli(p) trials until obtaining the r -th ~~success~~ success

$X \sim \text{NegBinom}(r, p), r \in \mathbb{N}_0$

$$P(X=x) = \binom{x-1}{r-1} p^r (1-p)^{x-r}, \quad x=r, r+1, \dots$$

Corollary: Let $Y = X-r$ represent the number of failures that occur before the r -th success in a sequence of independent Bernoulli(p) trials

then $P(Y=y|p_r) = P(X=y+r|p_r) = \binom{y+r-1}{r-1} p^r (1-p)^{y+r}, \quad y=0, 1, \dots$

$$= (-1)^y \binom{y+r-1}{r-1} p^r (1-p)^y$$

Remark: (1) motivates the "negative binomial" label

• Here, $n = y+r-1 \Rightarrow y = n-r+1$

• $E(Y) = \frac{r(1-p)}{p}, \quad V(Y) = \frac{r(1-p)}{p^2}$

Prop: Let $uX = E(Y) = \frac{r(1-p)}{p}$, then
 $\text{Var}(Y) = \frac{r(1-p)}{p} \cdot r(1-p) + uX^2$

Remarks: conceptually the 'quadratic relation'

Prop(Poisson Approx.): If $r \rightarrow \infty$, $p \rightarrow 1$ such that $r(1-p) \rightarrow \lambda$,
then $\begin{cases} E(Y) \rightarrow \lambda \\ V(Y) \rightarrow \lambda \end{cases}$
Moreover, $Y \xrightarrow{d} \text{Poisson}(\lambda)$

Def(Geometric Dist.): Let $X \sim \text{GEOMETRIC}(p)$, $P(X=1) = p$.
 $P(X=x|p) = p(1-p)^{x-1}$, $x=1, 2, \dots$

Remarks: This is a special case of $\text{NegativeBinomial}(r, p=1)$:
Hence $\begin{cases} E(X) = (1-p)/p \\ V(X) = (1-p)/p^2 \end{cases}$

$\Rightarrow X$: # trials until the first success

Prop('memoryless Property'): If $s > t$, then
 $P(X>s|X>t) = P(X>s-t)$

Remarks: Given $X > t$ (no success in the first t trials), the dist. of additional waiting time
until the first success does not depend on t , i.e.
past failures do not change the tail prob's for the future

3.2.a: Continuous Distributions

- continuous distributions put prob. mass on cont. spaces
- Formally, a continuous RV is one whose CDF is a continuous function

Def(Uniform Dist.): Let $X \sim \text{Uniform}([a, b])$, $a < b$.
 $f(x|a, b) = \begin{cases} \frac{1}{b-a}, & x \in [a, b] \\ 0, & \text{otherwise} \end{cases}$

Remarks: $E(X) = \frac{a+b}{2}$, $V(X) = \frac{(b-a)^2}{12}$

- uniform RV can either be continuous or discrete
depends on continuous range or a finite, discrete set of values

Def(Gamma Dist.): Let $X \sim \text{Gamma}(k, \beta)$ for $k > 0$: shape
 $\beta > 0$: scale

$$f(x|a, \beta) = \frac{1}{\Gamma(a)\beta^a} x^{a-1} e^{-x/\beta}, x > 0$$

Corollary(Gamma Function): $\Gamma(k) = \int_0^\infty t^{k-1} e^{-t} dt, k > 0$
 $\Gamma'(k) \beta^k = \int_0^\infty x^{k-1} e^{-x\beta} dx$
 $\Gamma'(k+1) = \Gamma'(k) \cdot k$
 $\Gamma'(n) = (n-1)! = (n-1) \cdots \cdot 1$

Remarks: sometimes we use $\delta = 1/\beta$ instead as scale parameter

$$E(X) = a\beta, V(X) = a\beta^2$$

$$X \sim \text{Gamma}(a, 1)$$

$$M_X(t) = \left(\frac{1}{1-\beta t}\right)^a, t < \frac{1}{\beta}$$

Prop: $\forall x > 0$ $P(X \leq x) = P(Y \leq x)$ for $\begin{cases} Y \sim \text{Gamma}(a, \beta) \\ Y \sim \text{Poisson}(\lambda \beta) \end{cases}$

$$\Leftrightarrow \int_0^x \text{pdf}_{\text{Gam}}(t) dt = \sum_{j=0}^{\infty} \text{pmf}_{\text{Pois}}(j)$$

Prop(Special cases):

a. Let $X \sim \text{Beta}(a, b)$: $\begin{cases} n = pr, p = \text{integer} \\ \beta = 1 - \frac{n}{p} \end{cases}$
Then $f(x|p) = \frac{1}{\Gamma(n+p)} x^{n-1} e^{-nx/p}, x > 0 \sim \chi_{2(p-n)}^2$

Remarks: we will later see that, for $Z_i \sim N(0, 1)$:

$$\sum_{i=1}^p Z_i^2 \sim \chi_p^2$$

b. Let $X \sim \text{Exponential}(\beta)$, $\beta > 0$, and note $E(X) = 1$
 $\text{Then } f(x) = \frac{1}{\beta} e^{-x/\beta}, x \geq 0, \text{ i.e.}$
the exponential pdf.
Note that the exponential dist. (like the geometric dist.)
has the memoryless property, i.e.
 $P(X > s | X > t) = P(X > s-t) \text{ for } s > t \geq 0$

c. If $X \sim \text{Exp}(\beta)$, $\beta > 0$, letting $Y = X^{1/\beta}$, $y \geq 0$,
then
 $F(y) = P(Y \leq y) = P(X \leq y^\beta) = 1 - e^{-y^\beta/\beta}$
where $\frac{d}{dy} e^{-y^\beta/\beta} = e^{-y^\beta/\beta} \cdot \beta y^{\beta-1} e^{-y^\beta/\beta}$
 $f_y(y) = \frac{1}{\beta y} F'(y) = \frac{\beta}{\beta} y^{\beta-1} e^{-y^\beta/\beta}, \text{ or simply}$
where $y \sim \text{Weibull}(\beta, \beta)$

Recall: Let $X \sim \text{Exp}(\beta)$, $F(x) = \frac{1}{\beta} e^{-x/\beta}, x \geq 0$
 $F(x) = \int_0^x \frac{1}{\beta} e^{-t/\beta} dt = \frac{1}{\beta} \left[-e^{-t/\beta} \right]_0^x$
 $= -e^{-x/\beta} \Big|_0^x$
 $= 1 - e^{-x/\beta}, x \geq 0.$

Remark: Weibull dist. useful for handling extreme rare events



3.2.5: Continuous Dist. (cont.)

Def(Normal Dist.): $X \sim \text{Normal}(\mu, \sigma^2)$, $\mu \in \mathbb{R}$, $\sigma^2 > 0$
 $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right), x \in \mathbb{R}$

Remark: Parameters μ : mean, σ^2 : variance

Prop (moments):
 $\frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} x e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx = \mu \quad (1)$
 $\frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} x^2 e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx = \mu^2 + \sigma^2 \quad (2)$
 $\frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} x^k e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx = \mu^k + \sigma^2 k + \text{higher terms} \quad (3) \quad \text{reformulating } \mathbb{E}(X^k) = \mathbb{E}(X^k) - \mathbb{E}^2(X)$

$\forall k \in \mathbb{N}$, $\sigma > 0$.

Def(Standardization): Subtracting the mean and rescaling by the SD

Does not always keep the RV to remain in the family of distributions)

Corollary: If $X \sim N(\mu, \sigma^2)$, then
Standard normal RV: $Z := \frac{X-\mu}{\sigma} \sim N(0, 1)$
 $f(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z^2\right), z \in \mathbb{R}$
 $F_z(z) = P(Z \leq z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2} dt$

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

$\Rightarrow \int_{-\infty}^{\infty} e^{-\frac{1}{2}t^2} dt = \sqrt{\frac{\pi}{2}} = \sqrt{\pi} \int_0^{\infty} e^{-\frac{1}{2}t^2} dt = \sqrt{\pi}, \text{ i.e. a definition for } \pi$

$\therefore \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi} \quad (\star)$

Prop: If $X \sim N(\mu, \sigma^2)$,
then $Y = X^2 \sim \chi^2_2$ s.t.
 $f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-y/2\sigma^2}$
 $= \frac{1}{\Gamma(1/2)} y^{-1/2} e^{-y/2\sigma^2} \text{ by (4)}$
i.e. a pdf for Gamma($\alpha=1/2, \beta=2\sigma^2$)

Remark: More generally, let $X \sim \chi^2_p$, then $f(x) = \frac{1}{\Gamma(p/2)} \frac{x^{p/2-1}}{2^{p/2}} e^{-x/2}, x \geq 0, p \in \mathbb{N}$

where $\chi^2_p \stackrel{d}{=} \text{Gamma}(p/2, 2)$

Prop(Normal Approx.): Let $X \sim \text{Binomial}(np)$ s.t. $\begin{cases} E(X) = np \\ \text{Var}(X) = np(1-p) \end{cases}$
Let $\begin{cases} np \rightarrow \infty \\ np(1-p) \rightarrow \infty \end{cases}$. Then $X \approx \text{Normal}(np, np(1-p))$

Remark: We will later show via CLT that

$\frac{1}{\sqrt{n}}(X-np) \xrightarrow{d} \text{Normal}(0, np(1-p))$ if $\begin{cases} n \rightarrow \infty \\ p \neq 0, 1 \end{cases}$

Corollary (continuity correction): Let $X \sim \text{Bin}(np)$ and $Y \sim N(np, np(1-p))$.

then,

$$\begin{aligned} P(X=k) &\approx P(Y \leq k+0.5) \\ P(X \geq k) &\approx P(Y \geq k+0.5) \approx P(X \geq k) \\ P(X \neq k) &\approx P(Y \neq k-0.5) \approx P(X \neq k) \end{aligned}$$

Remark: + cont. dist., e.g. normal, assign points zero prob; must look at interval around point to capture any approximated mass

+ w/o correction, we tend to underestimate the tails

3.2.c: Continuous Dist. (cont.)

Def (Beta Dist.): Let $K \sim \text{Beta}(\alpha, \beta)$, $\alpha, \beta > 0$

$$f(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} \text{ for } x \in (0, 1)$$

n-th moments: $E(X^n) = \frac{\Gamma(\alpha+\beta)/\Gamma(n+\alpha)}{\Gamma(\alpha)\Gamma(\beta+n)}$

$$E(X) = \frac{\alpha}{\alpha+\beta}$$

$$V(X) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$$

Remarks: Beta function: $B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1}$

$$\text{Fact: } B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \Rightarrow \text{Beta}(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

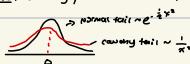
$\alpha=\beta=1 \Rightarrow \text{Beta}(1, 1) \equiv \text{Uniform}(0, 1)$

Beta pdfs for different (α, β) : $\begin{cases} \alpha > 1, \beta > 1 \Rightarrow \text{Beta is unimodal} \\ \alpha < 1, \beta < 1 \Rightarrow \text{Beta is } \text{U-shaped} \\ \alpha < 1, \beta > 1 \Rightarrow \text{Beta is } \text{M-shaped} \\ \alpha > 1, \beta < 1 \Rightarrow \text{Beta is } \text{W-shaped} \end{cases}$

Def (Cauchy Dist.): Let $X \sim \text{Cauchy}(\theta)$, $\theta \in \mathbb{R}$.

$$\text{Then: } f(x) = \frac{1}{\pi} \frac{1}{1+(x-\theta)^2} \text{ for } x \in \mathbb{R}.$$

Remarks: Cauchy pdf also bell-shaped:



θ is the median (constant) of X

Recall: $\begin{cases} E(X) = \theta \\ E(X^2) = D(X) + E(X)^2 \end{cases}$

Prop: Let $X, Y \stackrel{iid}{\sim} N(\mu, \sigma^2)$. Then $\frac{X-Y}{\sqrt{2}} \sim \text{Cauchy}(\mu)$ Proof (see later).

Def (Log-Normal Dist.): If $Y \sim N(\mu, \sigma^2)$, then $X = e^Y \sim \text{LogNormal}$, i.e.

$$Y \sim N(\mu, \sigma^2) \Rightarrow \log(Y) \sim N(\mu, \sigma^2)$$

$$F(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^x \exp\left(-\frac{(z-\mu)^2}{2\sigma^2}\right) dz$$

3.3: Exponential Families

Def: A family of pdf or pmf is called an exponential family if it has the form

$$\begin{aligned} f(x|\theta) &= c(\theta) h(x) \exp\left\{\sum_{i=1}^k w_i(\theta) t_i(x)\right\}, \quad x \in \mathcal{X} \\ &= c(\theta) h(x) \exp\left\langle w(\theta), t(x)\right\rangle \quad c(\theta) \end{aligned}$$

Remarks: where $w(\theta) = (w_1(\theta), \dots, w_k(\theta))$

$$t(x) = (t_1(x), \dots, t_k(x))$$

vector of sufficient statistics (coming later)

Here, $h(x) \geq 0$

$$t(x) = (t_1(x), \dots, t_k(x)) \text{ depends only on } \theta \text{ (constant)}$$

$$w(\theta) = (w_1(\theta), \dots, w_k(\theta)) \text{ depends only on } \theta \text{ (parameters)}$$

clearly, $c(\theta)$ is the reciprocal of the normalizing constant, i.e.

$$\int f(x|\theta) dx = 1 \Rightarrow \frac{1}{c(\theta)} \int_{\mathcal{X}} h(x) \exp\left\langle w(\theta), t(x)\right\rangle dx \quad (\text{or sum if discrete})$$

$$\begin{aligned} \text{Ex(Binomial)}: \text{Let } X \sim \text{Bin}(n, p); f(x) = \binom{n}{x} p^x (1-p)^{n-x}, x=0, 1, \dots, n \\ &= \frac{(1-p)^n p^n}{n!} \frac{x! (n-x)!}{x!(n-x)!} \log \frac{p}{1-p} \end{aligned}$$

$$\begin{aligned} \text{Ex(normal)}: \text{Let } X \sim N(\mu, \sigma^2); f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+ \\ &= \dots \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} x^2 + \frac{\mu}{\sigma^2} x\right\}. \end{aligned}$$

REMARKS: - other distributions in the exponential family:

Poisson, Geometric, NegBinom, Gamma, Exponential, log-normal

- Dirichlet is not in exponential families: causality, Mixtures

Then if X is a RV w/ pdf/pmf in the exponential form (e.g.), then

$$D E \left\{ \sum_{i=1}^k \frac{\partial w_i(\theta)}{\partial \theta_j} t_i(X) \right\} = \frac{\partial}{\partial \theta_j} \log(1/c(\theta)), j=1, \dots, k$$

★ check

$$\Rightarrow \text{Var} \left(\sum_{i=1}^k \frac{\partial w_i(\theta)}{\partial \theta_j} t_i(X) \right) = \frac{\partial^2}{\partial \theta_j^2} \log(1/c(\theta)) - E \left\{ \sum_{i=1}^k \frac{\partial w_i(\theta)}{\partial \theta_j} t_i(X) \right\}$$

Remark (high-level): Differentiating the logarithm of the normalizing constant w.r.t. a parameter results in sufficient expectations of sufficient statistics

Differentiating w.r.t. the natural parameter corresponding to the sufficient statistic of interest

when we differentiate w.r.t. a natural parameter $w_j(\theta)$, we're targeting the expectation of $t_j(X)$, i.e. expectation of a sufficient statistic

Corollary (Simplex): Suppose

$$w(\theta) = \theta = (\theta_1, \theta_2, \dots, \theta_n) \in \mathbb{R}^n$$

$$t_i(x) = t_i(x_1, x_2, \dots, x_n) \in \mathbb{R}^n$$

then $f(x|\theta) = c(\theta) h(x) \exp(-\theta_1 x_1 - \theta_n x_n)$ and

$$w_j := E(t_j(X)) = \frac{\partial}{\partial \theta_j} \log(1/c(\theta))$$

$$\Rightarrow \text{Var}(t_j(X)) = \frac{\partial^2}{\partial \theta_j^2} \log(1/c(\theta))$$

3.4: Location and Scale Families

Facts: If $f(x)$ is a valid pdf on \mathbb{R} , then $V \in \mathbb{R}$, $\sigma > 0$,

the function

$$g(x|V, \sigma) := \frac{1}{\sigma} f\left(\frac{x-V}{\sigma}\right), x \in \mathbb{R}$$

is a valid pdf on \mathbb{R} .

Def: Let $f(x)$ be any pdf on \mathbb{R} .

$$\text{The family of pdf } \left\{ g(x|V, \sigma) : \frac{1}{\sigma} f\left(\frac{x-V}{\sigma}\right) \mid \sigma > 0 \right\}$$

is called a location-scale family of distributions for $\left\{ \begin{array}{l} V: \text{location param} \\ \sigma: \text{scale param} \end{array} \right.$

Ex: Let $f(x) = \frac{1}{\sqrt{\pi}} \exp\left(-\frac{1}{2} x^2\right)$ s.t. $X \sim N(0, 1)$.

then, we obtain the $L-S$ family, i.e.

$$\left\{ g(x|V, \sigma) = \frac{1}{\sigma} f\left(\frac{x-V}{\sigma}\right) = \frac{1}{\sigma \sqrt{\pi}} \exp\left(-\frac{1}{2} \left(\frac{x-V}{\sigma}\right)^2\right) \mid \sigma > 0 \right\}.$$

Special cases: 1) $\sum g(x|V, \sigma) = f(x-V)$ ($\forall \sigma \in \mathbb{R}$) is location family

2) $\sum g(x|V, \sigma) = \frac{1}{\sigma} g\left(\frac{x}{\sigma}\right)$ ($\sigma > 0$) is scale family

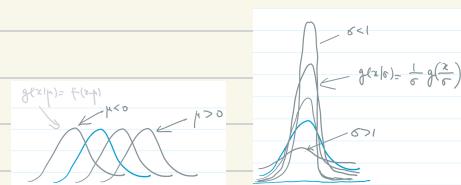
↳ why not $f(x-V)$?

Prop: If $X \sim f$, then $X+V \sim g(x) = f(x-V)$



If $X \sim f$, then $\sigma X \sim g(x) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right)$

If $X \sim f$, then $\sigma X + V \sim g(x) = \frac{1}{\sigma} f\left(\frac{x-V}{\sigma}\right)$



Thm: Suppose Y is a RV w/ pdf $f(y)$ and $E(Y)$, $V(Y)$ exist.

If X is a RV w/ pdf $\frac{1}{\sigma} f\left(\frac{x-V}{\sigma}\right)$, then

$$\begin{cases} E(X) = V(E(Y)) + V \\ V(X) = \sigma^2 V(Y) \end{cases}$$

3.5: Inequalities and Identities

- when we can't calculate probabilities, it is important to estimate them by inequalities to obtain bounds



Theorem (Chebyshev's Inequality): Let X be a RV, $g(x) \geq 0 \forall x$. Then, $P(g(X) \geq r) \leq \frac{E(g(X))}{r} \forall r > 0$

$$\Leftrightarrow E(g(X)) \geq r \cdot P(g(X) \geq r)$$

Corollary: Let $g(x) = \frac{(x-m)^2}{\sigma^2}$, where $\{m = E(X), \sigma^2 = \text{Var}(X)\}$. Then, $P(|X-m| \geq r\sigma) \leq \frac{1}{r^2}$. Alternatively, letting $t = \frac{r\sigma}{\sigma}$: $P(|X-m| \geq t\sigma) \leq \frac{1}{t^2}$ and $P(|X-m| \geq t\sigma) \geq 1 - \frac{1}{t^2}$.

Remarks: $t=2$: $P(|X-m| \geq 2\sigma) \leq 1/4 = 25\%$.
 $t=3$: $P(|X-m| \geq 3\sigma) \leq 1/9 = 11.1\%$.
 $t=4$: $P(|X-m| \geq 4\sigma) \leq 1/16 = 6.25\%$.

Theorem (Tighter Inequality for normal tails): If $Z \sim N(0,1)$, then for $t > 0$: $P(|Z| \geq t) \leq \sqrt{\frac{2}{\pi}} \frac{e^{-t^2/2}}{t}$



Remark: If $t=2$, $P(|Z| \geq 2) \leq \sqrt{\frac{2}{\pi}} \frac{e^{-1}}{2} = 0.054 << 1/4$ (from above, Chebyshev)

Corollary: In general, $\sqrt{\frac{2}{\pi}} e^{-t^2/2} \ll \frac{1}{t^2}$, exponentially small.

Identities

Poisson Identity: If $X \sim \text{Poisson}(\lambda)$, then $P(x) = e^{-\lambda} \frac{\lambda^x}{x!}$.
Set $\begin{cases} P(x+1) = e^{-\lambda} \lambda^{x+1} \\ P(x+2) = e^{-\lambda} \frac{\lambda^{x+2}}{(x+1)!} \end{cases}$
 $= P(x+1) \frac{\lambda}{x+1}$

Note: Recursion-like identities (like this one) may be useful in various situations that require such computations.

Gamma Identity: If $X_{\alpha, \beta} \sim \text{Gamma}(\alpha, \beta)$ with $\{ \alpha > 1, \beta > 0 \}$

Then $\forall a, b$

$$P(X_{\alpha, \beta} \in (a, b)) = \beta \left(F(a, \alpha, \beta) - F(b, \alpha, \beta) \right) + P(X_{\alpha+1, \beta} \in (a, b))$$

Remark: If $\alpha \in \mathbb{N}$, the above identity allows us to recurse to

$\text{Gamma}(\alpha-1, \beta), \text{Gamma}(\alpha-2, \beta), \dots$, and so on

to $\text{Gamma}(1, \beta) = E[X]$.

Stein's Identity for Normal RV: If $X \sim N(\mu, \sigma^2)$,

g is a differentiable function s.t. $E[g'(x)] < \infty$

then

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

Remark: see applications in practice problems.

Hwang's Identity (for Discrete Variables):

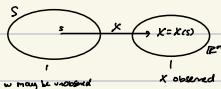
Let g : function with $|E(g)| < \infty$ and $|g(x)| < \infty$.

If $X \sim \text{Poisson}(\lambda)$, then $E(Xg(x)) = E(g(x))$

If $X \sim \text{Binomial}(n, p)$, then $E((1-p)g(X)) = E\left[\sum_{k=0}^n g(k)p^k\right]$

4.1: Joint and marginal distributions

DEF: An n -dimensional random vector is a function from a sample space into space \mathbb{R}^n (may change to other spaces)



Note: If $n=2$, $X = (X_1, X_2) \in \mathbb{R}^2$ is called a bivariate vector

DEF: Let (X, Y) be a discrete bivariate vector.

Then the function from $\mathbb{R}^2 \rightarrow \mathbb{R}$:

$$f_{X,Y}(x,y) := P(X=x, Y=y)$$

Joint prob. mass function (pmf) on (X, Y)

Remark: Denote by $f_{X,Y}(x,y)$ or $f_{Y|X}(y|x)$

where $f_{Y|X}(y|x)$ not necessarily $f_{X,Y}(x,y)$

Corollary: For $A \subseteq \mathbb{R}^2$, by AOP, we have $P((X, Y) \in A) := \sum_{(x,y) \in A} f_{X,Y}(x,y)$

Prop: Let $g(x,y)$ be a function from $\mathbb{R}^2 \rightarrow \mathbb{R}$.

Then $g(X,Y)$ is a real-valued RV: $E[g(X,Y)] := \sum_{(x,y) \in \mathbb{R}^2} g(x,y) f_{X,Y}(x,y)$

Prop (linearity of expectation): If g_1, g_2 are real-valued functions on \mathbb{R}^2 ; $a, b \in \mathbb{R}$, then

$$E[a g_1(X,Y) + b g_2(X,Y)] = a E[g_1(X,Y)] + b E[g_2(X,Y)]$$

Thm: Let (X, Y) be a discrete bivariate RV w/ joint pmf $f_{X,Y}(x,y)$.

Then, X and Y are discrete RV's w/ the following pmf's:

$$\begin{cases} f_X(x) = \sum_y f_{X,Y}(x,y) \\ \text{marginal pmf:} \\ f_Y(y) = \sum_x f_{X,Y}(x,y) \end{cases}$$

Remark: f_X and f_Y are called the marginal pmf's of X and Y

the distribution of X and Y are the marginal distributions (curr. dist. of (X, Y))

the dist. of random vector (X, Y) is the joint dist. of X and Y

Corollary: Joint dist./pmf completely determines its marginal dist./pmf
But marginal dist./pmf's do not determine the joint dist.

- CONTINUOUS BIVARIATE RV'S are described via joint prob. density functions

Def: a function $f_{X,Y}(x,y)$ from $\mathbb{R}^2 \rightarrow \mathbb{R}$ is called a joint pdf of the continuous bivariate random vector (X, Y) if, for every $A \subseteq \mathbb{R}^2$

$$P((X, Y) \in A) = \iint_A f_{X,Y}(x,y) dx dy$$

Remarks: $\cdot f_{X,Y}(x,y) \geq 0 \forall (x,y) \in \mathbb{R}^2$

$$\cdot \iint_{\mathbb{R}^2} f_{X,Y}(x,y) dx dy = 1$$

Def: If $g(X,Y)$ is a real-valued function,
then $g(X,Y)$ is a random variable w/ the expectation

$$E[g(X,Y)] := \iint_{\mathbb{R}^2} g(x,y) f_{X,Y}(x,y) dx dy$$

$$\text{marginal pmf's: } f_X(x) = \int_{\mathbb{R}} f_{X,Y}(x,y) dy$$

$$f_Y(y) = \int_{\mathbb{R}} f_{X,Y}(x,y) dx$$

$$\begin{aligned} \text{Joint pdf} &:= F_{XY}(x,y) := P(X \leq x, Y \leq y) \\ &= \int_{-\infty}^y \int_{-\infty}^x f_{XY}(s,t) ds dt \end{aligned}$$



$$\text{Corollary: } F_{XY}(x,y,z) := P(X \leq x, Y \leq y, Z \leq z) = \int_{-\infty}^z \int_{-\infty}^y \int_{-\infty}^x f_{XYZ}(s,t,u) ds du dz$$

Remarks: $\cdot 0 \leq F_{XY}(x,y) \leq 1 \quad \forall (x,y) \in \mathbb{R}^2$

$\cdot F_{XY}(y,z) \uparrow \text{ w.r.t. } x \in \mathbb{R}$

$\cdot F_{XY}(x,y) \uparrow \text{ w.r.t. } y \in \mathbb{R}$

\cdot Fundamental thm. of calculus (bivariate case) gives

$$\frac{dF_{XY}(x,y)}{dx dy} = f_{XY}(x,y)$$
 at continuous points (x,y)

of function for

4.2a: conditional distributions

Given a bivariate RV $(X,Y) \in \mathbb{R}^2$, we are interested in, say,

$$P(Y \in B | X \in A) \text{ for } A, B \subset \mathbb{R}$$

Prop: Let (X,Y) be discrete.

$$P(Y \in B | X \in A) = \frac{P(A \cap Y \in B)}{P(A)}$$

$$\text{where } P(A \cap Y \in B) = \sum_{x \in A, y \in B} P_{XY}(x,y)$$

$$\text{and } P(A) = \sum_{x \in A} P_{XY}(x,x) = \sum_{x \in A} f_{XY}(x,x)$$

$$\text{so, } P(Y \in B | X \in A) = \sum_{x \in A} \sum_{y \in B} P_{XY}(x,y) / \sum_{x \in A} f_{XY}(x,x)$$

$$P(Y \in B | X \in A) = \sum_{y \in B} \frac{f_{Y|X}(y|x)}{\sum_{x \in A} f_{Y|X}(y|x)}$$

Corollary: Let $A = \{x\}, B = \{y\}$. Then $P(Y=y | X=x) = \frac{f_{Y|X}(y|x)}{f_{X|X}(x)}$

Def: Let (X,Y) be a discrete random vector w/ joint pmf $f_{XY}(x,y)$. Then, for any $x \in \mathbb{R}$, if $f_X(x) = P(X=x) > 0$, the conditional pmf is:

$$\begin{aligned} f_{Y|X}(y|x) &= P(Y=y | X=x) \\ &= \frac{f_{XY}(x,y)}{f_X(x)} = \frac{f_{XY}(x,y)}{\sum_y f_{XY}(x,y)} \end{aligned}$$

Remarks: $\cdot f(y|x) \geq 0 \quad \forall y$

$$\sum_y f(y|x) = \sum_y \frac{f_{XY}(x,y)}{\sum_y f_{XY}(x,y)} = \frac{\sum_y f_{XY}(x,y)}{\sum_y f_{XY}(x,y)} = 1.$$

\cdot denote $F_{Y|X}(y|x)$

Note: By definition, the marginal pmf (f_X, f_Y) and the conditional pmf ($f_{Y|X}, f_{X|Y}$) are completely determined by the joint pmf (f_{XY})

$$\begin{aligned} \text{Corollary: } f_{Y|X}(y|x) &= \frac{f_{XY}(x,y)}{f_X(x)} \\ \Rightarrow f_{XY}(x,y) &= f_X(x) f_{Y|X}(y|x) \\ \text{similarly, } f_{XY}(x,y) &= f_Y(y) f_{X|Y}(x|y) \\ \text{wherever } f_X(x) > 0 \text{ and } f_Y(y) > 0 \end{aligned}$$

Remark: the marginal pmf and conditional pmf are contained to completely determine the joint pmf

$$\begin{aligned} \text{Def (conditional pdf): } &\text{consider the continuous bivariate random vector } (X,Y) \in \mathbb{R}^2. \\ \text{then } P(Y \in B | X=x) &:= \lim_{\epsilon \downarrow 0} P(Y \in B | X \in (x-\epsilon, x+\epsilon)) \\ &= \int_B \frac{f_{XY}(x,y)}{f_X(x)} dy \end{aligned}$$

assuming $f_{XY}(x,y)$ continuous at (x,y)

$$\begin{aligned} \text{Def (conditional pdf): } &\text{let } (X,Y) \text{ be a bivariate random vector w/ joint pdf } f_{XY}(x,y). \\ \text{then, for any } x \in \mathbb{R}, f_{X|Y}(y|x), & \\ f_{Y|X}(y|x) &:= \frac{f_{XY}(x,y)}{f_X(x)} = \frac{f_{XY}(x,y)}{\int f_{XY}(x,z) dz} \end{aligned}$$

Remarks: $\cdot f_{Y|X}(y|x) \geq 0 \quad \forall y, x \quad (f_X(x) > 0)$

$$\int f_{Y|X}(y|x) dy = 1$$

\cdot identities hold: $f_{XY}(x,y) = \begin{cases} f_X(x)f_{Y|X}(y|x) & \text{when } f_X(x) > 0 \\ f_Y(y)f_{X|Y}(x|y) & \text{when } f_Y(y) > 0 \end{cases}$

• By varying g , we actually have a collection of distributions for RV $Y: \{f_{Y|X}(y|x)\}_{x \in \mathbb{R}}$

Def (Cond. Expectation): Given $(X,Y) \sim f_{X,Y}$, $g: \mathbb{R} \rightarrow \mathbb{R}$ a function

$$\text{then, } E[g(Y)|X=x] := E[g(Y)|X=x] \\ := \int g(y) f_{Y|X}(y|x) dy$$

Remarks: two equivalent ways to define a joint dist. for a bivariate random vector (X,Y) :

a) joint pdf: $f_{X,Y}(x,y) \forall x,y$

$$b) f_{X,Y}(x,y) := \begin{cases} f_X(x)f_{Y|X}(y|x), & \text{if } x,y \in \mathbb{R} \\ 0, & \text{otherwise} \end{cases}$$

$$\text{where: } f_X(x) = \int f_{X,Y}(x,y) dy$$

Note: (b) often easier in applications

4.2.b: Independence

Def: Let (X,Y) be a bivariate random vector, with joint pdf/pdf $f_{X,Y}(x,y)$ and marginal pdf/pdf $f_X(x)$ and $f_Y(y)$.

Then X and Y are independent RV's if, $\forall x, y \in \mathbb{R}$

$$f_{X,Y}(x,y) = f_X(x)f_Y(y)$$

Remarks: 1) If $X \perp\!\!\!\perp Y$ then $f_{X,Y}(x,y) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_X(x)f_{Y|X}(y|x)}{f_X(x)} = f_{Y|X}(y|x) \cdot f_Y(y)$
(does not depend on x)

Moreover, $\forall A, B \subset \mathbb{R}$,

$$P(X \in A, Y \in B) = \frac{\int \int f_{X,Y}(x,y) dx dy}{\int f_{X,Y}(x,y) dx} = \frac{\int_B \int_A f_{X,Y}(x,y) dx dy}{\int_A f_{X,Y}(x,y) dx} = \frac{\left(\int_A f_X(x) dx \right) \left(\int_B f_{Y|X}(y|x) dy \right)}{\int_A f_X(x) dx} = \int_B f_Y(y) dy = P(Y \in B)$$

Hence the event $\{Y \in B\}$ is independent of $\{X \in A\}$, $\forall A, B$

Remark: To verify independence, need to check the above identity for all x,y or all A,B .

To show non-independence, need to identify a pair (x,y) or (A,B) where identity not satisfied.

Lemma: Let $(X,Y) \sim f_{X,Y}$. Then $X \perp\!\!\!\perp Y$ i.f.f. \exists functions $g(x)$ and $h(y)$ s.t.

$$f_{X,Y}(x,y) = g(x)h(y) \forall x, y$$

Then, if $X \perp\!\!\!\perp Y$, then

$$\Rightarrow \forall A \subset \mathbb{R}, B \subset \mathbb{R} \\ P(X \in A, Y \in B) = P(X \in A)P(Y \in B)$$

2) if function $g(x) = \text{any of } x$
 $h(y) = \text{any of } y$

$$E[g(x)h(y)] = E[g(x)]E[h(y)]$$

Then, if $X \perp\!\!\!\perp Y$, w.l.o.g. M_x and M_y , then

$$\text{RV } Z := X+Y \text{ has the mgf} \\ M_Z(t) = M_X(t)M_Y(t) \quad \forall t$$

$$\begin{aligned} \mathbb{P}[M_Z(t) = E[e^{tZ}]] &= E[e^{tX+tY}] \\ &= E[e^{tX}]e^{tY} \\ &= E[e^{tX}]E[e^{tY}], \quad X \perp\!\!\!\perp Y \\ &= M_X(t)M_Y(t). \end{aligned}$$

Then, the sum of two independent normal RV's is again normal!

4.3: Bivariate Transformation

Let (X,Y) be a bivariate random vector

$$g: \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$g(x,y) := (g_1(x,y), g_2(x,y)) \in \mathbb{R}^2$$

Then $(U,V) := g(X,Y)$ is a bivariate random vector

$$\forall A \subset \mathbb{R}^2, P((U,V) \in A) = P((X,Y) \in g^{-1}(A))$$

$$\text{where } g^{-1}(A) = \{(x,y) | g_1(x,y) \in A\}$$

Discrete case: If (X,Y) is discrete, then so is (U,V)

$$f_{U,V}(u,v) = \sum_{\substack{x,y: \\ g_1(x,y)=u \\ g_2(x,y)=v}} f_{X,Y}(x,y)$$

CONTINUOUS CASE: change of var-formulae.

Let (X,Y) -continuous bivariate vector
 $C(X,Y) \sim f_{X,Y}$

Let $g: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ be a one-to-one mapping
i.e. $g^{-1}(u,v)$ has at most one element

Define $h(u,v) = g(X,Y)$; what is $f_{h(u,v)}$?

Let $A = \{ \mathbb{R}^2 \times \mathbb{R}^2 \mid f_{X,Y}(x,y) > 0 \}$, support of $f_{X,Y}$

$$B = g(A) = \{ \mathbb{R}^2 \times \mathbb{R}^2 \mid f_{X,Y}(x,y) > 0 \}$$

$$\text{Write } \begin{pmatrix} u \\ v \end{pmatrix} : g(x,y) = \begin{pmatrix} h_1(x,y) \\ h_2(x,y) \end{pmatrix}$$

Which has the inverse function

$$\begin{pmatrix} x \\ y \end{pmatrix} = h^{-1}(u,v) = g^{-1}(h_1(u,v), h_2(u,v))$$

Then the pdf $f_{h(u,v)}$ is given by the

$$\text{change-of-var formula: } f_{h(u,v)}(u,v) = f_{X,Y}(h_1(u,v), h_2(u,v)) / |J|$$

where J : determinant of the Jacobian matrix:

$$J = \begin{vmatrix} \frac{\partial h_1}{\partial x} & \frac{\partial h_1}{\partial y} \\ \frac{\partial h_2}{\partial x} & \frac{\partial h_2}{\partial y} \end{vmatrix} = \begin{vmatrix} \frac{\partial h_1}{\partial u} & \frac{\partial h_1}{\partial v} \\ \frac{\partial h_2}{\partial u} & \frac{\partial h_2}{\partial v} \end{vmatrix} = \frac{\partial h_1}{\partial u} \frac{\partial h_2}{\partial v} - \frac{\partial h_1}{\partial v} \frac{\partial h_2}{\partial u}$$

Under the transformation $(x,y) \rightarrow g(x,y)$ \mathbb{R}^2 is many-to-one.

We may partition the support of $f_{X,Y}$

$$A = \{ (x,y) \mid f_{X,Y}(x,y) > 0 \}$$

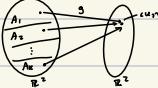
into disjoint subsets

$$A = A_1 \cup \dots \cup A_n$$

s.t. g is one-to-one from each A_i to $g(A_i)$, $i=1, \dots, n$

Prop: Let $h_i: h_1(u,v) \rightarrow h_2(u,v)$ be the inverse function of g restricted to domain $A_i \rightarrow g(A_i)$ and J_i the corresponding Jacobian. Then

$$f_{h(u,v)}(u,v) = \sum_{i=1}^n f_{X,Y}(h_1(u,v), h_2(u,v)) / |J_i|$$



4.4: Mixture & Hierarchical Models

-compute families of dist.'s in terms of simpler models of dist.'s, which are constructed in a hierarchical way

Then: If X and Y are any RV's, then	Proof: Suppose $(x,y) \sim f(x,y)$ -constant setting. Then
$EX = E[E(X Y)]$	$EX = \iint x f_{X,Y}(x,y) dx dy$ $= \iint x f_{X Y}(x y) f_Y(y) dx dy$ $= \int \left(\int x f_{X Y}(x y) dx \right) f_Y(y) dy$ $= E[X(Y)]$

$$\begin{aligned} \text{provided that the expectations exist.} \\ (\text{"law of iterated expectation"}) \\ \text{Proof: suppose } (x,y) \sim f(x,y) \text{ -constant setting. Then} \\ EX = \iint x f_{X,Y}(x,y) dx dy \\ = \iint x f_{X|Y}(x|y) f_Y(y) dx dy \\ = \int \left(\int x f_{X|Y}(x|y) dx \right) f_Y(y) dy \\ = \int E[X(Y)] f_Y(y) dy \\ = E[X(Y)]. \end{aligned}$$

Def: A RV X is said to have a mixture dist. if the dist. of X depends on a quantity that is also random, e.g. a Poisson-Binomial mixture

$$\begin{aligned} X|Y \sim \text{Binomial}(Y, p) \\ Y \sim \text{Poisson}(\lambda) \end{aligned} \quad \begin{cases} \text{mixture of Binomial distributions} \\ \text{mixing mechanism given by Poisson} \end{cases}$$

Adding more "models" to obtain arbitrarily complex distributions:

Prop:	
Suppose $Y \in \{1, \dots, K\}$, $P(Y=i) = p_i$, $i=1, \dots, K$	
and	
$X Y=i \sim \text{Normal}(\mu_i, \sigma_i^2)$	
then	
$f_X(x) = \sum_{i=1}^K p_i N(x \mu_i, \sigma_i^2)$	
$= \sum_{i=1}^K p_i \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2\sigma_i^2}(x-\mu_i)^2\right)$	
(a mixture of K normal components)	

If a prob. dist. for X may be obtained by multiple stages of cond. dist.'s, then we obtain a hierarchical model, e.g.

$$\begin{cases} X|Y \sim \text{Binomial}(Y, p) \\ Y|A \sim \text{Poisson}(\lambda) \\ A \sim \text{Exponential}(\beta) \end{cases}$$

where the randomness of A captures the variation across time (insect) mothers

THE UNCORRELATED VARIANCE FORMULA: For any two RV's X, Y ,

$$\text{Var}(X) = E[\text{Var}(Y|X)] + \text{Var}[E(Y|X)]$$

$$\begin{aligned} P. & \text{ Recall } \text{Var}(X) = E(X^2) - E^2(X) \\ \Rightarrow & \text{Var}(X) = E(X^2|Y) - E^2(X|Y) \quad (i) \\ \text{where } & E(E(X|Y)) = E(X) \\ \Rightarrow & E(\text{Var}(X|Y)) = E(X^2) - E(E(X|Y))^2 \quad (ii) \\ \text{thus } & \text{Var}(X) = (i) + (ii) = 0 \end{aligned}$$

4.3: COVARIANCE AND CORRELATION

Let (X, Y) be a bivariate vector $\{C(X), E(Y)\}$

$$\text{e.g. } (X, Y) = (\text{height}, \text{weight})$$

Want to make statement like:

"If X ↑ then Y will ↑ and vice versa"

Covariance: = cov of RV's X and Y

$$\text{cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$$

REMARKS: - If $X = Y$, then $\text{cov}(X, Y) = \text{var}(X) = E(X - E(X))^2$

- If $E(X) + E(Y) = 0$, then $\text{cov}(X, Y) = E(XY)$

$$\text{Correlation: } \text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}}$$

i.e. $\text{corr}(X, Y)$ = covariance of standardized versions of X and Y

$$\text{Then: } \text{cov}(X, Y) = E(XY) - E(X)E(Y)$$

$$\begin{aligned} \text{Pf: } \text{cov}(X, Y) &= E[(X - E(X))(Y - E(Y))] \\ &= E(XY - XE(Y) - YE(X) + E(X)E(Y)) \\ &\stackrel{L.O.E}{=} E(XY) - E(X)E(Y) - (EY)E(X) + E(Y)E(X) \\ &= E(XY) - E(X)E(Y) \end{aligned}$$

REMARKS: - generates $\text{var}(X) = E(X^2) - E^2(X)$

- conversely, $\text{cov}(X, Y) = 0$ does not imply $X \perp Y$

(unless (X, Y) is bivariate normal - see later)

Then: $\forall a, b \in \mathbb{R}$,

$$\text{Var}(aX + bY) = a^2\text{Var}(X) + b^2\text{Var}(Y) + 2ab\text{cov}(X, Y)$$

COROLLARY SCHWARTZ INEQUALITY: $E[XY] \leq \sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}$

$$\Leftrightarrow [E(XY)]^2 \leq E(X^2)E(Y^2)$$

Letting $U = X - E(X)$, $V = Y - E(Y)$,

$$[E(XY)]^2 \leq E(U^2)E(V^2)$$

$$[\text{cov}(X, Y)]^2 \leq \text{var}(X)\text{var}(Y)$$

$$\frac{|\text{cov}(X, Y)|}{\sqrt{\text{var}(X)\text{var}(Y)}} \leq 1$$

$$|\text{corr}(X, Y)| \leq 1$$

$$E[XY] \leq \sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}$$

Equality holds i.f.s. $Y = cX$ a.s., i.e.

$$P(Y = cX) = 1 \text{ (i.e. almost all)}$$

Letting $U = X - E(X)$, $V = Y - E(Y)$,

equally holds i.f.s. $P((Y - EY) = a(X - EX)) = 1$, a: constant a.s.

$$Y - EY = a(X - EX) \text{ a.s.}$$

$$\Rightarrow \text{cov}(X, Y) = E[(X - EX)(Y - EY)] = E[(X - EX)a(X - EX)] = aE[(X - EX)^2] = a\text{var}(X)$$

$$\text{var}(Y) = E[(Y - EY)^2] = E[a^2(X - EX)^2] = a^2\text{var}(X)$$

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}} = \frac{a}{\sqrt{a^2\text{var}(X)}} = \frac{a}{|a|} = \begin{cases} 1, & a > 0 \\ -1, & a < 0 \end{cases}$$

$$\Rightarrow |\text{corr}(X, Y)| = 1$$

Then. Suppose $\text{var}(X) = \text{var}(Y) = 1$.

(i) If $\text{corr}(X, Y) = 1$, then there are constants b, c s.t.

$$Y = bX + c \text{ w.p. 1}$$

(ii) If $\text{corr}(X, Y) = -1$, then

$$Y = bX + c \text{ w.p. 1}$$

for some constants $b, c, 0$.

REMARKS:



If $X \perp Y$, then

$$\text{cov}(X, Y) = E(XY) - E(X)E(Y) = 0$$

$$\text{corr}(X, Y) = 1 \Rightarrow Y = bX + c, b \neq 0$$

Remark: $\text{cov}(x,y) \geq 0$ does not imply $X \perp\!\!\! \perp Y$.

possible that X and Y are strongly dependent on one another (suggests nonlinear relationship)

Bivariate Normal Distribution

Def: $Z = (X,Y) \sim N(\mu, \Sigma)$ $\left| \begin{array}{l} \mu \in \mathbb{R}^2 \\ \Sigma \text{ positive definite, symmetric } \in \mathbb{R}^{2 \times 2} \end{array} \right.$

$$f_{X,Y}(x,y) := f_Z(z) := \frac{1}{(2\pi)^2 |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (z - \mu)^T \Sigma^{-1} (z - \mu) \right\}$$

Remark: $\mu = (\mu_1, \mu_2)^T$, $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$

μ : mean vector, Σ : covariance matrix of Z :
 $\left\{ \begin{array}{l} E\bar{Z} = \mu \\ \text{cov}\bar{Z} = \Sigma \end{array} \right.$

$E\bar{Z} = \mu$ means $E \begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$ $\left\{ \begin{array}{l} EX = \mu_1 \\ EY = \mu_2 \end{array} \right.$

$\text{cov}\bar{Z} = \Sigma$ means $E[(Z - E\bar{Z})(Z - E\bar{Z})^T] = \Sigma$

$$\Leftrightarrow E \begin{pmatrix} X - E\bar{X} \\ Y - E\bar{Y} \end{pmatrix} \begin{pmatrix} X - E\bar{X} \\ Y - E\bar{Y} \end{pmatrix}^T = \Sigma$$

$$\Leftrightarrow E \begin{bmatrix} (X - E\bar{X})^2 & (X - E\bar{X})(Y - E\bar{Y}) \\ (Y - E\bar{Y})(X - E\bar{X}) & (Y - E\bar{Y})^2 \end{bmatrix} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

$$\Leftrightarrow \begin{cases} \Sigma_{11} = \text{var} X \\ \Sigma_{21} = \text{cov}(X,Y) \\ \Sigma_{12} = \Sigma_{21} = \text{cov}(Y,X) \\ \Sigma_{22} = \text{var} Y \end{cases} \quad \text{Hence } \text{cov}\bar{Z} = E[(Z - E\bar{Z})(Z - E\bar{Z})^T]$$

a remarkable property of the normal dist.

due to the identities:

$$\text{if } f_{X,Y}(x,y) = f_Z(z) \Rightarrow \frac{1}{\sqrt{(2\pi)^2 |\Sigma|}} \exp \left\{ -\frac{1}{2} (z - \mu)^T \Sigma^{-1} (z - \mu) \right\}$$

then:

$$\begin{aligned} \iint_{\mathbb{R}^2} f_{X,Y}(x,y) dx dy &= 1 \\ \iint_{\mathbb{R}^2} x f_{X,Y}(x,y) dx dy &= \mu_1 \\ \iint_{\mathbb{R}^2} y f_{X,Y}(x,y) dx dy &= \mu_2 \\ \iint_{\mathbb{R}^2} f_{X,Y}(x,y) dx dy &= \text{tr} \Sigma \\ \iint_{\mathbb{R}^2} (z - \mu)^T (z - \mu) f_{X,Y}(x,y) dx dy &= \Sigma \end{aligned}$$

Remark: follows from properties of (multivariate) normal dist.

translating #5 into univariate forms (don't need to remember):

$$f_{X,Y}(x,y) = \frac{1}{\sqrt{2\pi \sigma_x \sigma_y \sqrt{1 - \rho_{xy}^2}}} \exp \left\{ -\frac{1}{2} \left[\frac{(x - \mu_1)^2}{\sigma_x^2} - 2\rho_{xy} \left(\frac{x - \mu_1}{\sigma_x} \right) \left(\frac{y - \mu_2}{\sigma_y} \right) + \frac{(y - \mu_2)^2}{\sigma_y^2} \right] \right\}$$

Then if $\text{cov}(X,Y) = 0$ and
 $(X,Y) \sim \text{bivariate Normal}$
then $X \perp\!\!\! \perp Y$

$P_{\text{cov}(X,Y) = 0} \Rightarrow \rho_{xy} = 0$
From (4), $f_{X,Y}(x,y)$ factorizes into
a product of $g(x)$ and $h(y)$, so $X \perp\!\!\! \perp Y$.

Fact: If $Z = (X,Y) \sim N(\mu, \Sigma)$ $\left| \begin{array}{l} \mu \in \mathbb{R}^2, \mu \neq \emptyset \\ \Sigma \in \mathbb{R}^{2 \times 2} \end{array} \right.$

then the marginal of X and Y are normal too:

$$X \sim N(\mu_1, \Sigma_{11})$$

$$Y \sim N(\mu_2, \Sigma_{22})$$

$$\text{and } \text{cov}(X,Y) = \Sigma_{12} = \Sigma_{21}$$

the conditional distributions are also normal:

$$Y|X=x \sim N(E(Y|X=x), \text{var}(Y|X=x))$$

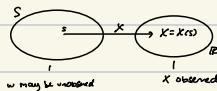
$$\text{where } E(Y|X=x) = \mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (x - \mu_1)$$

$$\text{var}(Y|X=x) = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$$

4.6: Multivariate distributions

Def: A random vector denoted $X = (X_1, \dots, X_n)$ is a function that maps elements of a

Sample space S into \mathbb{R}^n



$$X(s) = (X_1(s), \dots, X_n(s)) \in \mathbb{R}^n$$

NOW, we may specify probabilities, for $A \subset \mathbb{R}^n$:

$$P(X \in A) := P(s \in S : X(s) \in A)$$

If X takes countably many values, then we say X is a discrete RV, which is associated w/ a

Joint probability mass function (pmf) := $f(x) = f(x_1, \dots, x_n) = P(X_1 = x_1, \dots, X_n = x_n)$

$$\text{Hence for } A \subset \mathbb{R}^n, P(X \in A) = \sum_{x \in A} f(x)$$

Similarly, pdf := $f(x) = f(x_1, \dots, x_n)$

$$\text{where } P(X = (x_1, \dots, x_n) \in A) = \int_A f(x) dx = \iint_{\substack{\text{region } A \\ x_1 = x_1, \dots, x_n = x_n}} f(x_1, \dots, x_n) dx_1 \dots dx_n$$

$$\text{Expectation} := E[g] = \begin{cases} \int_{\mathbb{R}^n} g(x) f(x) dx & , X-\text{cont.} \\ \sum_{x \in \text{range}(X)} g(x) p(x) & , X-\text{discrete} \end{cases}$$

$$\text{marginal dist.} := f_{X_1, \dots, X_k} (x_1, \dots, x_k) = \begin{cases} \int_{\mathbb{R}^{n-k}} \dots \int_{\mathbb{R}^k} f_{X_1, \dots, X_n} (x_1, \dots, x_n) dx_{k+1} \dots dx_n \\ \sum_{x_{k+1}, \dots, x_n} f_{X_1, \dots, X_n} (x_1, \dots, x_n) \end{cases}$$

$$\text{e.g.: } (x_1, \dots, x_n) \mapsto f_{X_1, \dots, X_n}(x_1, \dots, x_n) \\ f_{X_1, X_2}(x_1, x_2) = \int_{\mathbb{R}^{n-2}} f_{X_1, X_2, X_3, \dots, X_n}(x_1, x_2, x_3, \dots, x_n) dx_3 \dots dx_n$$

$$\text{conditional dist.} := f_{X_1, \dots, X_n | X_1, \dots, X_k} (x_1, \dots, x_n | x_1, \dots, x_k) = \frac{f_{X_1, \dots, X_n}(x_1, \dots, x_n, x_1, \dots, x_k)}{f_{X_1, \dots, X_k}(x_1, \dots, x_k)}$$

REMARKS: admits factorization: $f_{X_1, \dots, X_n}(x_1, \dots, x_n) = f_{X_1, \dots, X_k}(x_1, \dots, x_k) f_{X_{k+1}, \dots, X_n}(x_{k+1}, \dots, x_n)$

closed formulae: $f_{X_1, \dots, X_n}(x_1, \dots, x_n) = f_{X_1}(x_1) f_{X_2, \dots, X_n}(x_2, \dots, x_n) f_{X_2 | X_1}(x_2 | x_1) \dots f_{X_n | X_1, \dots, X_{n-1}}(x_n | x_1, \dots, x_{n-1})$

$$\text{e.g. } f_{X,Y}(x,y) = f_{X|Y}(x|y) f_{Y|X}(y|x) = f_X(x) f_{Y|X}(y|x) f_{X|Y}(x|y)$$

Prop. Suppose $\begin{cases} X, Y \text{-discrete} \\ Z, T \text{-continuous} \end{cases}$ and let $f_{X,Y,Z,T}(x,y,z,t) = \underbrace{f_{XY}(x,y)}_{\text{part}} \underbrace{f_{Z|XY}(z,t|xy)}_{\text{conditional pdf}}$

$$= \underbrace{f_{X|Y}(x|y)}_{\text{part}} \underbrace{f_{Y|Z,T}(y|z,t)}_{\text{conditional pdf}}$$

Then, $\forall A, B, C, D \subset \mathbb{R}$,

$$P(X \in A, Y \in B, Z \in C, T \in D) = \sum_{x \in A} \sum_{y \in B} \int_C \int_D f_{X,Y,Z,T}(x,y,z,t) dz dt$$

Bernoulli(nominal dist.): $\begin{cases} n \text{ independent trials} \\ n \text{ possible outcomes w/} \\ \text{"cell prob." } (p_1, \dots, p_n) : \sum_i p_i = 1 \end{cases}$

$X_1 + \dots + X_n = m$

$X_i : i \text{th outcome}$

Now, joint pmf for $X_1, \dots, X_m \sim N^m$, $\sum_{i=1}^m X_i = m$

$$\begin{aligned} f(x_1, \dots, x_m) &= \left(\frac{1}{\sqrt{2\pi}} \right)^m p_1^{x_1} \cdots p_m^{x_m} \\ &= \frac{m!}{x_1! \cdots x_m!} p_1^{x_1} \cdots p_m^{x_m} \end{aligned}$$

and 0 otherwise

Remarks: By multinomial formula

$$\sum_{x_1, \dots, x_m} (x_1, \dots, x_m) p_1^{x_1} \cdots p_m^{x_m} = (p_1 + \cdots + p_m)^m = 1$$

$$\text{marginal distribution: } f_{X_1}(x_1) = \frac{m!}{x_1!(m-x_1)!} p_1^{x_1} (1-p_1)^{m-x_1}$$

$\Rightarrow X_1 \sim \text{Binomial}(m, p_1)$

$$\text{conditional distribution: } f_{X_1, \dots, X_m | X_m=k_m}(x_1, \dots, x_{m-1}, k_m) = \frac{(m-k_m)!}{x_1! \cdots x_{m-1}!} \left(\frac{p_1}{1-p_1} \right)^{x_1} \cdots \left(\frac{p_{m-1}}{1-p_{m-1}} \right)^{x_{m-1}}$$

$$\Rightarrow X_1, \dots, X_{m-1} | X_m=k_m \sim \text{Multinomial}(m-k_m \text{ trials, cell prob } \left(\frac{p_1}{1-p_1}, \dots, \frac{p_{m-1}}{1-p_{m-1}} \right))$$

Notice: kind of invariance property by marginalization and conditioning

this holds for quite a few families of distributions (often in exponential family)

for continuous distributions, one such family is the multivariate normal

Recall (independence): X and Y independent if $\forall x, y \in \mathbb{R}$, $f_{X,Y}(x,y) = f_X(x)f_Y(y)$

Def: Let (X_1, \dots, X_n) be a random vector w/ joint pdf/pmf f_{X_1, \dots, X_n} .
Let $f_{X_i}(x_i)$ be the marginal pdf of X_i .

Then X_1, \dots, X_n are mutually independent R.V.'s if $\forall (x_1, \dots, x_n)$

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i)$$

Remarks: each X_i may be itself a random vector

Many properties established for bivariate random vector can be extended naturally to the multivariate setting

Theorem: If X_1, \dots, X_n are continuous (independent) R.V.'s,
let $g_1(x_1), \dots, g_n(x_n)$ be real-valued functions on the domain
of X_1, \dots, X_n respectively, then

$$E[g_1(x_1) \cdots g_n(x_n)] = E[g_1(x_1)] \cdots E[g_n(x_n)]$$

Theorem: If X_1, \dots, X_n - mutually independent R.V.'s, let $\mathbf{z} := z_1, z_2, \dots, z_n$
then $M_{\mathbf{z}}(\mathbf{c}) = M_{X_1}(c_1) \cdots M_{X_n}(c_n)$

Corollary: in particular, if X_1, \dots, X_n i.i.d., then $M_{\mathbf{z}}(\mathbf{c}) = (M_{X_1}(c_1))^n$

Properties

\Rightarrow If $X_i \sim \text{Gamma}(k_i, \beta)$, X_i - independent.

Then $X_1 + \cdots + X_n \sim \text{Gamma}(k_1 + \cdots + k_n, \beta)$

\Rightarrow If $X_i \stackrel{\text{ind}}{\sim} N(\mu_i, \sigma_i^2)$, let $a_i, b_i \in \mathbb{R}$.

$$\begin{aligned} \text{then } \mathbf{z} := \sum_{i=1}^n (a_i X_i + b_i) &\sim N\left(\sum_{i=1}^n (a_i \mu_i + b_i), \sum_{i=1}^n a_i^2 \sigma_i^2\right) \\ &= \mathbf{a}^T \mathbf{x} + \mathbf{b}^T \mathbf{1} \end{aligned}$$

Lemma: Suppose $(X_1, \dots, X_n) \sim f_{X_1, \dots, X_n}$.

X_1, \dots, X_n are mutually independent i.e.

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = g_1(x_1) \cdots g_n(x_n)$$

for some function g_1, \dots, g_n

Theorem: If X_1, \dots, X_n are mutually independent,

let $g(x, \mathbf{z})$: function of $X_1, \dots, X_n, \mathbf{z}$

then $g(X_1, \dots, X_n, \mathbf{z})$ are also mutually independent

Proof: Similar to bivariate case

Proof: Follows directly from the factorization of the joint pdf/pmf

Note: the change-of-variable formula can be extended to the multivariate case too (See textbook)

5.1: iid samples

Def: Suppose $X_1, \dots, X_n \stackrel{iid}{\sim} f$ (pmf or pmf).

Then we say (X_1, \dots, X_n) is a random sample from the population with pmf/pdf f .

Remarks: Common variations:

- n-sample of F
- n/iid sample of F
- A-sample of size n

Def: Let (x_1, \dots, x_n) be a n-sample from a population.

Let $T(x_1, \dots, x_n)$ be a real-valued function.

Then $Y = T(x_1, \dots, x_n)$ is called a statistic,
i.e. a function of a random sample

Remarks: Statistic is a function

- telling us something about an underlying population via pmf/pdf F.
- does so only on random samples created

Modern viewpoint: $x_1, \dots, x_n \xrightarrow{\text{function}} T \rightarrow Y = T(x_1, \dots, x_n)$

algorithm/application

Examples: 1) $\bar{X} := \frac{1}{n}(X_1 + \dots + X_n)$; sample mean

$$\Rightarrow S^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

$$2) X_{(1)} := \min\{x_1, \dots, x_n\}$$

$$X_{(2)} := \min\{x_1, \dots, x_n \setminus X_{(1)}\}$$

⋮

$$X_{(k)} := \min\{x_1, \dots, x_n \setminus X_{(1)}, \dots, X_{(k-1)}\}$$

i.e. the kth smallest number of the sample

order statistics of size n-sample := $(X_{(1)}, \dots, X_{(n)})$

Then: Let (X_1, \dots, X_n) be an n-iid sample from a population with mean μ and variance $\sigma^2 < \infty$

$$\begin{cases} E\bar{X} = \mu \\ \text{Var}(\bar{X}) = \sigma^2/n \\ E S^2 = \sigma^2 \end{cases}$$

Remarks: provides statistical justification for using \bar{X} and S^2 as estimates of μ and σ^2 respectively; \bar{X} and S^2 are unbiased estimates

Distribution of \bar{X}

two main methods: { method of moments
2) change-of-var formula
(convolution formula)

$$\begin{aligned} \text{1) } M_{\bar{X}}(t) &= E e^{t\bar{X}} \\ &= E e^{\frac{t}{n}(X_1 + \dots + X_n)} \\ &= E e^{\frac{t}{n}X_1} \dots e^{\frac{t}{n}X_n} \\ &= (M_X(t/n))^n \end{aligned}$$

Remark: If M_X can be recognized as MGF of a known family,
then we can find the dist. of \bar{X} easily, e.g.

$$\begin{aligned} \text{if } X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2) \\ \text{then } \bar{X} \stackrel{iid}{\sim} N(\mu, \sigma^2/n) \end{aligned}$$

2) CONVOLUTION Formula

Then: If $X \sim f_X$, $Y \sim f_Y$,
then $Z = X+Y$ has the pdf which is the convolution
of f_X and f_Y :

$$\begin{aligned} f_Z(z) &= f_X(z-y) f_Y(y) \\ &:= \int_{-\infty}^{\infty} f_X(z-y) f_Y(y) dy \end{aligned}$$

Proof: use the change-of-var formula for the mapping $(X, Y) \rightarrow (X+Y, Y)$

Remarks: useful for deriving pdf/pmf when the transformation $T(X_1, \dots, X_n)$ does not belong to a known/well-recognized family (e.g. exponential, location-scale, etc.)

5.2: USEFUL CLASSICAL FACTS

LEM. IF $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$,

$$\text{let } \bar{X} = \frac{1}{n}(X_1 + \dots + X_n)$$

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

then

a) \bar{X} and S^2 are independent RV's

b) $\bar{X} \sim N(\mu, \sigma^2/n)$

c) $(n-1) \frac{S^2}{\sigma^2} \sim \chi^2_{n-1}$

Proof (c): seen in previous section

Recall: χ_p^2 w/p df:

$$f(x) = \frac{1}{\Gamma(p/2)^{1/2}} x^{p/2-1} e^{-x/2}, x > 0$$

and $\chi_p^2 \equiv \text{Gamma}(p/2, 1)$

OTHER STATISTICS AND THEIR DISTRIBUTIONS

IF $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, μ unknown

then \bar{X} is a statistic that tells us about μ

Q: how so? we know $\frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \sim N(0, 1)$

this gives us a way to quantify the uncertainty about μ , using statistic \bar{X} as an estimate

Q: what if σ is unknown too?

we may want to consider $\frac{\bar{X}-\mu}{S/\sqrt{n}}$

since S is a statistic that can be obtained from the sample

to be useful, we need to know the distribution of

$$T = \frac{\bar{X}-\mu}{S/\sqrt{n}}$$

LEM. $T \sim t_{n-1}$: student's t dist. w/ $n-1$ df

tp has the pdf

$$f_T(t) = \frac{\Gamma((n+1)/2)}{\Gamma(n/2)} \frac{1}{(t^2/(n-1))^{n/2}} \frac{1}{(1+t^2/(n-1))^{(n+1)/2}} = \text{student's t}$$

Remarks: t is the ratio of two independent random variables as

$$\bar{X}-\mu \sim \text{Normal}, \sqrt{(n-1)S^2} \sim \sqrt{\chi^2_{n-1}}$$

If $n=2$, then $T \sim \text{Ratio of 2 ind. Normal RV's} \Rightarrow T \sim \text{Cauchy}$

5.3: CONVERGENCE CONCEPTS

Given a sequence of RV's X_1, \dots, X_n , we want to study various notions of convergence to a rv X :

$\Leftrightarrow X_n \xrightarrow{P} X$ convergence in probability

$\Leftrightarrow X_n \xrightarrow{a.s.} X$ convergence almost surely (w/ prob 1)

$\Leftrightarrow X_n \xrightarrow{d} X$ convergence in dist.

we already encountered: $X_n \xrightarrow{d} X$ if $F_{X_n}(x) \rightarrow F_X(x)$ at all points where F_X is continuous

DEF: $X_n \xrightarrow{P} X$ /f

$$\forall \epsilon > 0, P(X_n - X \geq \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty$$

Notice: $\forall \epsilon > 0, P(X_n - X \geq 2\epsilon) \rightarrow 0$ as $n \rightarrow \infty$

$$= P(|X_n - X| \geq \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty$$

where S generally depends on n

Theorem (Law of Large Numbers - LLN): Let X_1, \dots, X_n be RV's with $E[X_i] = \mu$ and $\text{Var}[X_i] < \infty$

Define $\bar{X}_n = \frac{1}{n} (X_1 + \dots + X_n)$

Then $\bar{X}_n \xrightarrow{P} \mu$

Def: $X_n \xrightarrow{a.s.} X$ (w.r.t prob. 2) if

$P(\lim_{n \rightarrow \infty} |X_n - x| = 0) = 1 \equiv P(\lim_{n \rightarrow \infty} X_n = x) = 1$

Remark: Equivalently def: $\forall \epsilon > 0, P(\lim_{n \rightarrow \infty} |X_n - x| < \epsilon) = 1$

Notice: $X_n \xrightarrow{a.s.} X$ if

$P(\{\omega : X_n(\omega) \rightarrow X(\omega)\}) = 1$

where S: states is dependent on n



where, $X_n(\omega) \rightarrow X(\omega) \Leftrightarrow S^n(\omega) \rightarrow S^1(\omega) \Rightarrow P(S^n \rightarrow S^1) = 1$

Prop: If $X_n \xrightarrow{a.s.} X$ then $X_n \xrightarrow{P} X$.

If $X_n \xrightarrow{P} X$ then $X_n \xrightarrow{a.s.} X$

Theorem (Law of Large Numbers - SLN): Let X_1, \dots, X_n be RV's with $E[X_i] = \mu$ and $\text{E}[X_i^2] < \infty$

Define $\bar{X}_n = \frac{1}{n} (X_1 + \dots + X_n)$

Then $\bar{X}_n \xrightarrow{a.s.} \mu$

Proof: Similar in spirit but technically more involved than LLN

Note: Condition $\text{E}[X_i^2] < \infty$ is very mild

Theorem (CLT): If $X_1, X_2, \dots \xrightarrow{a.s.}$ RV whose MGF exist in a neighborhood of 0.

Let $\left(\frac{a_n - E[X_i]}{\sigma_n}, \frac{a_n^2 - E[X_i^2]}{\sigma_n^2} \right)$ and $Z = N(0, 1)$

then $\left(\frac{a_n - E[X_i]}{\sigma_n} \right) = \sqrt{n} \left(\frac{a_n - E[X_i]}{\sigma_n} \right) \xrightarrow{d} Z$

Remark: · perhaps most celebrated thm. in prob.
· widely applicable; only $\text{Var}[X_i] < \infty$ is required

Corollary: Let $Y_i = cX_i - m > 0$. Then $E[Y_i^2] < \infty$

and we may write $\frac{1}{\sqrt{n}} (Y_1 + \dots + Y_n) \xrightarrow{d} N(0, 1)$.

$$P\left(\sqrt{n} \left(\frac{a_n - E[X_i]}{\sigma_n} \right) = \sqrt{n} \left(\frac{Y_1 + \dots + Y_n - m}{\sigma_n} \right) = \frac{1}{\sqrt{n}} (Y_1 + \dots + Y_n) \xrightarrow{d} N(0, 1)\right).$$

· convergence (a.s., in prob., in dist.) preserves in nice transformations:

Lemma (1): If $\begin{cases} X_n \xrightarrow{a.s.} X \\ X_n \xrightarrow{a.s.} b: \text{constant} \end{cases}$

then, $aX_n + b \xrightarrow{a.s.} aX + b$

Lemma (2): If $\begin{cases} X_n \xrightarrow{P} X \\ X_n \xrightarrow{P} b: \text{constant} \end{cases}$

then, $aX_n + b \xrightarrow{P} aX + b$

Lemma (3): If $\begin{cases} X_n \xrightarrow{a.s.} X \\ X_n \xrightarrow{P} b: \text{const.} \end{cases}$

then, $aX_n + b \xrightarrow{a.s.} aX + b$

Remark: · L1 and L2 are almost immediate from definition

· L3 is known as Slutsky's theorem; P: based on characterization: $X_n \xrightarrow{P} X$

\Leftrightarrow for all continuous and bounded function $f(x)$

$E[f(X_n)] \rightarrow E[f(x)]$ as $n \rightarrow \infty$ (higher-dim def. as opposed to 1-dim - not covered)