

5420 Theory Notes

Note: Finished Lecture 31

Introduction

Calculus Review

Basic Integrals

$$\int x^m dx = \frac{1}{m+1} x^{m+1},$$

$$\int e^x dx = e^x.$$

Integration by Parts

$$\int u dv = uv - \int v du.$$

Gamma Integral

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

Derivatives

Chain rule:

$$(f(g(x)))' = g'(x)f'(g(x)).$$

Exponential:

$$(e^{f(x)})' = f'(x)e^{f(x)}.$$

Logarithm:

$$(\log f(x))' = \frac{f'(x)}{f(x)}.$$

Quotient rule:

$$\left(\frac{f(x)}{g(x)} \right)' = \frac{f'(x)g(x) - g'(x)f(x)}{g(x)^2}.$$

Power rule:

$$(x^m)' = mx^{m-1}.$$

Function powers:

$$(f(x))^{m'} = m f'(x) [f(x)]^{m-1}.$$

Even and Odd Functions

- A function f is **even** if

$$f(x) = f(-x),$$

e.g., $f(x) = x^2$.

- A function f is **odd** if

$$f(-x) = -f(x),$$

e.g., $f(x) = x^3$, $f(x) = \sin x$.

For integrals:

- If f is even,

$$\int_{-a}^a f(x) dx = 2 \int_0^a f(x) dx.$$

- If f is odd,

$$\int_{-a}^a f(x) dx = 0.$$

Trigonometric Derivatives

$$(\sin u)' = u' \cos u,$$

$$(\cos u)' = -u' \sin u.$$

Beta Function

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}.$$

Equivalently,

$$B(a, b) = \int_0^1 x^{a-1} (1-x)^{b-1} dx = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}.$$

Probability

- **Probability** is a branch of mathematics concerned with the study of *random phenomenon* (e.g., experiments, models of populations).
- We are primarily interested in probability as it relates to **statistical inference**, the science of drawing inferences about populations based on only a part of the population (i.e., a sample).

Some Definitions

1. **population:** the entire set of objects that we are interested in studying
e.g., all ISU students
2. **sample:** the subset of the population available for observation
e.g., STAT 542 students

Note: population and sample are crucial terms in understanding statistics (i.e., STAT 543), but will not occur very often in our discussions of probability theory (i.e., STAT 542).

3. **experiment:** process of obtaining an observed result of a random phenomenon
4. **sample space S :** the set of all possible outcomes of the experiment
 - elements $s \in S$ of a sample space are called **sample points** (s)
 - a sample space may be
 - **discrete**
(finite or countably infinite, i.e., listable as a finite/infinite sequence)

$$S = \{s_1, s_2, \dots, s_n\}$$

or

$$S = \{s_1, s_2, s_3, \dots\}$$

- or **continuous**
(uncountably infinite, i.e., a continuum of sample points like
 $S = [0, \infty)$)

5. **event** (e.g., A, B, \dots): subset of the sample space S

- **set:** A is a collection of elements
(in our case, A is a collection of outcomes)
- **membership:** $x \in A$ or $x \notin A$
(x is in A or x is not in A)
- **complement:**

$$A^c = \{x : x \notin A\}$$

(x such that x is not in A)

- **union:**

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

(x is in A or B or both)

- **intersection:**

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

- **subset:** $A \subset B$ means that A is contained in B
(formally, $x \in A \Rightarrow x \in B$)

- **equality:** $A = B$ if $A \subset B$ and $B \subset A$

- **empty set:** \emptyset

Algebraic Laws

- **commutativity:**

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

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

- **associativity:**

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

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

- **distributive law:**

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

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

- **DeMorgan's laws:**

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

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

Aside on disjoint and partitions

- events A and B are **disjoint** (mutually exclusive) if

$$A \cap B = \emptyset$$

- For a sequence A_1, A_2, \dots of events, we say A_1, A_2, \dots are **pairwise disjoint** if

$$A_i \cap A_j = \emptyset \quad \text{for all } i \neq j$$

- A_1, A_2, \dots is a **partition** of S if the A_i 's are pairwise disjoint and exhaustive, that is,

$$\bigcup_{i=1}^{\infty} A_i = S \quad \text{and} \quad A_i \cap A_j = \emptyset \quad \text{for all } i \neq j$$

Probability Functions

- A **probability function** is a function P defined on a Borel field \mathcal{B} of the sample space S that satisfies:
 1. $P(A) \geq 0$ for all $A \in \mathcal{B}$
 2. $P(S) = 1$
 3. If $A_1, A_2, \dots \in \mathcal{B}$ are *pairwise disjoint*, then

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

- Any function satisfying the above is a legitimate probability function.

Theorem 1.2.8.

If P is a probability function and A is any set in \mathcal{B} , then:

(a)

$$P(\emptyset) = 0$$

(b)

$$P(A) \leq 1$$

(c)

$$P(A^c) = 1 - P(A)$$

Proof of (c) (parts (a) and (b) follow from (c) and the axioms):

Since

$$S = A \cup A^c,$$

and A and A^c are disjoint, by the axioms of probability,

$$P(S) = P(A \cup A^c) = P(A) + P(A^c).$$

Because $P(S) = 1$, we have

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

which implies

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

Theorem 1.2.9.

If P is a probability function and A, B are sets in \mathcal{B} , then:

(a)

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

(b)

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

(c) If $A \subset B$, then

$$P(A) \leq P(B).$$

Theorem 1.2.11.

If P is a probability function, then

(a) For any partition $C_1, C_2, \dots \in \mathcal{B}$ (i.e., disjoint C_i 's and $\bigcup_{i=1}^{\infty} C_i = S$),

$$P(A) = \sum_{i=1}^{\infty} P(A \cap C_i).$$

(b) For any sets $A_1, A_2, \dots \in \mathcal{B}$,

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

Principle of Inclusion–Exclusion.

For any sets A_1, \dots, A_n ,

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{k=1}^n (-1)^{k-1} \left(\sum_{1 \leq i_1 < \dots < i_k \leq n} P(A_{i_1} \cap \dots \cap A_{i_k}) \right).$$

Equivalently,

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

This generalizes

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

and is proven by induction.

Bonferroni's Inequalities.

For any sets A_1, \dots, A_n and any $m \in \{1, \dots, n\}$,

- if m is odd,

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{k=1}^m (-1)^{k-1} \left(\sum_{1 \leq i_1 < \dots < i_k \leq n} P(A_{i_1} \cap \dots \cap A_{i_k}) \right),$$

- if m is even,

$$P\left(\bigcup_{i=1}^n A_i\right) \geq \sum_{k=1}^m (-1)^{k-1} \left(\sum_{1 \leq i_1 < \dots < i_k \leq n} P(A_{i_1} \cap \dots \cap A_{i_k}) \right).$$

In particular,

$$\sum_{i=1}^n P(A_i) - \sum_{1 \leq i < j \leq n} P(A_i \cap A_j) \leq P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i).$$

Combinatorics

Permutations / ordered arrangements II.

When selecting r objects from n objects (without replacement), the number of ordered arrangements possible is

$$n(n-1)\cdots(n-r+1) = \frac{n!}{(n-r)!}.$$

Combinations / unordered selections.

The number of ways to choose r objects from n objects (without replacement), where the ordering doesn't matter, is

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

Summary table: number of ways to select r objects from a group of n

	objects chosen without replacement	objects chosen with replacement
ordered	$\frac{n!}{(n-r)!}$	n^r
unordered	$\binom{n}{r}$	$\binom{n+r-1}{r}$

Conditional Probability

- **Definition:** If A, B are events in S with $P(B) > 0$, then

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

- In conditioning, B can be thought of as the **updated sample space**, i.e., not all of S is relevant since we know B has occurred.

$P(\cdot | B)$ is a probability function that satisfies the usual axioms and properties.

Axioms:

- $P(A | B) \geq 0$ for all events A
- $P(B | B) = 1$
(B is the updated sample space)
- If A_1, A_2, \dots are pairwise disjoint events, then

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

Some properties:

$$P(A^c \mid B) = 1 - P(A \mid B)$$

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

It also follows from our definition of conditional probability that

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

More generally, for events A_1, A_2, \dots, A_n ,

$$P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1) P(A_2 \mid A_1) P(A_3 \mid A_1 \cap A_2) \dots P(A_n \mid A_1 \cap \dots \cap A_{n-1}).$$

It is possible to reverse the conditioning of A and B to obtain **Bayes' rule**:

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

More generally, if A_1, A_2, \dots is a partition of the sample space S , then we obtain a general version of Bayes' rule:

$$P(A_i \mid B) = \frac{P(B \mid A_i) P(A_i)}{\sum_{j=1}^{\infty} P(B \mid A_j) P(A_j)}.$$

Independence

If $P(A \mid B) = P(A)$, then the occurrence of B does not affect the probability of A . It then follows that

$$P(A \cap B) = P(A)P(B) \quad \text{and} \quad P(B \mid A) = P(B).$$

We define two events A and B as **independent** if

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

More than two events.

A_1, \dots, A_n are **independent** if and only if, for any subcollection $\{i_1, \dots, i_k\} \subset \{1, \dots, n\}$ of distinct indices (with any $2 \leq k \leq n$), it holds that

$$P\left(\bigcap_{j=1}^k A_{i_j}\right) = \prod_{j=1}^k P(A_{i_j}).$$

- If A_1, \dots, A_n are independent, then

$$P(A_i \cap A_j) = P(A_i)P(A_j) \quad \text{for any } i \neq j.$$

- However,

$$P(A_i \cap A_j) = P(A_i)P(A_j) \text{ for } i \neq j$$

does **not** imply that A_1, \dots, A_n are independent.

If A_1, \dots, A_n are independent, then

$$P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1)P(A_2) \cdots P(A_n).$$

However,

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

holding does **not** imply that A_1, \dots, A_n are independent.

The assumption of independence of events allows the computation of joint occurrences of events through simple calculations.

Random Variables

Definition: A **random variable** (r.v.) X is a function defined on a sample space S that associates a real number with each outcome in S .

That is, for each $s \in S$, we have

$$X(s) \in \mathbb{R}.$$

In function notation,

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

We usually suppress the dependence of X on $s \in S$ and write

$$X = X(s).$$

We have $P(A)$ defined on events $A \subset S$, which can be used to assign probabilities for events concerning a random variable X on \mathbb{R} ($X : S \rightarrow \mathbb{R}$).

Define $P_X(\cdot)$ for events $B \subset \mathbb{R}$ as follows:

$$P_X(B) = P_X(X \in B) = P(\{s \in S : X(s) \in B\}).$$

$P_X(\cdot)$ satisfies the axioms and is therefore a legitimate probability function.

CDF

Definition.

The **cumulative distribution function** (cdf) of a random variable X , denoted by $F(\cdot)$, is defined by

$$F(x) = P(X \leq x), \quad x \in \mathbb{R}.$$

Sometimes written with subscript as $F_X(x)$.

A function $F(x)$, $x \in \mathbb{R}$, is a cdf for some random variable if and only if the following hold:

1. $F(x)$ is a nondecreasing function of x .

2.

$$\lim_{x \rightarrow -\infty} F(x) = 0 \quad \text{and} \quad \lim_{x \rightarrow \infty} F(x) = 1.$$

3. $F(x)$ is right continuous, i.e.,

$$\lim_{x \downarrow x_0} F(x) = F(x_0) \quad \text{for any } x_0 \in \mathbb{R}.$$

Discrete Random Variables

Definition.

If a cdf F is a step function (with jumps at a countable collection of points $x_i \in \mathbb{R}$), then we say the distribution described by F is **discrete** (with support or range $x_i \in \mathbb{R}$).

If a random variable X has a cdf $F = F_X$ which is a step function, then we say X is a **discrete random variable**.

Besides the cdf, there are other (equivalent) ways to state the probability distribution for a discrete distribution / discrete r.v. X .

1. Probability mass function (pmf).

The pmf of a discrete random variable X is given by

$$f(x) = P(X = x) \geq 0, \quad \text{for any } x \in \mathbb{R}.$$

2. Equivalent characterization via the cdf.

The pmf of a discrete r.v. X can also be written as

$$f(x) = P(X \leq x) - P(X < x) = F(x) - \lim_{y \rightarrow x^-} F(y).$$

Continuous Random Variables and Probability Density Functions

- If a cdf F is such that there exists a nonnegative function f satisfying

$$F(x) = \int_{-\infty}^x f(t) dt, \quad \text{for any } x \in \mathbb{R},$$

then the distribution described by F is said to be (absolutely) **continuous** with **probability density function (pdf)** f .

A random variable X with an (absolutely) continuous cdf F , or a pdf f , is said to be a **continuous random variable**.

- If F is (absolutely) continuous, then its derivative at $x \in \mathbb{R}$ is its pdf $f(x)$:

$$F'(x) = \frac{dF(x)}{dx} = f(x).$$

- If X is a continuous random variable, then

$$P(X = x) = 0 \quad \text{for any } x \in \mathbb{R}.$$

For $a < b$,

$$P(a < X < b) = P(a \leq X \leq b) = P(a \leq X < b) = P(a < X \leq b) = F(b) - F(a) = \int_a^b f(t) dt.$$

Properties of Probability Density or Mass Functions

A function $f(x)$ is a pdf (or pmf) for some random variable if and only if

1. $f(x) \geq 0$ for any $x \in \mathbb{R}$

2.

$$\int_{-\infty}^{\infty} f(x) dx = 1 \quad (\text{or } \sum_x f(x) = 1)$$

Any nonnegative function having a finite integral (or sum) can be turned into a pdf (or pmf) f by dividing by its integral (or sum).

We will write $X \sim f_X(x)$ (or $X \sim F_X(x)$) to denote that X has a distribution given by f (or F).

Computing Probabilities Using a pmf or pdf

To find general probabilities using a pmf or pdf, note that for $A \subset \mathbb{R}$,

Discrete case (using pmf):

$$P(X \in A) = \sum_{x \in A} f_X(x) = \sum_{x \in A, f_X(x) > 0} f_X(x)$$

Continuous case (using pdf):

$$P(X \in A) = \int_A f_X(x) dx$$

Relating the CDF to the PMF / PDF

Discrete random variable case

$$P(a < X \leq b) = F(b) - F(a)$$

$$P(a \leq X \leq b) = F(b) - F(a) + f(a)$$

$$P(a \leq X < b) = F(b) - F(a) + f(a) - f(b)$$

$$P(a < X < b) = F(b) - F(a) - f(b)$$

Continuous random variable case

$$P(a < X \leq b) = F(b) - F(a)$$

$$P(a \leq X \leq b) = F(b) - F(a)$$

$$P(a \leq X < b) = F(b) - F(a)$$

$$P(a < X < b) = F(b) - F(a)$$

Equivalently,

$$P(a < X < b) = \int_a^b f(x) dx.$$

Functions of a Random Variable

Introduction

- Consider a random variable $X \sim F_X(\cdot)$ and a function

$$g : \mathbb{R} \rightarrow \mathbb{R}.$$

(Here, X is a random variable and g may be *any* function.)

- Then

$$Y = g(X)$$

is also a random variable, having its own cdf $F_Y(\cdot)$.

Since Y is a function of X , we can describe the probabilistic behavior of Y in terms of that of X .

- Formally, there is also an inverse mapping g^{-1} defined by

$$g^{-1}(A) = \{x \in \mathbb{R} : g(x) \in A\}, \quad \text{for any } A \subset \mathbb{R}.$$

Distribution of a Function of a Random Variable

- The distribution of $Y = g(X)$ is completely determined by the distribution of X and the function g .

For any set $A \subset \mathbb{R}$,

$$P_Y(Y \in A) = P_X(g(X) \in A) = P_X(X \in g^{-1}(A)).$$

That is, the distribution of Y depends on the cdf (or pdf/pmf) F_X of X together with the function g .

Support (Range) Under Transformations

- If X has pdf/pmf $f_X(x)$, then the **range (support)** of X is

$$\mathcal{X} = \{x \in \mathbb{R} : f_X(x) > 0\}.$$

- If $Y = g(X)$ has pdf/pmf $f_Y(y)$, then the **range (support)** of Y is

$$\mathcal{Y} = \{y \in \mathbb{R} : f_Y(y) > 0\} = \{g(x) : x \in \mathcal{X}\}.$$

Discrete Case

Result.

If X is a discrete random variable with pmf $f_X(x)$ (i.e., X has range

$$\mathcal{X} = \{x \in \mathbb{R} : f_X(x) > 0\},$$

which is either finite or countably infinite), then

$$Y = g(X)$$

is also a discrete random variable with pmf

$$f_Y(y) = P(Y = y) = \begin{cases} \sum_{x \in g^{-1}(\{y\})} f_X(x) & y \in \mathcal{Y}, \\ 0 & y \notin \mathcal{Y}, \end{cases}$$

where the range (support) of Y is

$$\mathcal{Y} = \{g(x) : x \in \mathcal{X}\} = \{y \in \mathbb{R} : f_Y(y) > 0\}.$$

Continuous Case

For a continuous random variable X , the random variable

$$Y = g(X)$$

will *typically* (but not always) be continuous.

To determine the distribution of Y , one can use either of the following two approaches.

CDF Method

Compute the cdf $F_Y(\cdot)$ of Y :

$$\begin{aligned} F_Y(y) &= P(Y \leq y) \\ &= P(g(X) \leq y) \\ &= P(\{x \in \mathbb{R} : g(x) \leq y\}) \\ &= \int_{\{x \in \mathbb{R} : g(x) \leq y\}} f_X(x) dx. \end{aligned}$$

This is a general approach, but its success depends on being able to evaluate the integral.

PDF (Transformation) Method

Alternatively, one may compute the pdf $f_Y(\cdot)$ directly using a transformation technique.

This method is **only valid** when the function g is **monotone** or **piecewise monotone**.

Key Result

Theorem 2.1.5 (Monotone Transformation)

If X has pdf $f_X(x)$ and

$$Y = g(X),$$

where the function $g(\cdot)$ has either a **strictly positive** or a **strictly negative** derivative on

$$\mathcal{X} = \{x \in \mathbb{R} : f_X(x) > 0\},$$

then the pdf of Y has support

$$\mathcal{Y} = \{g(x) : x \in \mathcal{X}\},$$

and is given by

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{dg^{-1}(y)}{dy} \right| > 0, \quad \text{for } y \in \mathcal{Y},$$

with

$$f_Y(y) = 0, \quad \text{for } y \notin \mathcal{Y}.$$

Note that unless g is **strictly monotone** (or at least there is a way to break up

$$\mathcal{X} = \{x \in \mathbb{R} : f_X(x) > 0\}$$

into several intervals on each of which g is strictly increasing or strictly decreasing), X being a continuous random variable does **not** necessarily imply that

$$Y = g(X)$$

will be a continuous random variable.

Probability Integral Transform (PIT)

This is a famous (and for some purposes very useful) transformation connected with continuous cdfs.

If F is a continuous cdf, then

$$F(x) = \int_{-\infty}^x f(t) dt, \quad t \in \mathbb{R}.$$

If X has a continuous cdf $F(\cdot)$, then the random variable

$$Y = F(X)$$

is uniformly distributed on $(0, 1)$.

That is, Y has pdf

$$f_Y(y) = \begin{cases} 1, & 0 < y < 1, \\ 0, & \text{otherwise,} \end{cases}$$

and cdf

$$F_Y(y) = \begin{cases} 0, & y \leq 0, \\ y, & 0 \leq y \leq 1, \\ 1, & y \geq 1. \end{cases}$$

Expected Value of a Function of a Random Variable

Definition.

The expected value (or mean) of a random variable $g(X)$, denoted by $Eg(X)$, $E[g(X)]$, or $E(g(X))$, is defined as follows.

Discrete case:

$$Eg(X) = \sum_x g(x) f_X(x).$$

Continuous case:

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

Existence of the Expectation

The expectation $Eg(X)$ is defined **provided that**

Discrete case:

$$\sum_x |g(x)| f_X(x) < \infty,$$

Continuous case:

$$\int_{-\infty}^{\infty} |g(x)| f_X(x) dx < \infty.$$

(That is, we require $E[g(X)]$ to be a real, finite number.)

Nonexistence of the Expectation

We say that the expected value (or mean) $Eg(X)$ **does not exist** if

Discrete case:

$$\sum_x |g(x)| f_X(x) = \infty,$$

Continuous case:

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

Theorem 2.2.5 (Properties of Expectation)

Theorem.

Suppose X is a random variable such that

$$E|g_1(X)| < \infty \quad \text{and} \quad E|g_2(X)| < \infty,$$

and let $a, b, c \in \mathbb{R}$ be fixed constants. Then:

1.

$$E[ag_1(X) + b] = a Eg_1(X) + b.$$

2.

$$E[ag_1(X) + bg_2(X) + c] = a Eg_1(X) + b Eg_2(X) + c.$$

3. If $g_1(x) \geq a$ for all x , then

$$Eg_1(X) \geq a.$$

4. If $g_1(x) \leq b$ for all x , then

$$\mathbb{E}g_1(X) \leq b.$$

5. If $g_1(x) \geq g_2(x)$ for all x , then

$$\mathbb{E}g_1(X) \geq \mathbb{E}g_2(X).$$

Invariance of Expectation Under Transformation

Expectations are invariant under transformation.

If

$$Y = g(X),$$

then

$$\mathbb{E}Y = \sum_y y f_Y(y) = \sum_y y P(Y = y) = \sum_x g(x) f_X(x) = \mathbb{E}g(X)$$

in the discrete case.

(In the continuous case, replace sums with integrals.)

That is,

$$\mathbb{E}(Y) = \int_{-\infty}^{\infty} y f_Y(y) dy = \int_{-\infty}^{\infty} g(x) f_X(x) dx = \mathbb{E}g(X).$$

Variance

An important instance of the $\mathbb{E}g(X)$ notation arises when

$$g(X) = (X - \mathbb{E}X)^2.$$

Definition.

The **variance** of a random variable X , denoted $\text{Var}(X)$ or σ_X^2 , is

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

the expected squared distance between X and its mean $\mathbb{E}X$.

Two Important Variance Facts

1. For any real numbers a, b ,

$$\text{Var}(a + bX) = b^2 \text{Var}(X).$$

2.

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

Other Moments and Distributional Summaries

Moments are an important summary of a distribution.

1.

$$\mu = \mu_X = \mathbb{E}X$$

is often called the **mean**.

2.

$$\mu'_n = EX^n$$

is the n th (raw) moment, provided EX^n exists, i.e.,

Discrete case:

$$\sum_x |x^n| f_X(x) < \infty,$$

Continuous case:

$$\int_{-\infty}^{\infty} |x^n| f_X(x) dx < \infty.$$

3.

$$\mu_n = E[(X - \mu)^n]$$

is the n th **central moment**, provided EX^n exists.

(a)

$$\text{Var}(X) = \sigma_X^2 = E[(X - \mu)^2] = \mu_2$$

is the **variance**.

(b)

$$\sigma_X = \sqrt{\text{Var}(X)}$$

is the **standard deviation**.

(c)

$$\mu_3$$

is **skewness** (i.e., measures distributional balance around μ).

(d)

$$\mu_4$$

is **kurtosis** (i.e., a measure of how long the distributional tails are).

Regarding Moments

1. If EX^r exists for some $r > 0$, then EX^s exists for all

$$0 \leq s \leq r.$$

2. If EX^r does not exist for some $r > 0$, then EX^s will not exist for any

$$s > r.$$

3.

EX^2 exists if and only if $\text{Var}(X)$ exists.

4. For $r > 0$, the existence of EX^r is a matter of the distribution of X not having **heavy tails**, i.e., X does not assume large values with large probability.

Convergence (more on this later)

- Suppose a random variable X has mgf $M_X(t)$ and suppose X_1, X_2, \dots are a sequence of random variables, where X_n has mgf $M_{X_n}(t)$ for each $n \geq 1$.

If

$$\lim_{n \rightarrow \infty} M_{X_n}(t) = M_X(t)$$

holds for all $t \in (-h, h)$ for some $h > 0$, then the sequence X_1, X_2, \dots converges (in distribution) to X .

Moment Generating Functions

Other Generating Functions

There is nothing particularly illuminating about mgfs: these are just a technical device that are sometimes useful for proving theorems. In fact, they are not the only transforms for technical purposes or even the most useful (e.g., Fourier transforms / characteristic functions are often more useful in STAT 642).

Recall that the moment generating function (mgf) is defined as

$$M_X(t) = \mathbb{E}[e^{tX}].$$

Cumulant Generating Function

The function $\log M_X(t)$ is called the **cumulant generating function**.

The n th cumulant is given by

$$\left. \frac{d^n}{dt^n} \log M_X(t) \right|_{t=0}.$$

- The **first cumulant** is $\mathbb{E}[X]$
- The **second cumulant** is $\text{Var}(X)$

Factorial Moment Generating Function

The function

$$\mathbb{E}[t^X]$$

is called the **factorial moment generating function (fmgf)**.

The n th factorial moment is given by

$$\left. \frac{d^n}{dt^n} \mathbb{E}[t^X] \right|_{t=1} = \mathbb{E}[X(X-1)\cdots(X-n+1)].$$

For a discrete random variable, the fmfg is also called the **probability generating function**.

Interchanging Orders of Limits

We have several times in lecture up to this point interchanged “orders of limits” (i.e., switching the order of derivatives and expectations, or derivatives and summations, or derivatives and integrals).

This interchange is also implicit in using mgfs to compute moments:

$$\begin{aligned}
\frac{d^n}{dt^n} M_X(t) \Big|_{t=0} &= \frac{d^n}{dt^n} \mathbb{E}[e^{tX}] \Big|_{t=0} \\
&= \mathbb{E}\left(\frac{d^n}{dt^n} e^{tX} \Big|_{t=0}\right) \\
&= \mathbb{E}\left(X^n e^{tX} \Big|_{t=0}\right) \\
&= \mathbb{E}[X^n].
\end{aligned}$$

We don't need to worry about the validity of these interchanges here (it's covered in STAT 642, where it makes more sense with the right technical material in hand).

Inequalities

Markov inequality: Suppose X is a random variable and $g(x) \geq 0$. Then, for any $r > 0$,

$$\mathbb{P}(g(X) \geq r) \leq \frac{\mathbb{E}[g(X)]}{r}.$$

Chebyshev inequality: Suppose X is a random variable with mean $\mathbb{E}[X] = \mu$ and variance σ^2 . Then, for any $k > 0$,

$$\mathbb{P}(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2},$$

and equivalently,

$$\mathbb{P}(|X - \mu| < k\sigma) \geq 1 - \frac{1}{k^2}.$$

Convex Functions

Definition: A function $g(x)$ is **convex** if

$$g(\lambda x + (1 - \lambda)y) \leq \lambda g(x) + (1 - \lambda)g(y)$$

for all x, y and all $0 < \lambda < 1$.

Second Derivative Characterization

If g is twice differentiable, then $g(x)$ is convex if

$$g''(x) \geq 0 \quad \text{for all } x.$$

Example:

If $g(x) = x^2$, then $g'(x) = 2x$ and $g''(x) = 2 > 0$, so $g(x)$ is convex.

A function $g(x)$ is **concave** if $-g(x)$ is convex.

Jensen's inequality: Suppose X is a random variable and $g(x)$ is a convex function. Then,

$$\mathbb{E}[g(X)] \geq g(\mathbb{E}[X]).$$

For a **concave** function, the reverse inequality holds:

$$g(\mathbb{E}[X]) \geq \mathbb{E}[g(X)].$$

Common Distributions

Introduction

- Often it is useful to consider structural forms for a pdf f or cdf F , especially for modeling a population.
- In particular, it is possible to specify a family of distributions using a single functional form with one or more free **parameters**.

Discrete Distributions

Bernoulli Distribution

A **Bernoulli trial** is a random variable

$$X \sim \text{Bern}(p),$$

where p is a parameter.

Definition

- There are **two outcomes**, usually labeled 0 and 1 (failure / success).

$$X = \begin{cases} 0, & \text{with probability } 1 - p, \\ 1, & \text{with probability } p. \end{cases}$$

- The parameter satisfies $0 < p < 1$ (i.e., $X = 1$ with probability p).

Probability Mass Function

If $X \sim \text{Bern}(p)$, then

$$\mathbb{P}(X = 0) = 1 - p, \quad \mathbb{P}(X = 1) = p.$$

Equivalently, the pmf can be written as

$$f_X(x) = f_X(x | p) = \mathbb{P}(X = x | p) = p^x(1 - p)^{1-x}, \quad x = 0, 1.$$

Expectation

The expectation of X is

$$\mathbb{E}[X] = \sum_{x=0}^1 x \mathbb{P}(X = x) = (0)(1-p) + (1)(p) = p.$$

Second Moment

The second moment is

$$\mathbb{E}[X^2] = \sum_{x=0}^1 x^2 \mathbb{P}(X = x) = (0)^2(1-p) + (1)^2(p) = p.$$

Variance

The variance of X is

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = p - p^2 = p(1-p).$$

Binomial Distribution

The **Binomial distribution** is defined as

$$X \sim \text{Binom}(n, p), \quad 0 < p < 1,$$

where n is the number of trials and p is the success probability.

Probability Mass Function

The pmf is given by

$$f_X(x) = f_X(x \mid n, p) = \mathbb{P}(X = x \mid n, p) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n.$$

Motivation

The Binomial distribution arises as the distribution of the **number of successes** in n independent Bernoulli(p) trials.

Let

$$Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Bern}(p),$$

where $Y_i = 1$ if the i th trial is a “success” and 0 if it is a “failure.” Define

$$X = \sum_{i=1}^n Y_i.$$

Then $X \sim \text{Binom}(n, p)$.

Derivation of the pmf

For a given $x = 0, 1, \dots, n$,

$$\mathbb{P}(X = x) = \mathbb{P}\left(\sum_{i=1}^n Y_i = x\right).$$

Since Y_1, \dots, Y_n are independent,

$$\mathbb{P}(X = x) = \sum_{\substack{(y_1, \dots, y_n) \in \{0,1\}^n \\ \sum_{i=1}^n y_i = x}} \mathbb{P}(Y_1 = y_1, \dots, Y_n = y_n) = \sum_{\substack{(y_1, \dots, y_n) \in \{0,1\}^n \\ \sum_{i=1}^n y_i = x}} \prod_{i=1}^n \mathbb{P}(Y_i = y_i).$$

Because each $Y_i \sim \text{Bern}(p)$,

$$\prod_{i=1}^n \mathbb{P}(Y_i = y_i) = p^{\sum_i y_i} (1-p)^{n - \sum_i y_i} = p^x (1-p)^{n-x}.$$

Thus,

$$\mathbb{P}(X = x) = p^x (1-p)^{n-x} \times (\text{number of ways to choose exactly } x \text{ components of } (y_1, \dots, y_n) \text{ to be 1}).$$

The number of such configurations is $\binom{n}{x}$, giving

$$\mathbb{P}(X = x) = \binom{n}{x} p^x (1-p)^{n-x}.$$

Mean and Variance

- Mean:

$$\mathbb{E}[X] = np \quad (\text{proof on next slide})$$

- Variance:

$$\text{Var}(X) = np(1-p).$$

Moment Generating Function

The moment generating function of X is

$$M_X(t) = \mathbb{E}[e^{tX}] = (pe^t + (1-p))^n, \quad t \in \mathbb{R}.$$

This follows from

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

and the binomial expansion of $(a+b)^n$.

Negative Binomial

Let

$$X \sim \text{Neg-Binom}(r, p), \quad r \in \mathbb{Z}, \quad r \geq 1, \quad 0 < p < 1.$$

Probability Mass Function

The pmf is given by

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

Motivation

The Negative Binomial distribution describes the distribution of the **number of independent Bernoulli(p) trials needed to obtain r successes**.

Equivalently, consider a sequence of successes (S) and failures (F) of total length x such that the r th success occurs on the x th trial.

Alternative Parameterization

Define

$$Y = X - r,$$

the **number of failures prior to the r th success**. Then Y is also commonly referred to as Negative Binomial.

The pmf of Y is

$$\mathbb{P}(Y = y) = f_Y(y \mid r, p) = \binom{y+r-1}{y} p^r (1-p)^y, \quad y = 0, 1, 2, \dots$$

(Note that $\binom{y+r-1}{y} = \binom{y+r-1}{r-1}$.)

Caution

Be careful: **both** random variables X and Y (which are different) are often called “negative binomial” in the literature.

Mean and Variance

For Y (failures before the r th success),

$$\mathbb{E}[Y] = \frac{r(1-p)}{p}.$$

Since $X = Y + r$,

$$\mathbb{E}[X] = \mathbb{E}[Y] + r = \frac{r}{p}.$$

The variance is

$$\text{Var}(Y) = \frac{r(1-p)}{p^2} \quad \text{and} \quad \text{Var}(X) = \text{Var}(Y).$$

Moment Generating Function

The mgf of Y is

$$M_Y(t) = \mathbb{E}[e^{tY}] = \left[\frac{p}{1 - (1-p)e^t} \right]^r, \quad t < -\log(1-p).$$

Since $X = Y + r$,

$$M_X(t) = \mathbb{E}[e^{t(X)}] = \mathbb{E}[e^{t(Y+r)}] = e^{rt} M_Y(t).$$

Geometric

Let

$$X \sim \text{Geom}(p), \quad 0 < p < 1.$$

Relationship to the Negative Binomial

The Geometric distribution is a **special case** of the Negative Binomial distribution with $r = 1$.

Motivation

The Geometric distribution describes the distribution of the **number of independent Bernoulli(p) trials needed to obtain the first success**.

Probability Mass Function

The pmf is

$$f_X(x) = f_X(x \mid p) = p(1-p)^{x-1}, \quad x = 1, 2, 3, \dots$$

Mean and Variance

The mean is

$$\mathbb{E}[X] = \frac{1}{p}.$$

The variance is

$$\text{Var}(X) = \frac{1-p}{p^2}.$$

Moment Generating Function

The mgf of X is

$$M_X(t) = \mathbb{E}[e^{tX}] = \frac{pe^t}{1 - (1-p)e^t}, \quad t < -\log(1-p).$$

Memoryless Property of the Geometric Distribution The Geometric distribution has the famous **memoryless** property: for any integer $x_0 \geq 0$,

$$\mathbb{P}(X = x_0 + x \mid X > x_0) = \mathbb{P}(X = x).$$

Interpretation

The conditional distribution of the remaining waiting number of trials until the first success, given that we have already waited x_0 trials, is the same as the original distribution of the number of trials until the first success.

Given that each trial is an independent Bernoulli trial, this makes sense: whether we start counting trials at the beginning, or we start counting trials after x_0 trials without success, the distribution of the remaining number of trials needed until the first success should be the same.

Hypergeometric

Let

$$X \sim \text{Hypergeometric}(N, M, K),$$

where N, M, K are integers.

Probability Mass Function

The pmf is given by

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

Motivation

Choose K objects **without replacement** from a total population of size N that contains M “special” objects.

- N = total population size
- M = number of special objects
- $N - M$ = number of non-special objects
- X = number of special objects among the K chosen

Support Conditions

The pmf is nonzero only when

$$0 \leq x \leq K, \quad x \leq M, \quad K - x \leq N - M.$$

Typically, when $N > 2M$ and $M > K$, only the condition

$$0 \leq x \leq K$$

matters.

Mean

The mean of a Hypergeometric random variable is

$$\mathbb{E}[X] = K \frac{M}{N}.$$

(Interpretation: sample size \times proportion of special objects.)

Variance

The variance is

$$\text{Var}(X) = \frac{KM(N-M)(N-K)}{N^2(N-1)}.$$

Connection between Hypergeometric and Binomial

Hypergeometric (sampling without replacement)

Choose x special objects in a sample of size K drawn **without replacement** from a population where M objects are “special” and $N - M$ are not.

Binomial (sampling with replacement)

Choose x special objects in a sample of size n , where each selected item has probability p of being a special object.

In particular, if sampling *with replacement* from a population of size N containing M special objects, then

$$X \sim \text{Binomial}\left(n = K, p = \frac{M}{N}\right).$$

Result

As the population size N becomes large, the **Hypergeometric distribution tends to the Binomial distribution**.

Poisson

Let

$$X \sim \text{Poisson}(\lambda), \quad \lambda > 0.$$

Probability Mass Function

The pmf is given by

$$\mathbb{P}(X = x) = f_X(x) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

Motivation

The Poisson distribution is widely used to model **rare-event count data**, for example: - number of car accidents in a county, - number of system failures in a fixed time period.

Useful Identity

Recall that for any real a ,

$$e^a = \sum_{k=0}^{\infty} \frac{a^k}{k!}.$$

This identity ensures that the Poisson pmf satisfies: - $f_X(x) \geq 0$ for all x , - $\sum_x f_X(x) = 1$.

Mean

The mean of a Poisson random variable is

$$\mathbb{E}[X] = \lambda.$$

This follows from

$$\mathbb{E}[X] = \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!}.$$

Variance

The variance of a Poisson random variable is

$$\text{Var}(X) = \lambda.$$

One way to derive this is by noting that

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \mathbb{E}[X(X-1)] + \mathbb{E}[X] - (\mathbb{E}[X])^2,$$

and showing that

$$\mathbb{E}[X(X-1)] = \lambda^2.$$

Moment Generating Function

The moment generating function of X is

$$M_X(t) = \mathbb{E}[e^{tX}] = \exp(\lambda(e^t - 1)), \quad t \in \mathbb{R}.$$

Continuous Distributions

Uniform

Let

$$X \sim \text{Uniform}(a, b), \quad -\infty < a < b < \infty.$$

Probability Density Function

The pdf is given by

$$f_X(x) = \frac{1}{b-a}, \quad a < x < b.$$

Motivation

The Uniform distribution models **equally likely outcomes** over a finite range (a, b) .

- a is the lower endpoint of the range
- b is the upper endpoint

Moments

For $r > 0$,

$$\mathbb{E}[X^r] = \frac{1}{b-a} \int_a^b x^r dx = \frac{b^{r+1} - a^{r+1}}{(r+1)(b-a)}.$$

Mean

$$\mathbb{E}[X] = \frac{a+b}{2}.$$

Variance

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \frac{(b-a)^2}{12}.$$

Important Case: Uniform(0, 1)

Let

$$U \sim \text{Uniform}(0, 1).$$

Then:

1. $f_U(u) = 1$, for $0 < u < 1$
2. $\mathbb{E}[U] = \frac{1}{2}$
3. $\text{Var}(U) = \frac{1}{12}$

Probability Integral Transform (PIT)

1. If Y has a **continuous cdf** $F_Y(y)$, then

$$U = F_Y(Y) \sim \text{Uniform}(0, 1).$$

2. Conversely, if $U \sim \text{Uniform}(0, 1)$ and $F_Y(y)$ is a continuous cdf, then

$$Y = F_Y^{-1}(U)$$

has distribution F_Y .

(This is useful for simulating random variables.)

Gamma

Let

$$X \sim \text{Gamma}(\alpha, \beta), \quad \alpha > 0, \beta > 0.$$

Probability Density Function

The pdf is given by

$$f_X(x) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-x/\beta}, \quad 0 < x < \infty.$$

Motivation

The Gamma distribution is a **flexible family for positive-valued random variables**.

Parameters

- $\alpha > 0$ is the **shape parameter**
 - If $\alpha < 1$, the density is unbounded near $x = 0$
 - If $\alpha > 1$, the density is zero at $x = 0$
- $\beta > 0$ is the **scale parameter**

If $X \sim \text{Gamma}(\alpha, \beta)$, then

$$Z = \frac{X}{\beta} \sim \text{Gamma}(\alpha, 1).$$

The Gamma Function

The Gamma function is defined as

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx, \quad \alpha > 0,$$

which ensures that $f_X(x)$ integrates to 1.

Properties of the Gamma Function

1. $\Gamma(1 + \alpha) = \alpha\Gamma(\alpha)$, for $\alpha > 0$
2. $\Gamma(\alpha) = (\alpha - 1)!$, for integers $\alpha \geq 1$
3. $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$

Moments

For $r > 0$,

$$\mathbb{E}[X^r] = \beta^r \frac{\Gamma(\alpha + r)}{\Gamma(\alpha)}.$$

Mean

$$\mathbb{E}[X] = \alpha\beta.$$

Variance

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

Moment Generating Function

The mgf of X is

$$M_X(t) = \mathbb{E}[e^{tX}] = (1 - \beta t)^{-\alpha}, \quad t < \frac{1}{\beta}.$$

Relationship Between Gamma and Poisson (Integer Shape) For integer α ,

$$F_X(x | \alpha, \beta) = \mathbb{P}(X \leq x) = \mathbb{P}(Y \geq \alpha), \quad \text{where } Y \sim \text{Poisson}(x/\beta).$$

Continuous Distributions: Gamma — Special Cases

Chi-squared Distribution If $p > 0$ is an integer, then

$$\chi_p^2 \sim \text{Gamma}\left(\alpha = \frac{p}{2}, \beta = 2\right),$$

where p is called the **degrees of freedom** parameter.

Exponential Distribution The Exponential distribution is a special case of the Gamma distribution:

$$\text{Exp}(\beta) = \text{Gamma}(\alpha = 1, \beta).$$

The pdf is

$$f_X(x) = \frac{1}{\beta} e^{-x/\beta}, \quad 0 < x < \infty.$$

The Exponential distribution is commonly used to model **failure times** and has the **memoryless property**:

$$\mathbb{P}(X > s + t | X > t) = \mathbb{P}(X > s).$$

Weibull Distribution If

$$X \sim \text{Exp}(\beta)$$

and $\gamma > 0$, then

$$W = X^{1/\gamma} \sim \text{Weibull}(\gamma, \beta).$$

The pdf of W is

$$f_W(w) = \frac{\gamma}{\beta} w^{\gamma-1} e^{-w^\gamma/\beta}, \quad 0 < w < \infty.$$

The Weibull distribution is a **general failure time distribution**.

Inverse-Gamma Distribution If

$$X \sim \text{Gamma}(\alpha, \beta),$$

then

$$Y = \frac{1}{X}$$

has the **inverse-gamma distribution**, with pdf

$$f_Y(y) = \frac{1}{\Gamma(\alpha)\beta^\alpha} \left(\frac{1}{y}\right)^{1+\alpha} e^{-1/(\beta y)}, \quad 0 < y < \infty.$$

Beta

Let

$$X \sim \text{Beta}(\alpha, \beta), \quad \alpha > 0, \beta > 0.$$

Probability Density Function

The pdf is

$$f_X(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad 0 < x < 1.$$

Motivation

The Beta distribution is a flexible family, often used for modeling **proportions**.

Shape Parameters

Both α and β are **shape parameters**:

1. α determines behavior near $x = 0$
 - If $\alpha < 1$, density is unbounded near 0
 - If $\alpha > 1$, density is zero at 0
2. β determines behavior near $x = 1$
 - If $\beta < 1$, density is unbounded near 1
 - If $\beta > 1$, density is zero at 1
3. If $\alpha = \beta$, the distribution is symmetric.

Beta Function

The Beta function is defined as

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}.$$

Moments

For $r > 0$,

$$\mathbb{E}[X^r] = \frac{B(\alpha + r, \beta)}{B(\alpha, \beta)} = \frac{\Gamma(\alpha + \beta)\Gamma(\alpha + r)}{\Gamma(\alpha + \beta + r)\Gamma(\alpha)}.$$

Mean

$$\mathbb{E}[X] = \frac{\alpha}{\alpha + \beta}.$$

Variance

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

Related Distribution

If $\alpha = \beta = 1$, then

$$X \sim \text{Uniform}(0, 1).$$

Normal (Gaussian)

Let

$$X \sim N(\mu, \sigma^2), \quad -\infty < \mu < \infty, \sigma > 0.$$

Probability Density Function

The pdf is

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \left(\frac{x - \mu}{\sigma}\right)^2\right), \quad -\infty < x < \infty.$$

Motivation

The Normal distribution is the **single most important distribution** in statistics:

- widely used and analytically tractable
- bell-shaped density arises naturally
- Central Limit Theorem (normal distribution is extremely relevant in large samples)

Parameters

- $\mu \in \mathbb{R}$ is the mean: $\mathbb{E}[X] = \mu$
- $\sigma^2 = \text{Var}(X)$ is the variance; σ is the standard deviation

Standard Normal Distribution

Many properties of the Normal distribution are most easily derived using the standard normal distribution $N(0, 1)$.

1. If $X \sim N(\mu, \sigma^2)$, then

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1).$$

2. If $Z \sim N(0, 1)$, then for constants $a, b \in \mathbb{R}$,

$$X = a + bZ \sim N(\mu = a, \sigma^2 = b^2).$$

From Standard Normal to Normal

If, for $\mu \in \mathbb{R}$ and $\sigma > 0$, a random variable X has the same distribution as

$$X = \mu + \sigma Z,$$

where $Z \sim N(0, 1)$ (standard normal), then we say that X has a normal distribution with mean μ and variance σ^2 , denoted

$$X \sim N(\mu, \sigma^2).$$

Some facts / properties of $X \sim N(\mu, \sigma^2)$, $\mu \in \mathbb{R}$, $\sigma > 0$: **CDF**

For $-\infty < x < \infty$,

$$\begin{aligned} F_X(x) &= P(X \leq x) \\ &= P(\mu + \sigma Z \leq x) \\ &= P\left(Z \leq \frac{x - \mu}{\sigma}\right) \\ &= \Phi\left(\frac{x - \mu}{\sigma}\right), \end{aligned}$$

where $\Phi(\cdot)$ denotes the CDF of the standard normal distribution.

PDF

For $-\infty < x < \infty$,

$$\begin{aligned} f_X(x) &= \frac{d}{dx} F_X(x) \\ &= \frac{d}{dx} \Phi\left(\frac{x - \mu}{\sigma}\right) \\ &= \frac{1}{\sigma} \phi\left(\frac{x - \mu}{\sigma}\right) \\ &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \left(\frac{x - \mu}{\sigma}\right)^2\right\}, \end{aligned}$$

where $\phi(\cdot)$ denotes the standard normal density.

Mean

$$\begin{aligned}\mathbb{E}[X] &= \mathbb{E}(\mu + \sigma Z) \\ &= \mu + \sigma \mathbb{E}[Z] \\ &= \mu.\end{aligned}$$

Variance

$$\begin{aligned}\text{Var}(X) &= \text{Var}(\mu + \sigma Z) \\ &= \sigma^2 \text{Var}(Z) \\ &= \sigma^2.\end{aligned}$$

Also

Given μ , σ^2 , and $Z \sim N(0, 1)$,

$$X = \mu + \sigma Z \sim N(\mu, \sigma^2).$$

Moment Generating Function (MGF)

For $-\infty < t < \infty$,

$$\begin{aligned}M_X(t) &= \mathbb{E}[e^{tX}] \\ &= \mathbb{E}\left[e^{t(\mu+\sigma Z)}\right] \\ &= \mathbb{E}\left[e^{t\mu} e^{t\sigma Z}\right] \\ &= e^{t\mu} \mathbb{E}\left[e^{t\sigma Z}\right] \\ &= e^{\mu t} M_Z(\sigma t).\end{aligned}$$

Since the MGF of $Z \sim N(0, 1)$ is

$$M_Z(t) = \exp\left(\frac{t^2}{2}\right),$$

we obtain

$$M_X(t) = \exp\left(\mu t + \frac{\sigma^2 t^2}{2}\right).$$

Distributions Related to the Normal Log-normal Distribution

If

$$X \sim N(\mu, \sigma^2),$$

then

$$Y = e^X \sim \text{LogNormal}(\mu, \sigma^2).$$

Probability Density Function

The pdf of Y is given by

$$\begin{aligned} f_Y(y) &= \frac{1}{y} f_X(\log y) \\ &= \frac{1}{y\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2} \left(\frac{\log y - \mu}{\sigma} \right)^2 \right\}, \end{aligned}$$

for

$$0 < y < \infty.$$

Mean

$$\begin{aligned} \mathbb{E}[Y] &= \mathbb{E}[e^X] \\ &= M_X(t=1) \\ &= \exp \left(\mu + \frac{1}{2}\sigma^2 \right). \end{aligned}$$

Note:

By Jensen's inequality,

$$\mathbb{E}[Y] = \mathbb{E}[e^X] \geq e^{\mathbb{E}[X]} = e^\mu.$$

Shape and Applications

- The shape of the log-normal distribution is similar to that of a Gamma distribution with shape parameter $\alpha > 1$.
- Common in economics (e.g., assuming $\log(\text{income})$ is normally distributed).
- Also used as a failure-time distribution.

Chi-Square Connection

If

$$Z \sim N(0, 1),$$

then

$$Z^2 \sim \chi_1^2,$$

a chi-square distribution with 1 degree of freedom.

Other Continuous Distributions

Each of the following distributions has a **location parameter** $-\infty < \mu < \infty$ and a **scale parameter** $\sigma > 0$.

Cauchy Distribution

- μ is the **location parameter**
- σ is the **scale parameter**

Probability Density Function

$$f_X(x) = \frac{1}{\pi\sigma \left(1 + \left(\frac{x-\mu}{\sigma}\right)^2\right)}, \quad -\infty < x < \infty.$$

- A favorite **counter-example distribution**
- No finite moments
- No moment generating function (mgf)

In particular,

$$\mathbb{E}[X] = \infty, \quad \text{Var}(X) = \infty.$$

Logistic Distribution

Probability Density Function

$$f_X(x) = \frac{e^{(x-\mu)/\sigma}}{\sigma \left(1 + e^{(x-\mu)/\sigma}\right)^2}, \quad -\infty < x < \infty.$$

- Underlies **logistic regression**

Double Exponential (Laplace) Distribution

Probability Density Function

$$f_X(x) = \frac{1}{2\sigma} \exp\left(-\frac{|x-\mu|}{\sigma}\right), \quad -\infty < x < \infty.$$

- Similar to the normal distribution but with **heavier tails**

Extreme Value Distribution

Probability Density Function

$$f_X(x) = \frac{1}{\sigma} \exp\left(\frac{x-\mu}{\sigma}\right) \exp\left\{-\exp\left(\frac{x-\mu}{\sigma}\right)\right\}, \quad -\infty < x < \infty.$$

- Limiting distribution of **record highs / record lows**

Pareto Distribution

Probability Density Function

$$f_X(x) = \frac{\alpha}{\beta} \left(1 + \frac{x}{\beta}\right)^{-\alpha-1}, \quad 0 < x < \infty.$$

- $\alpha > 0$ is the **shape parameter**
- $\beta > 0$ is the **scale parameter**

Location-Scale Families of Continuous Random Variables

- We have already seen families with **location** and **scale** parameters (e.g., Normal, Cauchy).

Definition (Location-Scale Family)

Suppose $f_Z(z)$ is a given (fixed) pdf.
Then the collection of pdfs of the form

$$f_X(x | \mu, \sigma) = \frac{1}{\sigma} f_Z\left(\frac{x - \mu}{\sigma}\right),$$

for

$$-\infty < \mu < \infty, \quad \sigma > 0,$$

is called the **location-scale family** with standard pdf f_Z .

- μ is the **location parameter**
- $\sigma > 0$ is the **scale parameter**

Relationship via Transformation

If

$$Z \sim f_Z(z)$$

and

$$X = \mu + \sigma Z,$$

then

$$X \sim f_X(x | \mu, \sigma)$$

as defined above.

Distribution Function

$$\begin{aligned}
F_X(x) &= P(X \leq x) \\
&= P(\mu + \sigma Z \leq x) \\
&= P\left(Z \leq \frac{x - \mu}{\sigma}\right) \\
&= F_Z\left(\frac{x - \mu}{\sigma}\right).
\end{aligned}$$

Differentiating,

$$f_X(x) = \frac{d}{dx} F_X(x) = \frac{1}{\sigma} f_Z\left(\frac{x - \mu}{\sigma}\right).$$

Standard Distribution

The **standard pdf** f_Z belongs to the family as the special case

$$\mu = 0, \quad \sigma = 1.$$

Properties

Properties of the location-scale family can be obtained by studying the standard pdf f_Z .

For example, if moments exist,

$$\mathbb{E}[X] = \mu + \sigma \mathbb{E}[Z], \quad \text{Var}(X) = \sigma^2 \text{Var}(Z).$$

Location Family

The collection of pdfs of the form

$$f_X(x | \mu) = f_Z(x - \mu),$$

for $-\infty < \mu < \infty$, defines the **location family** with standard pdf f_Z .

Scale Family

The collection of pdfs of the form

$$f_X(x | \sigma) = \frac{1}{\sigma} f_Z\left(\frac{x}{\sigma}\right),$$

for $\sigma > 0$, defines the **scale family** with standard pdf f_Z .

Multivariate Distributions

Introduction

- We are generally interested in **more than one random variable at a time**.

1. **n observations of a single characteristic from some population**

$$(X_1, \dots, X_n)$$

2. k different characteristics from a single individual

For example, for one person:

- Y_1 : age
- Y_2 : height
- Y_3 : gender

Notation

Let

$$\mathbf{X} = (X_1, \dots, X_n)$$

denote an n -dimensional random vector (r.v.).

In other words, \mathbf{X} is a function from the sample space S to \mathbb{R}^n .

Cumulative Distribution Functions

The joint cumulative distribution function (cdf) of X and Y is

$$F_{X,Y}(x, y) = F(x, y) = P(X \leq x, Y \leq y), \quad x, y \in \mathbb{R}.$$

Discrete Case

- Rarely used in practice, but computed as

$$F(x, y) = P(X \leq x, Y \leq y) = \sum_{x_1 \leq x, y_1 \leq y} f(x_1, y_1).$$

Continuous Case

- Often useful, and computed as

$$F(x, y) = P(X \leq x, Y \leq y) = \int_{-\infty}^y \int_{-\infty}^x f(s, t) ds dt.$$

In the continuous case,

$$\frac{\partial^2 F(x, y)}{\partial x \partial y} = \frac{d}{dx} \frac{dF(x, y)}{dy} = f(x, y).$$

Properties of a Joint CDF

A function $F(x, y)$ is a cdf for some random vector (X, Y) if and only if:

1. **Limits at $-\infty$:**

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

2. **Limit at $+\infty$:**

$$\lim_{x, y \rightarrow \infty} F(x, y) = 1.$$

3. **Right continuity in each argument:**

$$\lim_{h \downarrow 0} F(x + h, y) = \lim_{h \downarrow 0} F(x, y + h) = F(x, y).$$

4. **Nondecreasing (nonnegative probability on rectangles):**

For all $x, y \in \mathbb{R}$ and any $\Delta_1, \Delta_2 > 0$,

$$\begin{aligned} P(x < X \leq x + \Delta_1, y < Y \leq y + \Delta_2) &= F(x + \Delta_1, y + \Delta_2) - F(x + \Delta_1, y) \\ &\quad - F(x, y + \Delta_2) + F(x, y) \\ &\geq 0. \end{aligned}$$

Continuous Random Vectors and Probability Density Functions

If there exists a function $f(x, y) \geq 0$ for $(x, y) \in \mathbb{R}^2$ such that

$$\iint_{\mathbb{R}^2} f(x, y) dx dy = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1,$$

and

$$P((X, Y) \in A) = \iint_A f(x, y) dx dy, \quad A \subset \mathbb{R}^2,$$

then we say that the distribution of X and Y is **jointly (absolutely) continuous**.

The function $f(x, y)$ is called the **joint probability density function (joint pdf)** of (X, Y) .

As in the univariate case, any function $f(x, y)$ satisfying:

1. **Nonnegativity**

$$f(x, y) \geq 0 \quad \text{for all } (x, y) \in \mathbb{R}^2,$$

2. Unit total probability

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1 \quad (\text{continuous case}),$$

or

$$\sum_{(x,y) \in \mathbb{R}^2} f(x, y) = 1 \quad (\text{discrete case}),$$

specifies the joint pdf (or pmf) of some bivariate random vector (X, Y) .

Marginal Distributions

Marginal Distributions in the Discrete Case

For a jointly discrete random vector (X, Y) with joint pmf $f(x, y)$,

$$f_X(x) = P(X = x) = \sum_y P(X = x, Y = y) = \sum_y f(x, y), \quad x \in \mathbb{R}.$$

Similarly,

$$f_Y(y) = P(Y = y) = \sum_x P(X = x, Y = y) = \sum_x f(x, y), \quad y \in \mathbb{R}.$$

Marginal Distributions in the Continuous Case

In the continuous (X, Y) case, for $x \in \mathbb{R}$,

$$\begin{aligned} P(X \leq x) &= P((X, Y) \in (-\infty, x] \times (-\infty, \infty)) \\ &= \int_{-\infty}^x \int_{-\infty}^{\infty} f(s, t) ds dt. \end{aligned}$$

This motivates the definition of marginal densities.

Marginal PDFs

For a jointly continuous random vector (X, Y) with joint pdf $f(x, y)$,

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy, \quad x \in \mathbb{R},$$

and

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx, \quad y \in \mathbb{R}.$$

We often say that we **integrate out** x or y to obtain a marginal pdf.

From Bivariate to Multivariate

Now consider the n -dimensional case: an \mathbb{R}^n -valued random vector (X_1, \dots, X_n) .

CDF

$$F(x_1, \dots, x_n) = P(X_1 \leq x_1, \dots, X_n \leq x_n), \quad x_1, \dots, x_n \in \mathbb{R}.$$

Properties:

1. If all $x_i \rightarrow \infty$, then

$$F(x_1, \dots, x_n) \rightarrow 1.$$

2. If any $x_i \rightarrow -\infty$, then

$$F(x_1, \dots, x_n) \rightarrow 0.$$

3. F is right-continuous in each argument.

4. A generalized **monotonicity condition** holds for n -dimensional rectangles.

PDF or PMF

A joint pdf or pmf satisfies

$$f(x_1, \dots, x_n) \geq 0$$

with integral (or sum) equal to 1.

For the continuous case,

$$P((X_1, \dots, X_n) \in A) = \int \cdots \int_{\{(x_1, \dots, x_n) \in A\}} f(x_1, \dots, x_n) dx_1 \cdots dx_n, \quad A \subset \mathbb{R}^n.$$

For the discrete case,

$$P((X_1, \dots, X_n) \in A) = \sum_{(x_1, \dots, x_n) \in A} f(x_1, \dots, x_n).$$

Marginal Distributions in the Multivariate Case

A marginal pdf or pmf can be obtained for **any subset** of (X_1, \dots, X_n) .

For example, to obtain the marginal distribution of (X_1, \dots, X_k) from (X_1, \dots, X_n) , fix $x_1, \dots, x_k \in \mathbb{R}$ and integrate (or sum) out the remaining variables:

$$f(x_1, \dots, x_k) = \int \cdots \int f(x_1, \dots, x_n) dx_{k+1} \cdots dx_n \quad (\text{continuous case}),$$

or

$$f(x_1, \dots, x_k) = \sum_{x_{k+1}, \dots, x_n} f(x_1, \dots, x_n) \quad (\text{discrete case}).$$

Expectations of Several Random Variables

Extension of the Univariate Case

Let (X_1, \dots, X_n) be a random vector and let $g : \mathbb{R}^n \rightarrow \mathbb{R}$.

Discrete Case

If (X_1, \dots, X_n) are jointly discrete with pmf f , define

$$\mathbb{E}[g(X_1, \dots, X_n)] = \sum_{(x_1, \dots, x_n)} g(x_1, \dots, x_n) f(x_1, \dots, x_n),$$

provided that

$$\sum_{(x_1, \dots, x_n)} |g(x_1, \dots, x_n)| f(x_1, \dots, x_n) < \infty.$$

Continuous Case

If (X_1, \dots, X_n) are jointly continuous with pdf f , define

$$\mathbb{E}[g(X_1, \dots, X_n)] = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(x_1, \dots, x_n) f(x_1, \dots, x_n) dx_1 \cdots dx_n,$$

provided that

$$\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} |g(x_1, \dots, x_n)| f(x_1, \dots, x_n) dx_1 \cdots dx_n < \infty.$$

Bivariate Case ($n = 2$)

These reduce to

$$\mathbb{E}[g(X, Y)] = \begin{cases} \sum_{(x,y)} g(x, y) f(x, y), & \text{discrete case,} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f(x, y) dx dy, & \text{continuous case.} \end{cases}$$

Covariance

The **covariance** of X_i and X_j is defined as

$$\text{Cov}(X_i, X_j) = \sigma_{X_i, X_j} = \mathbb{E}[(X_i - \mathbb{E}X_i)(X_j - \mathbb{E}X_j)].$$

Equivalently,

$$\text{Cov}(X_i, X_j) = \mathbb{E}[X_i X_j] - (\mathbb{E}X_i)(\mathbb{E}X_j).$$

Interpretation

- If $\text{Cov}(X_i, X_j) > 0$, then larger (or smaller) than average values of X_i tend to occur with larger (or smaller) than average values of X_j .

- If $\text{Cov}(X_i, X_j) < 0$, then larger (or smaller) than average values of X_i tend to occur with smaller (or larger) than average values of X_j .

Relationships and Properties

1. Variance as a special case

$$\text{Cov}(X_i, X_i) = \text{Var}(X_i).$$

Indeed,

$$\text{Cov}(X_i, X_i) = \mathbb{E}[X_i^2] - (\mathbb{E}X_i)^2 = \text{Var}(X_i).$$

2. Symmetry

$$\text{Cov}(X_i, X_j) = \text{Cov}(X_j, X_i).$$

3. Linearity with constants

For constants a, b, c, d ,

$$\text{Cov}(aX_i + c, bX_j + d) = ab \text{Cov}(X_i, X_j).$$

4. Variance of a linear combination (two variables)

$$\text{Var}(aX_i + bX_j + c) = a^2 \text{Var}(X_i) + b^2 \text{Var}(X_j) + 2ab \text{Cov}(X_i, X_j).$$

5. Covariance of linear combinations (general case)

For constants a_1, \dots, a_n and b_1, \dots, b_n ,

$$\text{Cov}\left(c + \sum_{i=1}^n a_i X_i, d + \sum_{j=1}^n b_j X_j\right) = \sum_{i=1}^n \sum_{j=1}^n a_i b_j \text{Cov}(X_i, X_j).$$

This can be written as

$$\sum_{i=1}^n a_i^2 \text{Var}(X_i) + 2 \sum_{1 \leq i < j \leq n} a_i b_j \text{Cov}(X_i, X_j).$$

6. Variance of a linear combination (general case)

(Equation 4 is a special case of the following.)

$$\text{Var}\left(c + \sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n \sum_{j=1}^n a_i a_j \text{Cov}(X_i, X_j).$$

Equivalently,

$$\text{Var}\left(c + \sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i^2 \text{Var}(X_i) + 2 \sum_{1 \leq i < j \leq n} a_i a_j \text{Cov}(X_i, X_j).$$

Correlation

Note: **Covariance has no natural scale.**

The **correlation** of X_i and X_j is a standardized measure of association, defined as

$$\text{Corr}(X_i, X_j) = \rho_{X_i, X_j} = \frac{\text{Cov}(X_i, X_j)}{\sqrt{\text{Var}(X_i)}\sqrt{\text{Var}(X_j)}} = \frac{\text{Cov}(X_i, X_j)}{\sigma_{X_i}\sigma_{X_j}},$$

provided that

$$\text{Var}(X_i) < \infty \quad \text{and} \quad \text{Var}(X_j) < \infty.$$

Theorem

1. Bounds

$$-1 \leq \rho_{X,Y} \leq 1.$$

2. Perfect correlation

$$\rho_{X,Y} = \pm 1 \quad \text{if and only if} \quad P(X = aY + b) = 1,$$

for some constants a and b , where

- $a > 0$ if $\rho_{X,Y} = +1$,
- $a < 0$ if $\rho_{X,Y} = -1$.

Moment Generating Functions

The **joint moment generating function (mgf)** of (X_1, \dots, X_n) is defined as

$$M_{X_1, \dots, X_n}(t_1, \dots, t_n) = \mathbb{E}[e^{t_1 X_1 + \dots + t_n X_n}], \quad t_1, \dots, t_n \in \mathbb{R},$$

provided the expectation exists for all

$$-h < t_1, \dots, t_n < h$$

for some $h > 0$.

Marginal MGFs from the Joint MGF

The joint mgf can be used to obtain **univariate mgfs**. Specifically,

$$M_{X_i}(t_i) = M_{X_1, \dots, X_n}(t_1 = 0, \dots, t_{i-1} = 0, t_i, t_{i+1} = 0, \dots, t_n = 0).$$

Applications (as before)

- **Characterizes distributions**

For example, if (X_1, \dots, X_n) and (Y_1, \dots, Y_n) have the same joint mgf, then these random vectors have the same distribution.

- **Transformations**

For example, the mgf of $(a_1 X_1, \dots, a_n X_n)$ is

$$M_{X_1, \dots, X_n}(a_1 t_1, \dots, a_n t_n).$$

Conditional Distributions

Recall that $P(A | B)$ is the probability that A occurs given that B occurs:

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

We want to apply this idea to **random variables**.

Example: Truncated Distributions

Fix x_0 . For $x > x_0$,

$$\begin{aligned} P(X \leq x | X > x_0) &= \frac{P(X \leq x, X > x_0)}{P(X > x_0)} \\ &= \frac{P(x_0 < X \leq x)}{P(X > x_0)} \\ &= \frac{F(x) - F(x_0)}{1 - F(x_0)}. \end{aligned}$$

Definitions: Conditional Distribution (Given a General Event A)

Suppose we observe an event A with

$$P(A) > 0.$$

Conditional CDF

The **conditional cdf** of X given A is

$$F(x | A) = P(X \leq x | A) = \frac{P(A, X \leq x)}{P(A)}.$$

Conditional PMF / PDF

- **Discrete case (pmf):**

$$f(x | A) = P(X = x | A) = \frac{P(A, X = x)}{P(A)}, \quad x \in \mathbb{R}.$$

- **Continuous case (pdf):**

$$f(x | A) = \frac{d}{dx} F(x | A), \quad x \in \mathbb{R}.$$

Discrete Case

For jointly discrete random variables X, Y with joint pmf $f(x, y)$, for each x with $f_X(x) > 0$, define

$$f(y | x) = P(Y = y | X = x) = \frac{f(x, y)}{f_X(x)},$$

called the **conditional pmf of Y given $X = x$** , which specifies a distribution of Y given $X = x$.

Similarly, for each y with $f_Y(y) > 0$, define

$$f(x | y) = P(X = x | Y = y) = \frac{f(x, y)}{f_Y(y)},$$

called the **conditional pmf of X given $Y = y$** .

Continuous Distributions

For jointly continuous random variables X, Y with joint pdf $f(x, y)$, for each x with $f_X(x) > 0$, define the **conditional pdf of Y given $X = x$** as

$$f(y | x) = \frac{f(x, y)}{f_X(x)},$$

which specifies the pdf for the continuous distribution of Y given $X = x$.

For each y with $f_Y(y) > 0$, define the **conditional pdf of X given $Y = y$** as

$$f(x | y) = \frac{f(x, y)}{f_Y(y)}.$$

Conditional CDF

The conditional cdf of X given $Y = y$ is

$$F(x | y) = P(X \leq x | Y = y) = \int_{-\infty}^x f(t | y) dt = \int_{-\infty}^x \frac{f(t, y)}{f_Y(y)} dt.$$

Technical note. For a continuous random variable Y , $P(Y = y) = 0$ for any y . Thus, $F(x | y) = P(X \leq x | Y = y)$ is defined as a limit; this is well-defined whenever $f_Y(y) > 0$.

Conditional Expectations

Definition

Suppose X, Y are jointly discrete or jointly continuous, $g(x, y)$ is a real-valued function, and x is such that $f(y | x)$ is defined.

The **conditional mean (conditional expected value)** of $g(X, Y)$ given $X = x$ is

$$\mathbb{E}[g(X, Y) | X = x] = \begin{cases} \int g(x, y) f(y | x) dy, & \text{continuous case,} \\ \sum_y g(x, y) f(y | x), & \text{discrete case,} \end{cases}$$

provided that

$$\int |g(x, y)| f(y | x) dy < \infty \quad \text{or} \quad \sum_y |g(x, y)| f(y | x) < \infty.$$

Note. $\mathbb{E}[g(X, Y) | X = x]$ is a **function of x** .

Common Conditional Expectations

1. Conditional mean

$$\mathbb{E}[X | Y = y].$$

2. Conditional variance

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

Properties of Conditional Expectations

Result 1

For X, Y either jointly continuous or jointly discrete with $f_X(x) > 0$,

$$\mathbb{E}[ag(X, Y) + bh(X, Y) + c | X = x] = a\mathbb{E}[g(X, Y) | X = x] + b\mathbb{E}[h(X, Y) | X = x] + c,$$

and

$$\mathbb{E}[g(X) h(X, Y) | X = x] = g(x) \mathbb{E}[h(X, Y) | X = x],$$

provided the expectations exist.

These results also hold for the **random-variable version** of conditional expectations.

Result 2

Provided the necessary conditional means exist for a set of x with probability 1, the following equalities hold with probability 1:

$$\mathbb{E}[ag(X, Y) + bh(X, Y) + c | X] = a\mathbb{E}[g(X, Y) | X] + b\mathbb{E}[h(X, Y) | X] + c,$$

and

$$\mathbb{E}[g(X) h(X, Y) | X] = g(X) \mathbb{E}[h(X, Y) | X].$$

That is, both sides are random variables equal with probability 1.

Law of Total Expectation If X and Y are random variables and $\mathbb{E}[Y]$ exists, then

$$\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y | X]].$$

Interpretation. If $m(x) = \mathbb{E}[Y | X = x]$, then $m(X) = \mathbb{E}[Y | X]$ is a random variable. Averaging over X yields $\mathbb{E}[Y]$.

Law of Total Variance (EVVE) If $\text{Var}(Y)$ exists, then

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

That is, the unconditional variance equals the **mean of the conditional variance plus the variance of the conditional mean**.

Hierarchical and Mixture Models

Key Idea Joint distributions can be specified **hierarchically** (or conditionally):

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

This perspective underlies: - hierarchical models, - mixture distributions, - Bayesian models, - mixed discrete-continuous models.

Bayesian Statistics (as a Hierarchical Model)

Let: - θ = parameter (random), - X = observed data.

Bayes' rule in density form:

$$f(\theta | x) = \frac{f(x, \theta)}{f_X(x)} = \frac{f(x | \theta)f(\theta)}{\int f(x | \theta)f(\theta) d\theta}.$$

Interpretation: - $f(\theta)$ = prior, - $f(x | \theta)$ = likelihood, - $f(\theta | x)$ = posterior.

Mixed Discrete-Continuous Models

It is possible to have: - X discrete, - $Y | X = x$ continuous.

Joint distribution still defined as:

$$f(x, y) = f_X(x) f_{Y|X}(y | x).$$

Independence of Random Variables

Definition

Random variables X and Y are **independent** if

$$P(X \in A, Y \in B) = P(X \in A)P(Y \in B) \quad \text{for all } A, B \subset \mathbb{R}.$$

Equivalent Characterizations

CDF Factorization

$$F_{X,Y}(x,y) = F_X(x)F_Y(y) \quad \text{for all } x, y.$$

Discrete Case

$$f(x, y) = f_X(x)f_Y(y).$$

Continuous Case

There exist versions of the densities such that

$$f(x, y) = f_X(x)f_Y(y).$$

Consequences for Conditional Distributions

If X and Y are independent, then:

$$f(y | x) = f_Y(y), \quad f(x | y) = f_X(x).$$

That is, conditioning does **nothing** under independence.

Independence and Functions of Random Variables

Theorem

If X and Y are independent, then:

$$g(X) \text{ and } h(Y) \text{ are independent}$$

for any measurable functions g, h .

Expectations Under Independence

If expectations exist:

$$\mathbb{E}[g(X)h(Y)] = \mathbb{E}[g(X)]\mathbb{E}[h(Y)].$$

Special case:

$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y].$$

Covariance and Independence

If X and Y are independent and $\mathbb{E}[XY]$ exists:

$$\text{Cov}(X, Y) = 0.$$

Important:

Zero covariance **does not imply** independence.

Variance of Sums Under Independence

If X and Y are independent and second moments exist:

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

More generally, if X_1, \dots, X_n are independent:

$$\text{Var}\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i^2 \text{Var}(X_i).$$

Moment Generating Functions (MGFs) and Independence

Univariate MGF

$$M_X(t) = \mathbb{E}[e^{tX}],$$

defined for t in a neighborhood of 0.

Joint MGF

For (X, Y) :

$$M_{X,Y}(t_1, t_2) = \mathbb{E}[e^{t_1 X + t_2 Y}].$$

Independence via MGFs

Corollary

If X and Y are independent and both MGFs exist:

$$M_{X+Y}(t) = M_X(t) M_Y(t).$$

Characterization Theorem

Suppose $M_X(t)$ and $M_Y(t)$ exist in neighborhoods of 0.

Then X and Y are independent **if and only if**

$$M_{X,Y}(t_1, t_2) = M_X(t_1) M_Y(t_2)$$

for all (t_1, t_2) near $(0, 0)$.

Sums of Independent Random Variables

The same technique shows that, for independent X_1, \dots, X_n :

1. Normal

If each $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$, then

$$S = \sum_{i=1}^n X_i \sim \mathcal{N}\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right).$$

2. Binomial (common p)

If each $X_i \sim \text{Binomial}(n_i, p)$, then

$$S = \sum_{i=1}^n X_i \sim \text{Binomial}\left(\sum_{i=1}^n n_i, p\right).$$

3. Gamma (common β)

If each $X_i \sim \text{Gamma}(\alpha_i, \beta)$, then

$$S = \sum_{i=1}^n X_i \sim \text{Gamma}\left(\sum_{i=1}^n \alpha_i, \beta\right).$$

Recall the Gamma density (shape-scale form):

$$f_X(x) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta}, \quad x > 0.$$

Special Cases

(a) Exponential

If each $X_i \sim \text{Exponential}(\beta) \sim \text{Gamma}(1, \beta)$, then

$$S = \sum_{i=1}^n X_i \sim \text{Gamma}(n, \beta).$$

(b) Chi-square

If $X_i \sim \chi_{\nu_i}^2 \sim \text{Gamma}\left(\frac{\nu_i}{2}, 2\right)$, then

$$S = \sum_{i=1}^n X_i \sim \chi_{\sum_{i=1}^n \nu_i}^2.$$

The mgf of a χ_{ν}^2 random variable is

$$M_X(t) = (1 - 2t)^{-\nu/2}.$$

4. Poisson

If each $X_i \sim \text{Poisson}(\lambda_i)$, then

$$S = \sum_{i=1}^n X_i \sim \text{Poisson}\left(\sum_{i=1}^n \lambda_i\right).$$

5. Negative Binomial (common p)

If each $X_i \sim \text{Neg-Binomial}(r_i, p)$, then

$$S = \sum_{i=1}^n X_i \sim \text{Neg-Binomial}\left(\sum_{i=1}^n r_i, p\right).$$

(a) Geometric

If each $X_i \sim \text{Geometric}(p) \sim \text{Neg-Binomial}(1, p)$, then

$$S = \sum_{i=1}^n X_i \sim \text{Neg-Binomial}(n, p).$$

Transformations of Random Variables

Assume the transformation is **one-to-one** with inverse functions:

$$x_i = u_i^{-1}(y_1, \dots, y_n), \quad i = 1, \dots, n.$$

For $n = 2$:

$$f_{U,V}(u, v) = f_{X,Y}(h(u, v), h(u, v)) | \text{Jacobian term}|.$$

In general,

$$f_{Y_1, \dots, Y_n}(y_1, \dots, y_n) = f_{X_1, \dots, X_n}(u^{-1}(y_1, \dots, y_n)) |J|.$$

Jacobian

Define the Jacobian as

$$J = \det \begin{pmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} & \cdots & \frac{\partial x_1}{\partial y_n} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} & \cdots & \frac{\partial x_2}{\partial y_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_n}{\partial y_1} & \frac{\partial x_n}{\partial y_2} & \cdots & \frac{\partial x_n}{\partial y_n} \end{pmatrix}.$$

Equivalently,

$$J = \det \begin{pmatrix} \frac{\partial u_1^{-1}}{\partial y_1} & \frac{\partial u_1^{-1}}{\partial y_2} & \cdots & \frac{\partial u_1^{-1}}{\partial y_n} \\ \frac{\partial u_2^{-1}}{\partial y_1} & \frac{\partial u_2^{-1}}{\partial y_2} & \cdots & \frac{\partial u_2^{-1}}{\partial y_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial u_n^{-1}}{\partial y_1} & \frac{\partial u_n^{-1}}{\partial y_2} & \cdots & \frac{\partial u_n^{-1}}{\partial y_n} \end{pmatrix}.$$

If J is continuous and $J \neq 0$ on B (except possibly on a set of probability zero), then

$$f_{Y_1, \dots, Y_n}(y_1, \dots, y_n) = f_{X_1, \dots, X_n}(u^{-1}(y_1, \dots, y_n)) |J|, \quad (y_1, \dots, y_n) \in B.$$

Partial Transformations

Often we are interested in only one transformation:

$$Y_1 = u_1(X_1, \dots, X_n).$$

Then choose convenient definitions to complete the transformation, e.g.

$$Y_2 = X_2, \dots, Y_n = X_n.$$

If the transformation is **not one-to-one**, partition the support A of (X_1, \dots, X_n) into sets A_i where the transformation is one-to-one, and sum:

$$f_Y(y) = \sum_{i=1}^k f_X(u_i^{-1}(y)) |J_i|.$$

Independence Notes**

- If X_1 and X_2 are independent, then $h(X_1)$ and $w(X_2)$ are independent.
- If X_1 and X_2 are independent, then $h(X_1, X_2)$ and $w(X_1, X_2)$ are **not necessarily independent**.

Multivariate Continuous Case: Convolutions

Example: Sum of Two Continuous Random Variables

Let (X_1, X_2) have joint density $f_{X_1, X_2}(x_1, x_2)$.

Define the transformation:

$$S = X_1 + X_2, \quad T = X_2.$$

Inverse transformation:

$$X_1 = S - T, \quad X_2 = T.$$

The Jacobian is

$$J = \det \begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix} = 1.$$

The joint density of (S, T) is

$$f_{S,T}(s, t) = f_{X_1, X_2}(s - t, t).$$

The marginal density of S is

$$f_S(s) = \int_{-\infty}^{\infty} f_{X_1, X_2}(s - t, t) dt.$$

If X_1 and X_2 are independent, then

$$f_S(s) = \int_{-\infty}^{\infty} f_{X_1}(s - t) f_{X_2}(t) dt,$$

which is called the **convolution formula**.

Multivariate Distributions

Multinomial Distribution

- Suppose $0 \leq p_1, p_2, \dots, p_k$ are probabilities such that

$$\sum_{i=1}^k p_i = 1.$$

- Consider a series of n identical trials where, on each trial, one can get **exactly one** of k possible outcomes o_1, \dots, o_k .
- Let $X_i =$ number of trials resulting in outcome o_i .
- Then

$$X = (X_1, \dots, X_k)$$

has a Multinomial(n, p_1, \dots, p_k) distribution with joint pmf

$$P(X_1 = x_1, \dots, X_k = x_k) = f(x_1, \dots, x_k) = \begin{cases} \frac{n!}{x_1!x_2!\cdots x_k!} p_1^{x_1} p_2^{x_2} \cdots p_k^{x_k}, & x_i \geq 0, \sum_{i=1}^k x_i = n, \\ 0, & \text{otherwise.} \end{cases}$$

- The quantity

$$\frac{n!}{x_1!x_2!\cdots x_k!}$$

is the **multinomial coefficient**, counting the number of arrangements of n trials resulting in x_1 outcomes of o_1 , x_2 of o_2 , \dots , x_k of o_k .

- Individual marginal distributions are

$$X_i \sim \text{Binomial}(n, p_i).$$

- Conditional distributions of subsets given others are again multinomial.

Dirichlet Distribution

- Suppose Y_1, \dots, Y_k are independent with

$$Y_i \sim \text{Gamma}(\alpha_i, 1), \quad i = 1, \dots, k.$$

- Define

$$X_i = \frac{Y_i}{\sum_{j=1}^k Y_j}, \quad i = 1, \dots, k.$$

- Then

$$X = (X_1, \dots, X_k)$$

has a $\text{Dirichlet}(\alpha_1, \dots, \alpha_k)$ distribution.

- Note that

$$0 < X_i < 1, \quad \sum_{i=1}^k X_i = 1.$$

- Individual marginal distributions are

$$X_i \sim \text{Beta}\left(\alpha_i, \sum_{j \neq i} \alpha_j\right).$$

- Conditionals of subsets given others are (scaled) Dirichlet distributions.

Building to MVN

Matrix–Vector Multivariate Notation

- For a $p \times q$ matrix

$$A = [A_{ij}]_{i=1, \dots, p; j=1, \dots, q}$$

of random variables, $E(A)$ denotes the matrix of componentwise expectations.

- For a $k \times 1$ random vector,

$$X = (X_1, \dots, X_k)'.$$

- The expected value of X is

$$E(X) = \mu_X = \mu = (E[X_1], \dots, E[X_k])'.$$

Variance–Covariance Matrix

Definition. The variance–covariance matrix of $X = (X_1, \dots, X_k)'$ is the $k \times k$ matrix

$$\text{Var}(X) = \Sigma = \begin{pmatrix} \sigma_{11} & \cdots & \sigma_{1k} \\ \vdots & \ddots & \vdots \\ \sigma_{k1} & \cdots & \sigma_{kk} \end{pmatrix},$$

where

$$\sigma_{ij} = \text{Cov}(X_i, X_j).$$

Result.

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

Covariance Between Two Random Vectors

Let

$$X = (X_1, \dots, X_k)', \quad Y = (Y_1, \dots, Y_m)'.$$

Then

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)'],$$

which is a $k \times m$ matrix whose (i, j) entry is

$$\text{Cov}(X_i, Y_j).$$

Linear Transformations: Expectation and Variance

Let $A_{r \times k}$ and $B_{s \times m}$ be fixed matrices, $a_r \in \mathbb{R}^r$, $b_s \in \mathbb{R}^s$, and $c, d \in \mathbb{R}$.

1. Expectation

$$E(a_r + AX) = a_r + AE(X).$$

In one dimension:

$$E\left(c + \sum_{i=1}^k a_i X_i\right) = c + \sum_{i=1}^k a_i E(X_i).$$

2. Variance

$$\text{Var}(a_r + AX) = A\Sigma A'.$$

In one dimension:

$$\text{Var}\left(\sum_{i=1}^k a_i X_i\right) = \sum_{i=1}^k \sum_{j=1}^k a_i a_j \text{Cov}(X_i, X_j).$$

3. Covariance

$$\text{Cov}(a_r + AX, b_s + BY) = A \text{Cov}(X, Y) B'.$$

Matrix Definitions

Let B be a $k \times k$ matrix.

1. **Non-singular** if $\text{rank}(B) = k$ (equivalently $\det(B) \neq 0$ or B^{-1} exists).
2. **Singular** if $\text{rank}(B) < k$.
3. **Non-negative definite** if

$$a' B a \geq 0 \quad \text{for all } a \in \mathbb{R}^k.$$

4. **Positive definite** if

$$a' B a > 0 \quad \text{for all } a \neq 0.$$

Facts. - Positive definite \Rightarrow non-negative definite. - Non-negative definite and non-singular \Rightarrow positive definite. - B is positive definite iff $\det(B) > 0$.

Lemma. If $\Sigma = \text{Var}(X)$, then Σ is symmetric and non-negative definite. If Σ is not positive definite, then X lies in a hyperplane

$$\{x \in \mathbb{R}^k : a' x = b\}$$

with probability 1 for some $a \neq 0$.

Multivariate Normal Distribution

Definition. A random vector

$$X = (X_1, \dots, X_k)'$$

is said to have a multivariate normal distribution

$$X \sim \text{MVN}_k(\mu, \Sigma)$$

if

$$X = \mu + P'Z,$$

where

- $Z = (Z_1, \dots, Z_s)'$ with $Z_i \stackrel{\text{iid}}{\sim} N(0, 1)$,
- P is an $s \times k$ matrix such that

$$P'P = \Sigma.$$

Then

$$E(X) = \mu, \quad \text{Var}(X) = \Sigma.$$

- $\mu \in \mathbb{R}^k$ is the mean vector.
- Σ is symmetric and non-negative definite.
- For any non-negative definite Σ , a matrix P satisfying $P'P = \Sigma$ exists.
- If Σ is positive definite, P is $k \times k$ and non-singular.

Properties of the Multivariate Normal

1. Linear combinations of an MVN vector are again normal.
2. If all linear combinations $a'X$ are normal, then X is MVN.
3. Subvectors of an MVN vector are MVN.
4. If

$$X = \begin{pmatrix} X^{(1)} \\ X^{(2)} \end{pmatrix}$$

is MVN, then

$$X^{(1)} \perp X^{(2)} \iff \text{Cov}(X^{(1)}, X^{(2)}) = 0.$$

(This is a **strong** property.)

5. If Σ is non-singular, the joint pdf of X exists.
6. Conditional distributions of subvectors are again MVN.

Moment Generating Function

If

$$X \sim \text{MVN}_k(\mu, \Sigma),$$

then the mgf is

$$M_X(t) = E(e^{t'X}) = \exp\left(t'\mu + \frac{1}{2}t'\Sigma t\right), \quad t \in \mathbb{R}^k.$$

Transformation Results

Result 1.

If

$$X \sim \text{MVN}_k(\mu, \Sigma)$$

and

$$Y = a + BX,$$

where $a \in \mathbb{R}^m$ and B is $m \times k$, then

$$Y \sim \text{MVN}_m(a + B\mu, B\Sigma B').$$

Result 2.

Partition

$$X = \begin{pmatrix} X^{(1)} \\ X^{(2)} \end{pmatrix}, \quad \mu = \begin{pmatrix} \mu^{(1)} \\ \mu^{(2)} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}.$$

Then

$$X^{(1)} \sim \text{MVN}_p(\mu^{(1)}, \Sigma_{11}), \quad X^{(2)} \sim \text{MVN}_{k-p}(\mu^{(2)}, \Sigma_{22}).$$

Result 3.

X is MVN if and only if

$$a'X = \sum_{i=1}^k a_i X_i$$

is normal for every $a \in \mathbb{R}^k$.

Independence and Covariance in the MVN

- Independence \Rightarrow zero covariance
- Zero covariance + MVN \Rightarrow independence

Result 4

Independence of Subvectors in the MVN

Suppose

$$X \sim \text{MVN}_k(\mu, \Sigma),$$

and partition X , μ , and Σ as

$$X = \begin{pmatrix} X^{(1)} \\ X^{(2)} \end{pmatrix}, \quad \mu = \begin{pmatrix} \mu^{(1)} \\ \mu^{(2)} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix},$$

where $X^{(1)}$ is $p \times 1$ and $X^{(2)}$ is $(k-p) \times 1$.

Then,

$$X^{(1)} \text{ and } X^{(2)} \text{ are independent} \iff \text{Cov}(X^{(1)}, X^{(2)}) = \Sigma_{12} = 0.$$

Independent Normals as a Special Case of MVN Sometimes we begin with X_1, \dots, X_k independent and

$$X_i \sim N(\mu_i, \sigma_i^2).$$

In this case,

$$X = \begin{pmatrix} X_1 \\ \vdots \\ X_k \end{pmatrix} \sim \text{MVN}_k(\mu, \Sigma),$$

where

$$\mu = (\mu_1, \dots, \mu_k)', \quad \Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_k^2).$$

Linear Combinations of Independent Normals Let

$$U = \sum_{i=1}^k a_i X_i, \quad V = \sum_{j=1}^k b_j X_j.$$

Then U and V are normal, and

$$\text{Cov}(U, V) = \sum_{i=1}^k \sum_{j=1}^k a_i b_j \text{Cov}(X_i, X_j) = \sum_{i=1}^k a_i b_i \sigma_i^2,$$

since

$$\text{Cov}(X_i, X_j) = \begin{cases} 0, & i \neq j, \\ \sigma_i^2, & i = j. \end{cases}$$

Thus,

$$U \text{ and } V \text{ are independent} \iff \sum_{i=1}^k a_i b_i \sigma_i^2 = 0.$$

If all $\sigma_i^2 = \sigma^2$, then

$$U \perp V \iff \sum_{i=1}^k a_i b_i = 0.$$

Result 5

PDF of the Multivariate Normal

If

$$X \sim \text{MVN}_k(\mu, \Sigma)$$

with **non-singular** Σ , then the joint pdf of X is

$$f_X(x) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu)' \Sigma^{-1} (x - \mu)\right\}, \quad x \in \mathbb{R}^k.$$

Result 6

Conditional Distribution of an MVN

Suppose

$$X \sim \text{MVN}_k(\mu, \Sigma),$$

and partition as

$$X = \begin{pmatrix} X^{(1)} \\ X^{(2)} \end{pmatrix}, \quad \mu = \begin{pmatrix} \mu^{(1)} \\ \mu^{(2)} \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix},$$

where $X^{(1)}$ is $p \times 1$ and $X^{(2)}$ is $(k-p) \times 1$.

Then the conditional distribution of $X^{(1)} | X^{(2)} = x^{(2)}$ is

$$X^{(1)} | X^{(2)} = x^{(2)} \sim \text{MVN}_p\left(\mu^{(1)} + \Sigma_{12}\Sigma_{22}^{-1}(x^{(2)} - \mu^{(2)}), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\right).$$

Result 7

Quadratic Forms in the MVN

Suppose

$$X \sim \text{MVN}_k(\mu, \Sigma)$$

with **non-singular** Σ . Then

$$Q(X) = (X - \mu)' \Sigma^{-1} (X - \mu) \sim \chi_k^2.$$

Supporting Lemma

If

$$Z \sim N(0, 1),$$

then

$$Z^2 \sim \chi_1^2.$$

Random Samples

Definition (Random Sample).

If X_1, \dots, X_n are independent and identically distributed (iid) with

$$X_i \sim f_X(x),$$

then we call X_1, \dots, X_n a **random sample** from the population with pdf (or pmf) $f_X(x)$.

Statistics and Sampling Distributions

Let

$$Y = T(X_1, \dots, X_n).$$

- Y is called a **statistic**, i.e., a function of the data computable from the sample.
- The distribution of a statistic Y is called the **sampling distribution** of Y .

Properties of the Sample Mean

Let X_1, \dots, X_n be a random sample from $f_X(x)$ with

$$\mu = E(X_i), \quad \sigma^2 = \text{Var}(X_i).$$

Define the sample mean

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

Important Results

1. Expectation

$$E(\bar{X}_n) = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \mu.$$

2. Variance

$$\text{Var}(\bar{X}_n) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{\sigma^2}{n},$$

since the X_i are independent.

MGF Approach for \bar{X}_n

Sometimes the exact distribution of \bar{X}_n can be found using moment generating functions.

$$M_{\bar{X}_n}(t) = E(e^{t\bar{X}_n}) = E\left(e^{\frac{t}{n}(X_1 + \dots + X_n)}\right) = \prod_{i=1}^n E\left(e^{\frac{t}{n}X_i}\right) = \left[M_X\left(\frac{t}{n}\right)\right]^n.$$

Distribution of the Sample Variance

Let X_1, \dots, X_n be a random sample with

$$\mu = E(X_i), \quad \sigma^2 = \text{Var}(X_i).$$

Define the sample variance

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}_n^2 \right).$$

- The **exact sampling distribution** of S^2 is difficult to obtain in general.
- If

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2),$$

then, after appropriate scaling,

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2.$$

Result.

$$E(S^2) = \sigma^2 = \text{Var}(X_i).$$

Order Statistics

Introduction - Distribution of the Maximum and Minimum

Let X_1, \dots, X_n be a random sample with common cdf

$$F_X(x) = P(X_1 \leq x).$$

Define

$$X_{(n)} = \max\{X_1, \dots, X_n\}, \quad X_{(1)} = \min\{X_1, \dots, X_n\}.$$

Important Results

1. CDF of the Maximum

$$F_{X_{(n)}}(x) = P(X_{(n)} \leq x) = P(X_1 \leq x, \dots, X_n \leq x) = [F_X(x)]^n.$$

2. CDF of the Minimum

$$F_{X_{(1)}}(x) = P(X_{(1)} \leq x) = 1 - P(X_1 > x, \dots, X_n > x) = 1 - [1 - F_X(x)]^n.$$

3. PDFs (Continuous Case)

If $F_X(x)$ is continuous with pdf $f_X(x)$, then

$$f_{X(n)}(x) = n f_X(x) [F_X(x)]^{n-1},$$

and

$$f_{X(1)}(x) = n f_X(x) [1 - F_X(x)]^{n-1}.$$

Generalizing

Definition.

The **order statistics** of a sample X_1, \dots, X_n are the ordered values

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}.$$

- Typically assume X_1, \dots, X_n are iid with a **continuous** distribution.
- We may be interested in:
 1. The distribution of a single order statistic $X_{(k)}$
 2. The joint distribution of two or more order statistics $(X_{(i)}, X_{(j)})$
 3. Functions of order statistics (e.g., the range $R = X_{(n)} - X_{(1)}$)
- Order statistics can be viewed as a (discontinuous) transformation of (X_1, \dots, X_n) .

Result 1: CDF of the k th Order Statistic

If X_1, \dots, X_n are a random sample with common cdf $F_X(x)$, then for $k = 1, \dots, n$,

$$F_{X(k)}(x) = P(X_{(k)} \leq x) = P(\text{at least } k \text{ of the } X_i \leq x) = \sum_{j=k}^n \binom{n}{j} [F_X(x)]^j [1 - F_X(x)]^{n-j}.$$

Result 2: PDF of the k th Order Statistic (Continuous Case)

If X_1, \dots, X_n are a random sample with continuous cdf $F_X(x)$ and pdf $f_X(x)$, then

$$f_{X(k)}(x) = \frac{n!}{(k-1)!(n-k)!} f_X(x) [F_X(x)]^{k-1} [1 - F_X(x)]^{n-k}.$$

Joint Distributions of Order Statistics

Result 1

Joint pdf of Two Order Statistics

Let X_1, \dots, X_n be iid with continuous cdf $F_X(x)$ and pdf $f_X(x)$.

For $i < j$, the joint pdf of $(X_{(i)}, X_{(j)})$ is

$$f_{X_{(i)}, X_{(j)}}(u, v) = \frac{n!}{(i-1)!(j-i-1)!(n-j)!} f_X(u)f_X(v) [F_X(u)]^{i-1} [F_X(v) - F_X(u)]^{j-i-1} [1 - F_X(v)]^{n-j},$$

for $u < v$.

Result 2

Joint pdf of All Order Statistics

Let X_1, \dots, X_n be iid with continuous pdf $f_X(x)$. Then the joint pdf of

$$(X_{(1)}, X_{(2)}, \dots, X_{(n)})$$

is

$$f_{X_{(1)}, \dots, X_{(n)}}(u_1, \dots, u_n) = n! \prod_{i=1}^n f_X(u_i),$$

for

$$u_1 < u_2 < \dots < u_n.$$

The factor $n!$ counts the number of permutations of (X_1, \dots, X_n) that produce the same ordered sample.

Sampling from the Normal Distribution

Joint Distribution of \bar{X}_n and S^2

If

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2),$$

then:

1. The sample mean satisfies

$$\bar{X}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right).$$

2. The sample variance satisfies

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2.$$

3. \bar{X}_n and S^2 are **independent**.

Derived Distributions

Student's t

Definition

Let

$$Z \sim N(0, 1), \quad V \sim \chi_{\nu}^2,$$

with Z and V independent. Define

$$T = \frac{Z}{\sqrt{V/\nu}}.$$

Then T has a **Student's t distribution** with ν degrees of freedom, denoted

$$T \sim t_{\nu}.$$

Properties of the t Distribution

1. The t density is symmetric about 0.
2. $E(T) = 0$ for $\nu > 1$.
3. The variance is

$$\text{Var}(T) = \frac{\nu}{\nu - 2}, \quad \nu > 2,$$

and is infinite for $\nu \leq 2$.

4. For $\nu = 1$, the t distribution reduces to the **Cauchy distribution** with density

$$f(t) = \frac{1}{\pi(1+t^2)}.$$

5. As $\nu \rightarrow \infty$, the t_{ν} distribution converges to $N(0, 1)$.

Snedecor's F

Definition

Let

$$V_1 \sim \chi_{\nu_1}^2, \quad V_2 \sim \chi_{\nu_2}^2,$$

with V_1 and V_2 independent. Define

$$X = \frac{V_1/\nu_1}{V_2/\nu_2}.$$

Then

$$X \sim F_{\nu_1, \nu_2}.$$

Note.

$$F_{\nu_1, \nu_2} \neq F_{\nu_2, \nu_1}, \quad X^{-1} \sim F_{\nu_2, \nu_1}.$$

Properties of the F Distribution

Let

$$X = \frac{\chi_{\nu_1}^2 / \nu_1}{\chi_{\nu_2}^2 / \nu_2} \sim F_{\nu_1, \nu_2}.$$

1. Expectation

$$E(X) = \frac{\nu_2}{\nu_2 - 2}, \quad \nu_2 > 2.$$

2. If $\nu_1 = 1$, then

$$X \sim t_{\nu_2}^2.$$

3. The reciprocal satisfies

$$X^{-1} \sim F_{\nu_2, \nu_1}.$$

4. If $\nu_2 \rightarrow \infty$, then

$$X \sim \frac{\chi_{\nu_1}^2}{\nu_1}.$$

5. The transformation

$$\frac{(\nu_1/\nu_2)X}{1 + (\nu_1/\nu_2)X}$$

has a

$$\text{Beta}\left(\frac{\nu_1}{2}, \frac{\nu_2}{2}\right)$$

distribution. Equivalently,

$$\frac{\chi_{\nu_1}^2}{\chi_{\nu_1}^2 + \chi_{\nu_2}^2} \sim \text{Beta}\left(\frac{\nu_1}{2}, \frac{\nu_2}{2}\right).$$

Inequalities (Again)

Convergence