

Lecture Notes for Math 472: Statistical Inference

Last updated: January 30, 2026

Contents

0	Tentative course outline	2
1	2025-01-12 Week 01 Lecture 01	3
1.1	What is probability?	3
1.1.1	A general framework: sample space, events, etc	3
1.1.2	Definition of probability measure	4
2	2026-01-14 Week 01 Lecture 02	5
2.1	Independent events and conditional probabilities	5
2.2	Random variables	6
3	2026-01-16 Week 01 Lecture 03	7
3.1	Random variables	7
3.1.1	Discrete vs continuous	7
3.1.2	Expected value	7
3.1.3	Joint distributions	8
4	2026-01-21 Week 02 Lecture 04	9
4.1	Independence of random variables	9
4.1.1	Definition and characterization	9
4.1.2	Independence and expectations/variances	9
4.1.3	Independence and long-term behavior: the weak law of large numbers	10
5	2026-01-23 Week 02 Lecture 05	11
5.1	Proof of the weak law of large numbers	11
5.2	Equality in distribution	11
5.3	Populations and samples	12
6	2026-01-26 Week 03 Lecture 06	13
6.1	Sampling	13
6.2	What is a statistic?	13
7	2026-01-28 Week 03 Lecture 07	15
7.1	Basic examples of statistics and estimators	15
8	2026-01-30 Week 03 Lecture 08	17
8.1	The Gaussian distribution	17
9	2026-02-02 Week 04 Lecture 09	19
9.1	Chi-squared distributions	19

0 Tentative course outline

This course is a problem-oriented introduction to the basic concepts of probability and statistics, providing a foundation for applications and further study.

1. **Weeks 1-2.** Sampling distributions (4 lessons).
chi-squared, t, and F distributions, distributions of sample mean and variance
2. **Weeks 3-4.** Point estimation (5 lessons)
properties and methods of point estimation
3. **Weeks 5-6.** Interval estimation (4 lessons)
Confidence intervals for means, variances, proportions and differences
4. **Weeks 7-12.** Hypothesis Testing (19 lessons)
Neyman-Pearson lemma, likelihood ratio test; tests concerning means and variances, tests based on count data, nonparametric tests, analysis of variance
5. **Weeks 13-14.** Regression and correlation (6 lessons)
regression, bivariate normal distributions, method of least squares

1 2025-01-12 | Week 01 | Lecture 01

- give syllabus
- do activity with why you're in this course

The nexus question of this lecture: What is a probability?

Reading assignment: Sections 1.1, 1.2, 1.3, 2.1, 2.4 of the textbook.

1.1 What is probability?

1.1.1 A general framework: sample space, events, etc

We begin with a general framework and some terminology to formalize the notions of probability. This is based on section 2.4 in the textbook.

- An **experiment** is an activity or process whose outcome is subject to uncertainty, and about which an observation is made.

Examples include flipping a coin, rolling a dice, measuring the size of a wave, or the amount of rainfall, conducting a poll, performing a diagnostic test, opening a pack of Pokemon cards, etc.

- The **sample space** S of an experiment is the set of all possible outcomes. The elements of the sample space are called **sample points**.

We think of each sample point as representing a unique outcome of the experiment. In the case of rolling a dice, the sample points are 1, 2, 3, 4, 5 and 6, and the sample space is $S = \{1, 2, 3, 4, 5, 6\}$.

- We use the term **event** to refer to a collection of outcomes, i.e., a subset of S .

Example: if our experiment is rolling a 6-sided dice, here are some events

$A = [\text{observe an odd number}]$	$E_2 = [\text{observe a 2}]$
$B = [\text{observe an even number}]$	$E_3 = [\text{observe a 3}]$
$C = [\text{observe a number less than 5}]$	$E_4 = [\text{observe a 4}]$
$D = [\text{observe a 2 or a 3}]$	$E_5 = [\text{observe a 5}]$
$E_1 = [\text{observe a 1}]$	$E_6 = [\text{observe a 6}]$

- There are two types of events: **compound events**, which can be decomposed into other events, and **simple events**, which cannot.

In the above example, the events A, B, C and D are compound events. E_1, \dots, E_6 are simple events.

- A sample space is **discrete** if it is countable (i.e., finite or countably infinite). In a discrete sample space S , the set of all possible events is the **power set** of S .¹

In the dice-rolling example, the set of all possible events is $\{E : E \subseteq \{1, 2, 3, 4, 5, 6\}\}$.

$A = [\text{observe an odd number}] = \{1, 3, 5\}$	$E_2 = [\text{observe a 2}] = \{2\}$
$B = [\text{observe an even number}] = \{2, 4, 6\}$	$E_3 = [\text{observe a 3}] = \{3\}$
$C = [\text{observe a number less than 5}] = \{1, 2, 3, 4\}$	$E_4 = [\text{observe a 4}] = \{4\}$
$D = [\text{observe a 2 or a 3}] = \{2, 3\}$	$E_5 = [\text{observe a 5}] = \{5\}$
$E_1 = [\text{observe a 1}] = \{1\}$	$E_6 = [\text{observe a 6}] = \{6\}$

¹If S is not discrete, a complication arises: in that case, some subsets of S are too wild and untameable for us to treat them mathematically as “events”. Resolving that issue requires introducing measure theory, which is beyond the scope of this class, so we will ignore it and simply steer clear of any setting where any issues might arise.

- Some observations about events:
 - The sample points are *elements* of S . The simple events are *singleton subsets* of S . In the dice example, we have:
 - * Sample points: 1,2,3,4,5,6.
 - * Simple events: $\{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{6\}$.
 - The empty set \emptyset and the whole sample space S are always both events: \emptyset is the event “nothing happens” and S is the event “something happens”.
 - Events satisfy the properties of a boolean algebra:
 - * “**And**”: If E and F are events, then $E \cap F$ is the event that E and F occur.
 - * “**Or**”: If E and F are events, then $E \cup F$ is the event that E or F occurs.
 - * “**Not**”: If E is an event, then $E^c = S \setminus E$ is the event that E does not occur.
- Two events E and F are **mutually exclusive** if $E \cap F = \emptyset$. This means that E and F cannot both happen at the same time.
 In the dice example, the events A and B are mutually exclusive, since the dice roll cannot be both even and odd. But A and C are not mutually exclusive because $A \cap C = \{1, 3\} \neq \emptyset$. If a 1 or a 3 is rolled, then both A and C occur.

1.1.2 Definition of probability measure

Definition 1 (Probability measure). Let S be a sample space associated with an experiment. A function \mathbb{P} is said to be a **probability measure** on S if it satisfies the following three axioms:

A.1 (Nonnegativity) For every event $E \subseteq S$,

$$\mathbb{P}[E] \geq 0.$$

A.2 (Total mass one) $\mathbb{P}[S] = 1$.

A.3 (Countable additivity) If E_1, E_2, \dots is a sequence of events which are pairwise mutually exclusive (meaning $E_i \cap E_j = \emptyset$ if $i \neq j$), then

$$\mathbb{P}[E_1 \cup E_2 \cup \dots] = \sum_{i=1}^{\infty} \mathbb{P}[E_i].$$

If \mathbb{P} is a probability measure, then for every event $E \subseteq S$, the number $\mathbb{P}[E]$ is called the **probability** of E .

The above definition only tells us the conditions an assignment of probabilities must satisfy; it doesn't tell us how to assign specific probabilities to events.

Probability measures satisfy some basic properties:

Proposition 2 (Basic properties of probability measure). *If \mathbb{P} is a probability measure, then the following properties hold:*

- (The null event has probability zero) $\mathbb{P}[\emptyset] = 0$.
- (Finite additivity) Let $\{E_1, \dots, E_n\}$ be a finite sequence of events. If the sequence is pairwise disjoint, then

$$\mathbb{P}[E_1 \cup E_2 \cup \dots \cup E_n] = \mathbb{P}[E_1] + \mathbb{P}[E_2] + \dots + \mathbb{P}[E_n].$$

- (“With probability one, an event E either does occur or doesn't”) $\mathbb{P}[E^c] = 1 - \mathbb{P}[E]$.

- (Excision Property) If A, B are events and $A \subseteq B$, then

$$\mathbb{P}[B \setminus A] = \mathbb{P}[B] - \mathbb{P}[A].$$

- (“The particular is less likely than the general”) If A, B are events and $A \subseteq B$, then $\mathbb{P}[A] \leq \mathbb{P}[B]$.

- (“Probabilities are between 0 and 1”) For any event E , $\mathbb{P}[E] \in [0, 1]$.

2 2026-01-14 | Week 01 | Lecture 02

The topic of this lecture: independent events, conditional probabilities, random variables

2.1 Independent events and conditional probabilities

This section is based on section 2.7 in the textbook.

Definition 3 (Independence). Two events A and B are said to be **independent** if $\mathbb{P}[A \cap B] = \mathbb{P}[A] \mathbb{P}[B]$. Otherwise, the events are said to be dependent.

Definition 4 (Conditional probability). Let A, B be events, and assume that $\mathbb{P}[B] > 0$. Then the **conditional probability of A , given B** , denoted $\mathbb{P}[A | B]$, is given by the formula

$$\mathbb{P}[A | B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]}.$$

Interpretation: $\mathbb{P}[A | B]$ is the probability of A when we know that event B happened.

Definition 5. We say that there exists a **positive relationship** between events A and B if

$$\mathbb{P}[A | B] > \mathbb{P}[A],$$

and a **negative relationship** if

$$\mathbb{P}[A | B] < \mathbb{P}[A].$$

Remark 6. Note the the conditions of Definition 5 are symmetric in the sense that

$$\mathbb{P}[A | B] > \mathbb{P}[A] \iff \mathbb{P}[B | A] > \mathbb{P}[B],$$

provided that both A and B have positive probability.

Example 7. Roll a 6-sided dice. Let A be the event that a ‘2’ was rolled, and B be the event that an even number was observed.

- The unconditional probability: $\mathbb{P}[A] = 1/6$.
- The conditional probability: $\mathbb{P}[A | B] = 1/3$.

Since $\frac{1}{3} > \frac{1}{6}$, we conclude there is a positive relationship between rolling a ‘2’ and rolling an even number.

End of Example 7. \square

The notion of independence formalizes the idea of “no relationship”.

Proposition 8. *If A, B are events with positive probabilities then the following are equivalent:*

- (i.) A and B are independent.
- (ii.) $\mathbb{P}[A | B] = \mathbb{P}[A]$ and $\mathbb{P}[B | A] = \mathbb{P}[B]$.

In words, independence means that the probabilities of each event are unaffected by whether or not the other event occurs. Proposition 8 simply formalizes this idea using conditional probabilities.

2.2 Random variables

Based on Sections 2.11, 4.2 in the textbook

Definition 9 (Random variable). A **random variable** (or **rv**) is a real-valued function whose domain is a sample space.

The value of a random variable is thought of as varying depending on the outcome of the experiment (the sample point). Random variables are usually denoted with capital letters, like X, Y, Z .

Example 10 (Sum 2d4). Roll a 4-sided dice twice (this is the **experiment**). There are 16 possible **outcomes**. The **sample space** is

$$S = \{(x, y) : x, y \in \{1, 2, 3, 4\}\}.$$

Let X be the sum of the two rolls. We can represent X by the following table:

		Dice 2			
		1	2	3	4
Dice 1	1	2	3	4	5
	2	3	4	5	6
	3	4	5	6	7
	4	5	6	7	8

Events are often defined using preimages of random variables. Most interesting take the form $[X \in B]$, where X is a random variable and $B \subseteq \mathbb{R}$. For example, the event that $X = 6$ is:

$$\begin{aligned} [X = 6] &= \{\omega \in S : X(\omega) = 6\} \\ &= \{(1, 5), (2, 4), (3, 3), (4, 2), (5, 1)\}. \end{aligned}$$

The textbook uses the notation $\{X = 6\}$ instead of $[X = 6]$.

Here's another example of an event. Let $E = \{2, 4, 6, 8\}$. Then

$$\begin{aligned} [X \text{ is even}] &= [X \in E] \\ &= \{\omega \in S : X(\omega) \in E\} \\ &= \{(1, 1), (1, 3), (2, 2), (2, 4), (3, 1), (3, 3), (4, 2), (4, 4)\}. \end{aligned}$$

When writing random variables, we usually suppress the arguments, e.g., writing X rather than $X(\omega)$.

End of Example 10. \square

3 2026-01-16 | Week 01 | Lecture 03

3.1 Random variables

3.1.1 Discrete vs continuous

Definition 11 (Discrete random variable). We say that a random variable X is a **discrete random variable** if it can assume only a finite or countably infinite number of distinct values.

Definition 12 (Probability mass function, pmf). Let X be a discrete random variable. The **probability mass function** (or **pmf**) of X is the function

$$p(x) = \mathbb{P}[X = x],$$

defined for every $x \in \mathbb{R}$.

Example 13. The pmf of X in Example 10 is

$$p(2) = 1/16 \quad p(3) = 2/16, \quad p(4) = 3/16, \quad p(5) = 4/16, \quad p(6) = 3/16, \quad p(7) = 2/16, \quad p(8) = 1/16$$

and $p(x) = 0$ for all other $x \in \mathbb{R}$.

End of Example 13. \square

Definition 14 (Distribution function - section 4.2). Let X be any random variable. The **cumulative distribution function** (or **cdf**) of X is the function

$$F(x) = \mathbb{P}[X \leq x],$$

defined for all $x \in \mathbb{R}$.

Remark 15. The domain of a cdf is always \mathbb{R} , and it is always a nondecreasing function with $F(-\infty) = 0$ and $F(+\infty) = 1$. The cdf of a discrete random variable is always a step function.

Definition 16 (Continuous rv). Let Y be a random variable with distribution function F . We say that Y is a **continuous random variable** if there exists a nonnegative function f such that

$$F(y) = \int_{-\infty}^y f(t)dt \tag{1}$$

for all $y \in \mathbb{R}$. The function f is called the **probability density function** (or **pdf**) of Y .

Remark 17. For continuous random variables, the distribution function F is always continuous. Moreover, for a continuous random variable Y , $\mathbb{P}[Y = b] = 0$ for all $b \in \mathbb{R}$.

Theorem 18 (Theorem 4.3 in textbook). *If Y is a continuous random variable with pdf f , then*

$$\mathbb{P}[a \leq Y \leq b] = \int_a^b f(t)dt$$

for all $-\infty \leq a \leq b \leq +\infty$.

3.1.2 Expected value

Definition 19 (Expectation of a continuous random variable). If Y is a random variable with pdf f , then the **expected value** of Y , denoted $\mathbb{E}[Y]$, is the quantity

$$\mathbb{E}[Y] = \int_{-\infty}^{\infty} yf(y)dy,$$

provided that $\int_{-\infty}^{\infty} |y|f(y)dy < \infty$.

Remark 20. $\mathbb{E}[Y]$ is the long-run average of Y , if we were to repeat the experiment many times.

The next theorem is called the *Law of the unconscious statistician (LOTUS)*.

Theorem 21 (LOTUS - single variable case). *Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be a function.*

(i.) *If X has pmf p , then*

$$\mathbb{E}[g(X)] = \sum_{\substack{x \in \mathbb{R} \\ p(x) > 0}} g(x)p(x).$$

(ii.) *If Y has pdf f , then*

$$\mathbb{E}[g(Y)] = \int_{-\infty}^{\infty} g(y)f(y)dy.$$

Remark 22. Often we wish to compute probabilities of functions of multiple random variables, for example:

- What is the probability that $\frac{X_1 + \dots + X_n}{n} \in (0, 1)$? Here, the function is $g(x_1, \dots, x_n) = \frac{x_1 + \dots + x_n}{n}$.
- What is the probability that $\max(X, Y) \leq 10$? Here the function is $g(x, y) = \max(x, y)$.
- Suppose we roll two dice and take the maximum. What is the expected value? In this case, our dice rolls are X, Y and we want to compute $\mathbb{E}[g(X, Y)]$, where $g(x, y) = \max(x, y)$.

To answer these sorts of questions, we need the notion of a “joint distribution”.

3.1.3 Joint distributions

This subsection is based on section 5.4 in the textbook. Everything in this section generalizes naturally to n variables, but the results are simpler to state for just 2 random variables.

Definition 23 (Joint pmf). Let X_1 and X_2 be discrete random variables. The **joint probability mass function** for X_1 and X_2 is the function

$$p(x_1, x_2) = \mathbb{P}[X_1 = x_1, X_2 = x_2],$$

defined for all $x_1, x_2 \in \mathbb{R}$.

Definition 24 (Joint pdf). Let Y_1 and Y_2 be continuous random variables. We say that Y_1 and Y_2 are **jointly continuous** if there exists a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}_{\geq 0}$ such that

$$\mathbb{P}[Y_1 \leq y_1, Y_2 \leq y_2] = \int_{-\infty}^{y_1} \int_{-\infty}^{y_2} f(t_1, t_2) dt_2 dt_1.$$

for all $y_1, y_2 \in \mathbb{R}$. The function f is called the **joint probability density function** for Y_1 and Y_2 .

Theorem 25 (LOTUS - multivariable case). *Let $g : \mathbb{R}^2 \rightarrow \mathbb{R}$.*

- *If X_1, X_2 have joint pmf $p(x_1, x_2)$, then*

$$\mathbb{E}[g(X_1, X_2)] = \sum_{\substack{(x_1, x_2) \in \mathbb{R}^2 \\ p(x_1, x_2) > 0}} g(x_1, x_2)p(x_1, x_2).$$

- *If Y_1, Y_2 are jointly continuous random variables with joint pdf $f(y_1, y_2)$, then*

$$\mathbb{E}[g(Y_1, Y_2)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(y_1, y_2)f(y_1, y_2)dy_1dy_2.$$

Remark 26. Theorem 25 generalizes to n variables. It gives us a way to answer questions like the third question posed in Remark 22.

4 2026-01-21 | Week 02 | Lecture 04

Question for this lecture: what does it mean for random variables to be independent, and what does it buy us?

4.1 Independence of random variables

4.1.1 Definition and characterization

The aim of this section is to define what it means for random variables to be independent.

Definition 27. We say that the random variables X_1, X_2, \dots, X_n are **independent** if the following holds for all possible values x_1, \dots, x_n in their range:

$$\mathbb{P}[X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n] = \mathbb{P}[X_1 \leq x_1] \mathbb{P}[X_2 \leq x_2] \cdots \mathbb{P}[X_n \leq x_n].$$

Theorem 28 (Factorization theorem – Theorem 5.4 in textbook). *For discrete/continuous random variables, independence is equivalent to factorizability of the joint pmf/pdf. More formally, we have:*

- **Discrete case:** Let X_1, \dots, X_n be discrete random variables. Then X_1, \dots, X_n are independent if and only if

$$\mathbb{P}[X_1 = x_1, X_2 = x_2, \dots, X_n = x_n] = \mathbb{P}[X_1 = x_1] \mathbb{P}[X_2 = x_2] \cdots \mathbb{P}[X_n = x_n]$$

for all $x_1, \dots, x_n \in \mathbb{R}$.

- **Continuous case:** Let Y_1, \dots, Y_n be continuous random variables with pdfs f_1, \dots, f_n , and joint pdf $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Then Y_1, \dots, Y_n are independent if and only if

$$f(y_1, \dots, y_n) = f_1(y_1) f_2(y_2) \cdots f_n(y_n)$$

for all $y_1, \dots, y_n \in \mathbb{R}$.

Remark 29. In this course, we will frequently work with “samples” of n independent random variables X_1, \dots, X_n . The intuition about independence for n variables is that observing any number of them doesn’t give you any information about the others.

4.1.2 Independence and expectations/variances

Independence splits the expectation of a product into a product of expectations.

Theorem 30. If X_1, X_2, \dots, X_n , are independent, then

$$\mathbb{E}[X_1 X_2 \cdots X_n] = \mathbb{E}[X_1] \mathbb{E}[X_2] \cdots \mathbb{E}[X_n],$$

provided that the expectations exist.

Proof. This can be proven for continuous/discrete random variables by an application of Theorem 25 and Theorem 28. \square

Definition 31 (Variance of a random variable). Let X be a random variable. The **variance** of X , denoted $\text{Var}(X)$, is the quantity

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

where $\mu = \mathbb{E}[X]$. The positive square root of the variance is the **standard deviation** of X .

Remark 32. Variance is often denoted by σ^2 . The textbook uses the notation $V(X)$ instead of $\text{Var}(X)$.

Theorem 33 (Properties of Variance). *Let X be a random variable and a, b be scalars. Then*

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

If X and Y are independent random variables then

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y).$$

Proof. Follows by direct computation using the definition of variance. □

The second part of Theorem 33 says that for independent random variables, variance is additive.

4.1.3 Independence and long-term behavior: the weak law of large numbers

Combining many independent sources of randomness tends to produce predictable phenomena. Here we introduce one example of that.

Theorem 34 (Weak Law of Large Numbers). *Let $(X_n)_{n=1}^\infty$ be a sequence of independent random variables. Assume that for all n , $\mu = \mathbb{E}[X_n]$ and $\text{Var}(X_n) \leq B$ for some fixed bound $B < \infty$. Let $S_n = X_1 + \dots + X_n$. Then for all $\epsilon > 0$,*

$$\mathbb{P} \left[\left| \frac{S_n}{n} - \mu \right| > \epsilon \right] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

5 2026-01-23 | Week 02 | Lecture 05

5.1 Proof of the weak law of large numbers

We will need the following theorem, which says that if the average male height is 5 feet tall, then no more than 10% of men can be more than 50 feet tall.

Theorem 35 (Markov's inequality). *Let X be a nonnegative random variable and let $a > 0$. Then*

$$\mathbb{P}[X \geq a] \leq \frac{\mathbb{E}[X]}{a} \quad (2)$$

Proof. Let $A = [X \geq a]$. Then $X \cdot \mathbf{1}_A \leq X$, so

$$\begin{aligned} \mathbb{E}[X] &\geq \mathbb{E}[X \cdot \mathbf{1}_A] \\ &\geq \mathbb{E}[a \mathbf{1}_A] \\ &\geq a \mathbb{E}[\mathbf{1}_A] \\ &= a \mathbb{P}[A] \\ &= a \mathbb{P}[X \geq a]. \end{aligned}$$

Dividing by a implies Eq. (2). □

Theorem 36 (Weak Law of Large Numbers). *Let $(X_n)_{n=1}^\infty$ be a sequence of independent random variables. Assume that for all n , $\mu = \mathbb{E}[X_n]$ and $\text{Var}(X_n) \leq B$ for some fixed bound $B < \infty$. Let $S_n = X_1 + \dots + X_n$. Then for all $\epsilon > 0$,*

$$\mathbb{P}\left[\left|\frac{S_n}{n} - \mu\right| > \epsilon\right] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Proof. First observe that by independence,

$$\text{Var}(S_n) = \text{Var}(X_1) + \dots + \text{Var}(X_n) \leq nB. \quad (3)$$

Next, let $\epsilon > 0$ be arbitrary.

$$\begin{aligned} \mathbb{P}\left[\left|\frac{S_n}{n} - \mu\right| > \epsilon\right] &= \mathbb{P}[|S_n - n\mu| > n\epsilon] \\ &= \mathbb{P}[(S_n - n\mu)^2 > (n\epsilon)^2] \\ &\leq \frac{\mathbb{E}[(S_n - n\mu)^2]}{n^2\epsilon^2} && \text{by Theorem 35} \\ &= \frac{\text{Var}(S_n)}{n^2\epsilon^2} && \text{by definition of variance} \\ &\leq \frac{B}{\epsilon^2 n} && \text{by Eq. (3).} \end{aligned}$$

The right hand side tends to zero as $n \rightarrow \infty$, proving the theorem. □

5.2 Equality in distribution

Definition 37 (Identically distributed). Two random variables X and Y are *identically distributed* if they have the same cdf.

Remark 38. Identically distributed random variables are said to “have the same distribution”. For discrete/continuous random variables, this is equivalent to them having the same pmf/pdf.

The next example shows that identically distributed random variables need not be equal as functions.

Example 39. Flip a coin. Define two random variables

$$X = \begin{cases} 1 & : \text{ coin is heads} \\ 0 & : \text{ coin is tails} \end{cases}$$

$$Y = \begin{cases} 0 & : \text{ coin is heads} \\ 1 & : \text{ coin is tails} \end{cases}$$

Then X and Y are identically distributed, since they have the same pmf:

$$\mathbb{P}[X = 1] = \mathbb{P}[Y = 1] = \frac{1}{2} \quad \text{and} \quad \mathbb{P}[X = 0] = \mathbb{P}[Y = 0] = \frac{1}{2}.$$

But of course X and Y are never equal, i.e., $\mathbb{P}[X \neq Y] = 1$.

End of Example 39. \square

5.3 Populations and samples

Here we introduce a conceptual framework for mathematical statistics.

A **population** is a large body of data that is the target of our interest. The subset collected from it is our **sample**.

Example 40 (Populations). A population can be real or theoretical. Here are some examples

- The set of people in Hawaii (real, finite)
- The set of voters in the 2026 Midterm elections (hypothetical, finite)
- The decimal expansion of π (countably infinite)
- The infinitely many observations that could be made during a laboratory experiment if the experiment were repeated over and over again (hypothetical)
- The lifetimes of light bulbs produced by a factory

Importantly, a population can also be a *probability distribution*, specified by a pdf, pmf, or cdf

- Observations made from an exponential distribution with mean $\lambda > 0$ (i.e., the distribution with pdf $f(x) = \lambda e^{-\lambda x} \mathbf{1}_{[x>0]}$.)

End of Example 40. \square

6 2026-01-26 | Week 03 | Lecture 06

6.1 Sampling

For discussions of “random sample”, see sections 2.12 and 6.1 in the textbook.

The simplest sampling procedure is called simple random sampling.

Definition 41 (Simple random sampling). Let N and n denote the numbers of elements in the population and sample, respectively. If the sampling is conducted in such a way that each of the $\binom{N}{n}$ samples has an equal probability of being selected, the sampling is called **simple random sampling**, and the result is a **simple random sample**.

More commonly in mathematical statistics, we think of the population as a distribution.

Definition 42 (Random sample from a distribution). Consider a given probability distribution on \mathbb{R} that can be represented by a pdf or pmf f . We say that the random variables X_1, \dots, X_n form a **random sample** from this distribution if these random variables are independent and the distribution of each is given by f . Such random variables are also said to be **independent and identically distributed** (iid). The number of random variables n is the **sample size**.

The objective of statistics is to make an inference about a population based on information contained in a sample from that population and to provide an associated measure of goodness for the inference.

6.2 What is a statistic?

Section 7.1 in the textbook.

Definition 43 (Statistic). A **statistic** is a function of the observable random variables in a sample and known constants.

In other words, if X_1, \dots, X_n is a random sample and $T : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function, then the random variable

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

is a statistic.² The probability distribution of such a statistic Y is called its **sampling distribution**.

Often, we have some quantity of interest, called a **target parameter**, and we want a single “best guess” of some quantity of interest. When this is the case, then we call a statistic an estimator:

Definition 44 (Estimator). An **estimator** is a statistic, that is a function $T(X_1, \dots, X_n)$ of a sample, that is used to approximate a target parameter. An **estimate** is the realized value of an estimator (e.g., a number) that is obtained when the sample is actually taken.

In words, an estimator is a rule, often expressed as a formula, that tells us how to calculate the value of an estimate based on the measurements contained in a sample. Here are two important examples of estimators.

Definition 45 (Sample Mean, Sample Variance). Let X_1, \dots, X_n be a random sample. The **sample mean**, denoted \bar{X} , is the statistic

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

The **sample variance**, denoted S^2 , is the statistic

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

The target parameters these are used to estimate are the mean and variance of the population. The following examples illustrate some of the terminology we’ve introduced.

²A statistic does not need to be real-valued, but we will focus on the case when it is.

Example 46 (Estimating average lightbulb lifetime). Suppose the lifetime of each light bulb produced in a certain factory is distributed according to the following pdf

$$f(y) = \lambda e^{-\lambda y} \mathbf{1}_{[y>0]}.$$

for some parameter $\lambda > 0$, where y is measured in minutes. Here, the **population** is either the set of lightbulbs produced by the factory, which we idealize as the probability distribution given by f . Suppose we do not know λ , and we wish to estimate it. Hence λ is the **target parameter**. A random sample Y_1, \dots, Y_n is taken. (So the Y_i 's are independent and each as pdf f .)

In this case f is an exponential distribution with rate $\lambda > 0$. The exponential distribution is used to modeling lifetimes of systems that are not subject to degradation (e.g., lifetimes of light bulbs, not diesel engines) as well as waiting times (e.g., the length of time intervals between cars passing by a tree on an idyllic country road). In the road example, if the rate of cars passing is 2 cars per minute, then on average you have to wait only 1/2 of a minute for a car to pass. That is, the *rate* and *expected values* are reciprocals of each other. Let's prove this fact in the following claim:

Claim 1: If Y has pdf f , then $\mathbb{E}[Y] = \frac{1}{\lambda}$.

Proof of Claim 1. By definition of expected value for a continuous random variable,

$$\begin{aligned} \mathbb{E}[Y] &= \int_{-\infty}^{\infty} y \lambda e^{-\lambda y} \mathbf{1}_{[y>0]} dy \\ &= \int_0^{\infty} y \lambda e^{-\lambda y} dy \\ &= \int_0^{\infty} y \frac{d}{dy} [-e^{-\lambda y}] dy \\ &= [-ye^{-\lambda y}]_{y=0}^{y=\infty} - \int_0^{\infty} \frac{d}{dy} [y] (-e^{-\lambda y}) dy && \text{by integration by parts} \\ &= \int_0^{\infty} e^{-\lambda y} dy \\ &= \left[\frac{1}{\lambda} e^{-\lambda y} \right]_{y=\infty}^{y=0} \\ &= \frac{1}{\lambda}. \end{aligned}$$

□ Claim

End of Example 46. □

7 2026-01-28 | Week 03 | Lecture 07

7.1 Basic examples of statistics and estimators

This lecture is based on section 7.1 in the textbook.

Example 47 (Example 46 continued). We have a random sample of Y_1, \dots, Y_n from the distribution with pdf

$$f(y) = \lambda e^{-\lambda y} \mathbf{1}_{[y>0]},$$

where λ is a target parameter.

In the setting of this example, we interpret the random variables Y_1, \dots, Y_n as lifetimes (in, say, years) of lightbulbs produced by a certain factory.

In the last lecture, we used integration by parts to show that if Y has pdf f , then

$$\mathbb{E}[Y] = \frac{1}{\lambda}.$$

Consider the **statistic**

$$\bar{Y} = \frac{Y_1 + \dots + Y_n}{n}.$$

By Theorem 36 (The Weak Law of Large Numbers), \bar{Y} is likely to be close to $\mathbb{E}[Y] = 1/\lambda$ when n is large. Therefore $\frac{1}{\bar{Y}}$ is likely to be close to λ . Therefore a reasonable **estimator** for λ is

$$\frac{1}{\bar{Y}} = \frac{n}{Y_1 + \dots + Y_n}.$$

If $n = 10$ and the observed values of Y_1, \dots, Y_{10} are

$$2, 0.8, 0.1, 2.3, 3.2, 5.4, 2, 0.4, 2.8, \text{ and } 3.4$$

then $\bar{Y} = 2.24$, and hence our **estimate** for λ would be $\frac{1}{2.24} \approx 0.45$.

End of Example 47. \square

Example 48 (Example 7.1 in textbook). Roll a dice three times. Let Y_1, Y_2, Y_3 be the values of the rolls. The average number observed in this sample of size 3 is

$$\bar{Y} = \frac{Y_1 + Y_2 + Y_3}{3}.$$

This is the **sample mean**, it is function of the sample (and the known constant n) and is therefore a **statistic**.

Question: What the mean $\mu_{\bar{Y}}$ and standard deviation $\sigma_{\bar{Y}}$ of the random variable \bar{Y} ?

First we compute the mean and variance of a single Y_i .

Claim 1: $\mathbb{E}[Y_i] = 3.5$ and $\text{Var}(Y_i) = 35/12 \approx 2.9167$.

Proof of Claim 1. Fix $i \in [3]$. For simplicity, let $Y = Y_i$. Using the definition of expected value,

$$\mathbb{E}[Y] = \sum_{k=1}^6 k \mathbb{P}[Y = k] = \sum_{k=1}^6 k \frac{1}{6} = 3.5.$$

Next we compute the variance of Y :

$$\begin{aligned} \text{Var}(Y) &= \mathbb{E}[(Y - 3.5)^2] && \text{Definition 31} \\ &= \sum_{y=1}^6 (y - 3.5)^2 \mathbb{P}[Y = y] && \text{By Theorem 21 (LOTUS) using } g(y) = (y - 3.5)^2 \\ &= \sum_{k=1}^6 (y - 3.5)^2 \frac{1}{6} && \text{since } \mathbb{P}[Y = y] = 1/6 \text{ for all } k \in [6] \\ &= 35/12. \end{aligned}$$

□ Claim

Claim 2: $\mu_{\bar{Y}} = \mathbb{E}[\bar{Y}] = 3.5$

Proof of Claim 2. We have:

$$\begin{aligned}
 \mathbb{E}[\bar{Y}] &= \mathbb{E}\left[\frac{Y_1 + Y_2 + Y_3}{3}\right] && \text{by definition of } \bar{Y} \\
 &= \frac{1}{3}(\mathbb{E}[Y_1] + \mathbb{E}[Y_2] + \mathbb{E}[Y_3]) && \text{by linearity of expectation} \\
 &= \frac{1}{3}(3.5 + 3.5 + 3.5) && \text{by Claim 1} \\
 &= 3.5.
 \end{aligned}$$

□ Claim

Claim 3: $\sigma_{\bar{Y}}^2 = \text{Var}(\bar{Y}) = .9722$

Proof of Claim 3.

$$\begin{aligned}
 \text{Var}(\bar{Y}) &= \text{Var}\left(\frac{Y_1 + Y_2 + Y_3}{3}\right) \\
 &= \frac{1}{9}\text{Var}(Y_1 + Y_2 + Y_3) && \text{by Theorem 33} \\
 &= \frac{1}{9}(\text{Var}(Y_1) + \text{Var}(Y_2) + \text{Var}(Y_3)) && \text{by Theorem 33} \\
 &= \frac{1}{9}(2.9167 + 2.9167 + 2.9167) && \text{by Claim 1} \\
 &= \frac{2.9167}{3} && (*)
 \end{aligned}$$

and this is approximately .97224.

□ Claim

Since the standard deviation is the square root of the variance, the standard deviation is

$$\sigma_{\bar{Y}} = \sqrt{.9722} = .9860$$

Claim 3.

Remark 49. In equation (*), we had $\text{Var}(\bar{Y}) = \frac{2.9167}{3}$. Notice how, if we computed the variance of four or five dice rolls, we'd get

$$\text{Var}(\bar{Y}) = \frac{2.9167}{4} \quad \text{or} \quad \text{Var}(\bar{Y}) = \frac{2.9167}{5}.$$

In general, for n dice rolls, we'd get

$$\text{Var}(\bar{Y}) = \frac{2.9167}{n}.$$

This quantity tends to zero as $n \rightarrow \infty$. In words, as the number of samples n grows, the variance of the sample mean decreases.

End of Example 49. □

8 2026-01-30 | Week 03 | Lecture 08

8.1 The Gaussian distribution

Begin section 7.2

Definition 50 (Normal / Gaussian Distribution). A random variable Y is said to be **normally distributed** with mean $\mu \in \mathbb{R}$ and variance $\sigma^2 > 0$ if it has pdf

$$f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y - \mu)^2}{2\sigma^2} \right\},$$

defined for all $y \in \mathbb{R}$. We abbreviate this by writing $Y \sim \mathcal{N}(\mu, \sigma^2)$. If $Y \sim \mathcal{N}(0, 1)$, we say Y has **standard normal distribution**.

Remark 51. Normal random variables/distributions are also called **Gaussian** random variables/distributions.

Example 52. Normal distributions are the classic “bell curve” distributions. Whenever a random quantity is the result of adding up many small, independent random effects, we expect it to be approximately normally distributed.

- A classic example is human height, which is influenced by the combined effects of thousands of genes together with environmental factors.
- Another example is repeated measurements of the same physical quantity (e.g., if we all went out and measured the temperature outside tomorrow at midday), this is because the total error is the sum of many small independent sources of noise (e.g., instrument precision, time, location, shade, wind, etc).

End of Example 52. \square

Remark 53 (Standardization). If $Y \sim \mathcal{N}(\mu, \sigma^2)$ then the random variable

$$Z = \frac{X - \mu}{\sigma}$$

is a standard normal random variable. The transformation

$$X \mapsto \frac{X - \mu}{\sigma}$$

in which we subtract off the mean of X and then divide by its standard deviation is called **standardization**.

Remark 54 (The Empirical Rule). The **Empirical Rule** (also called the **68-95-99.7 rule**; see textbook section 1.3) is summarized by the following figure, for normal distributions:

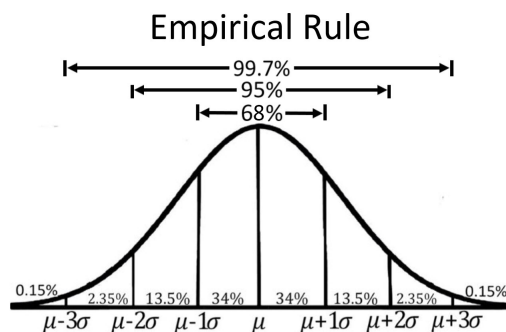


Image source: <https://andymath.com/wp-content/uploads/2019/12/empirical-rule-normdist2.jpg>

The next theorem says that if we have a sample of Y_1, \dots, Y_n normal random variables, then their sample mean is normally distributed as well.

Theorem 55. Let Y_1, \dots, Y_n be a random sample of size n from a normal distribution with mean μ and variance σ^2 . Then

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

is normally distributed with mean $\mu_{\bar{Y}} = \mu$ and variance $\sigma_{\bar{Y}}^2 = \sigma^2/n$.

Proof sketch. The proof follows from three facts:

- Any linear combination of normal random variables is a normal random variable (this can be proved using moment generating functions).
- That $\mathbb{E}[\bar{Y}] = \mu$ follows by linearity of expectation.
- The variance calculation follows by Theorem 33.

□

Remark 56. Applying the idea of standardization, observe that the random variable $Z = \frac{\bar{Y} - \mu_{\bar{Y}}}{\sigma_{\bar{Y}}} = \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}}$ has a standard normal distribution.

Example 57 (Example 7.2 in textbook). A bottling machine discharges an average of μ ounces per bottle. It has been observed that the amount of liquid discharged is normally distributed with standard deviation $\sigma = 1$ per ounce. A sample of $n = 9$ filled bottles is randomly selected, and is measured.

Question: Find the probability that the sample mean will be within .3 ounces of the true mean μ .

Solution: Let Y_1, \dots, Y_9 be the ounces in each of the bottles. Then $Y_i \sim \mathcal{N}(\mu, 1)$ for $i \in [9]$. By Theorem 55, $\bar{Y} \sim \mathcal{N}(\mu, \sigma^2/n)$, and $\sigma^2/n = 1/9$. We want to find

$$\begin{aligned} \mathbb{P}[|\bar{Y} - \mu| \leq .3] &= \mathbb{P}[-.3 \leq \bar{Y} - \mu \leq .3] \\ &= \mathbb{P}\left[-\frac{.3}{\sigma/\sqrt{n}} \leq \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \leq \frac{.3}{\sigma/\sqrt{n}}\right]. \end{aligned} \quad (4)$$

By Remark 56, $Z = \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1)$, and plugging in $\sigma = 1$ and $n = 9$ and simplifying gives

$$\begin{aligned} \mathbb{P}[|\bar{Y} - \mu| \leq .3] &= \mathbb{P}[-.9 \leq Z \leq .9] \\ &= 1 - 2\mathbb{P}[Z > .9] \\ &= .6318. \end{aligned}$$

End of Example 57. □

Example 58 (Example 7.3 in textbook). How many observations need to be included in the sample to ensure that \bar{Y} is within .3 ounces of μ with probability .95?

Now we want

$$\mathbb{P}[|\bar{Y} - \mu| \leq .3] = .95.$$

By Eq. (4),

$$\mathbb{P}[|\bar{Y} - \mu| \leq .3] = \mathbb{P}\left[-\frac{.3}{\sigma/\sqrt{n}} \leq \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \leq \frac{.3}{\sigma/\sqrt{n}}\right].$$

Recalling that $\sigma = 1$ and $Z = \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1)$, we have

$$\mathbb{P}[|\bar{Y} - \mu| \leq .3] = \mathbb{P} [.3\sqrt{n} \leq Z \leq .3\sqrt{n}].$$

So we need to find n such that $\mathbb{P} [.3\sqrt{n} \leq Z \leq .3\sqrt{n}] = .95$.

It is a well-known fact that

$$\mathbb{P}[-1.96 \leq Z \leq 1.96] = .95.$$

So we need $.3\sqrt{n} = 1.96$, or equivalently, $n = \left(\frac{1.96}{.3}\right)^2 = 42.68$.

So we need at least 43 samples, since 42 is not quite enough.

End of Example 58. □

9 2026-02-02 | Week 04 | Lecture 09

9.1 Chi-squared distributions

Chapter 7.2

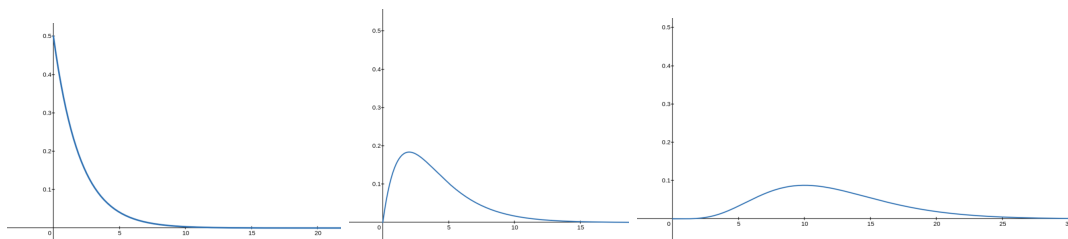
A special kind of gamma distribution:

Definition 59 (Chi-Squared Distribution). Let k be a positive integer. A nonnegative random variable X is said to have **chi-squared distribution with parameter d** if it has pdf

$$f(x) = \begin{cases} \frac{1}{2^{d/2}\Gamma(d/2)} x^{\frac{d}{2}-1} e^{-x/2} & : x \geq 0 \\ 0 & : x < 0 \end{cases}$$

In other words, if X is a Gamma distributed random variable with parameters $\alpha = d/2$ and $\beta = 2$. The parameter d is usually called the **degrees of freedom** of X .

From left-to-right, here's what this looks like for $d = 2$, $d = 4$, and $d = 12$:



Using properties of gamma distribution, we have:

Proposition 60. A χ^2 random variable X with d degrees of freedom has

$$\mathbb{E}[X] = \alpha\beta = d \quad \text{and} \quad \text{Var}(X) = \alpha\beta^2 = 2d$$

Theorem 61 (Chi-squared is sum of squared normals). Let $n \geq 1$. If Z_1, \dots, Z_n are independent standard normal random variables, then the random variable $V = Z_1^2 + \dots + Z_n^2$ has chi-squared distribution with n degrees of freedom.

Proof. Proved using moment generating functions. (See Example 6.11 in the textbook). \square

Remark 62. Taking $n = 1$, Theorem 61 implies that $Z \sim \mathcal{N}(0, 1)$ implies Z^2 is chi-squared with 1 degree of freedom.