

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Spring 2025



HAVERFORD
COLLEGE

- **Lab 5 grades** on Moodle
- **Lab 7 / Project proposal** due last night
- **Lab 8 posted**, due ~~Wednesday~~ Thursday
- **Final project** instructions/rubric posted
- Next week:
 - Finish statistics, then midterm review
 - **Midterm in-class Thursday April 17**

Outline for today

- Randomized trials for the null distribution
- Are the means of two samples different?
 - t-tests
 - Permutation testing
- Bootstrapping

Outline for today

- Randomized trials for the null distribution
- Are the means of two samples different?
 - t-tests
 - Permutation testing
- Bootstrapping

Central Limit Theorem

- Assumptions

- X_1, X_2, \dots, X_n are iid samples
- From a population with mean μ
- Finite variance σ^2

- THEN

$$Z = \lim_{n \rightarrow \infty} \sqrt{n} \left(\frac{\bar{X}_n - \mu}{\sigma} \right)$$

is a standard normal distribution (i.e. mean 0 and variance 1)

Central Limit Theorem

- Last time we saw that the central limit theorem could be used to estimate a p-value

$$Z = \lim_{n \rightarrow \infty} \sqrt{n} \left(\frac{\bar{X}_n - \mu}{\sigma} \right)$$

- We first obtain a Z-score, then compute the probability of observing a result *as or more* extreme **under the null hypothesis**

Recap Handout 17

+

Continuous \Leftrightarrow Discrete features

Handout 17

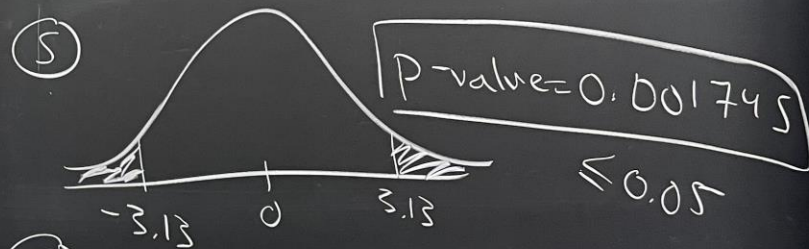
$$\textcircled{1} E[X] = \sum_x x \cdot p(x) = 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = \boxed{\frac{1}{2} = \mu}$$

$$\begin{aligned}\textcircled{2} \text{Var}(X) &= E[(X - \mu)^2] \\ &= \sum_x (x - \mu)^2 p(x) \\ &= \left(0 - \frac{1}{2}\right)^2 \cdot \frac{1}{2} + \left(1 - \frac{1}{2}\right)^2 \cdot \frac{1}{2} \\ &= \boxed{\frac{1}{4} = \sigma^2}\end{aligned}$$

$$\textcircled{3} \bar{X}_n = \frac{54}{80} = 0.675$$

$$\textcircled{4} z = \frac{\bar{X}_n - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{0.675 - 0.5}{\sqrt{\frac{0.25}{80}}}$$

$$\boxed{z \approx 3.13}$$



$\textcircled{6}$ reject $H_0!$

Continuous \rightarrow discrete

Lab 6

- Sort & find label break points

- ex age \Rightarrow $\text{age} \geq 50$
 $\text{age} < 50$

discrete \rightarrow continuous

- binary $\rightarrow 0 \text{ \& } 1$

- ordered $\rightarrow 0, 1, 2, \dots, V-1$

- non-ordered

\Rightarrow indicators

is apple?	is pear?	is orange?
0, 1	0, 1	0, 1

Better way? Randomized trials

- Die example
 - $n=10$ rolls
 - $[4, 2, 3, 1, 3, 1, 3, 3, 3, 1]$
 - $\bar{X}_n = 2.4$
- H_0 : null hypothesis (fair die)
 - What if we don't know mean & variance of null distribution?
- H_1 : is the die weighted toward lower values? (one-sided)

Randomized trials: general idea

1. Run T trials that *mimic* our data under the null hypothesis
roll a fair die
2. Record relevant information for each trial
mean of the rolls
3. Count how many times you observe a result *as or more extreme* than your data (N_e)
any trial with mean less than or equal to 2.4
4. $p\text{-value} = N_e/T$

Randomized trials: general idea

1. Run T trials that *mimic* our data under the null hypothesis

roll a fair die

Right now: each person does 1 trial!

2. Record relevant information for each trial

mean of the rolls

3. Count how many times you observe a result *as or more extreme* than your data (N_e)

any trial with mean less than or equal to 2.4

4. $p\text{-value} = N_e/T$

4.3 4.1 4.1

$$T = 18$$

3.8 4.1 3.2

4.7 3.3

$$N_e = 0$$

3.8 3.1

$$p\text{-value} = \frac{N_e}{T} = 0$$

3.2 3.9

3.8 3.1

\Rightarrow reject null

4.0 4.0

hypothesis of a
fair die

3.4 3.4

Handout 18

Handout 18

① $T = 20$

② $N_e = 8$

$N_e = 12$

one-sided

13 or fewer

two-sided

13 or fewer

or 17 or more

③ $p\text{-value} = \frac{12}{20} = 0.6 \gg 0.05$

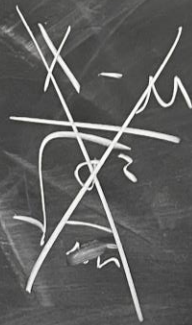
fail to reject H_0

④



$p\text{-value} = \frac{\overbrace{\text{shaded area}}^{Ne}}{\underbrace{\text{shaded area} + \text{unshaded area}}_T}$

⑤



don't know μ & σ^2

Outline for today

- Randomized trials for the null distribution
- Are the means of two samples different?
 - t-tests
 - Permutation testing
- Bootstrapping

Difference in means

(blood pressure)

example:

before drug: [117, 54, 96, 123, 157, ...]

n examples

$$\bar{X}_n = 112$$

m examples

$$\bar{X}_m = 96$$

H_0 : all #'s are drawn from the same distribution

H_1 : after the drug, blood pressure was lowered
(one-sided)

plot

permutation testing

- Simulate null distribution!
- permute the "labels" of the data (i.e. "before" & "after")

- for t in range(T):
permute labels \leftarrow helper

plot : $\underbrace{\bar{X}_m^{(t)} - \bar{X}_n^{(t)}}_{\text{difference in means}}$

1 trial

"before" : [98, 123, 105, 54 ...]

"after" : [82, 72, 117, 157, 96 ...]

as or more
extreme than
-16



actual obs

$$\bar{X}_m - \bar{X}_n = 96 -$$

$$\text{continuous} = \boxed{-16}$$

$$\bar{X}_m - \bar{X}_n$$

Say: $T=1000$, $N_e=4 \Rightarrow \boxed{p=0.004} < 0.05$

54 ...]

still (n)

$$\boxed{\bar{X}_n^{(1)} = 101}$$

57, 96 ...]

still (m)

$$\boxed{\bar{X}_m^{(1)} = 105}$$

actual obs

$$\bar{X}_m - \bar{X}_n = 96 - 112$$

$$\text{continuum} = \boxed{-16}$$

$$\bar{X}_m - \bar{X}_n$$

$$\boxed{p = 0.004} < 0.05$$

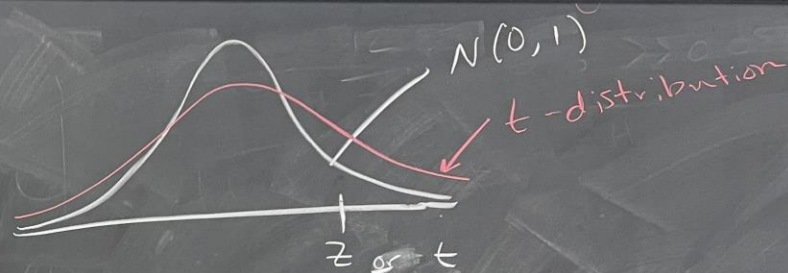
t-tests CLT-inspired test

CLT: $Z = \frac{\bar{X}_n - \mu}{\sqrt{\frac{\sigma^2}{n}}} \sim N(0, 1)$ ↖ drawn from

don't know σ ? \Rightarrow use sample variance

$$t = \frac{\bar{X}_n - \mu}{\sqrt{\frac{s^2}{n}}} \sim t\text{-distribution}$$

(s²)



Sample variance $\Rightarrow S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$

↑
Sample mean

difference in means

2 fields A & B

	A	B
\bar{X}_n	1.3	1.6
n	22	24
s	0.5	0.3

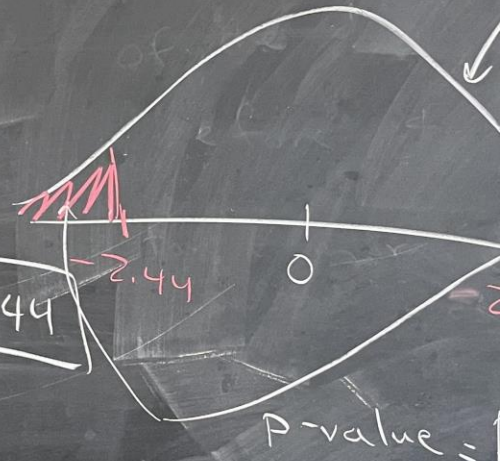
sample
std dev

t-test

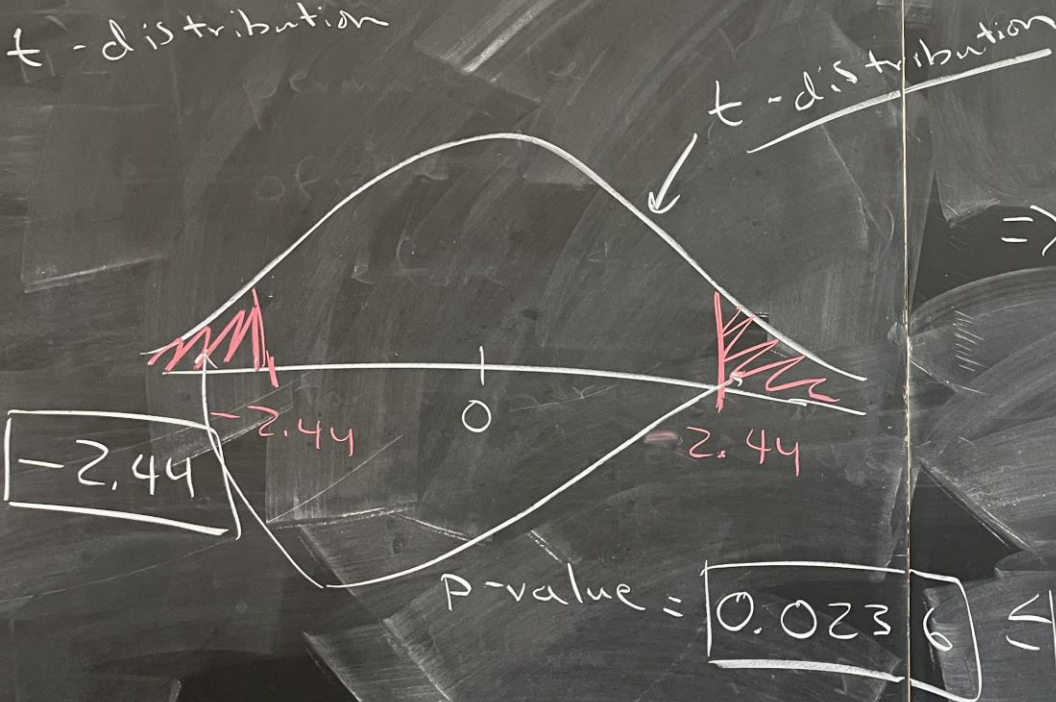
$$t = \frac{\bar{X}_A - \bar{X}_B}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}}$$

$\sim t$ -distribution

$$t = \frac{1.3 - 1.6}{\sqrt{\frac{0.25}{22} + \frac{0.09}{24}}} = -2.44$$



t-distribution



⇒ reject null!

$$H_0: \mu_A = \mu_B$$

$$H_1: \mu_A \neq \mu_B \quad (\text{2-sided})$$

Outline for today

- Randomized trials for the null distribution
- Are the means of two samples different?
 - t-tests
 - Permutation testing
- **Bootstrapping**

Next time!