

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023



HAVERFORD
COLLEGE

- **Lab 7 + project proposal** due Wednesday
 - During lab this week: project meetings with all groups
 - Try to come to the same lab section as your partner
 - We can also discuss Lab 7

Languages Left Behind:

Why Language Models Fail at Non-English
Content Moderation

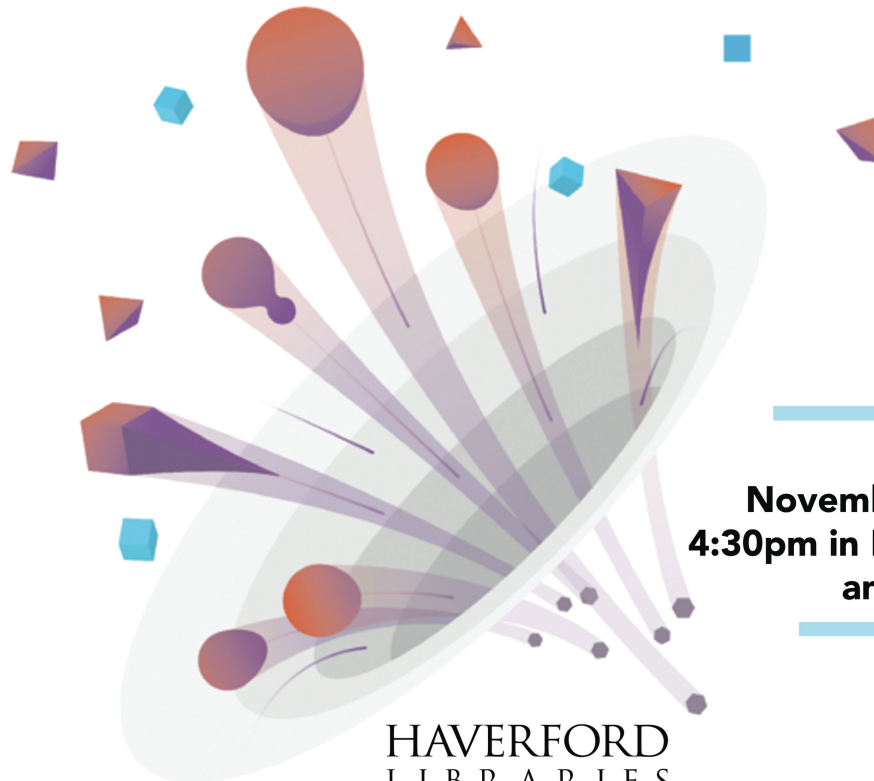


Aliya Bhatia
Policy Analyst
Free Expression Project



Gabriel Nicholas
Research Fellow

Center for Democracy & Technology



Thursday
November 9, 2023
4:30pm in Lutnick 200
and on Zoom

HAVERFORD
LIBRARIES

Outline for November 7

- Recap PCA and Handout 16
- Begin: statistics and hypothesis testing
- Central limit theorem
- Bootstrapping

Outline for November 7

- Recap PCA and Handout 16
- Begin: statistics and hypothesis testing
- Central limit theorem
- Bootstrapping

Handout 16

$$A = \begin{bmatrix} 3/10 & -3/10 \\ -3/10 & 3/10 \end{bmatrix}$$

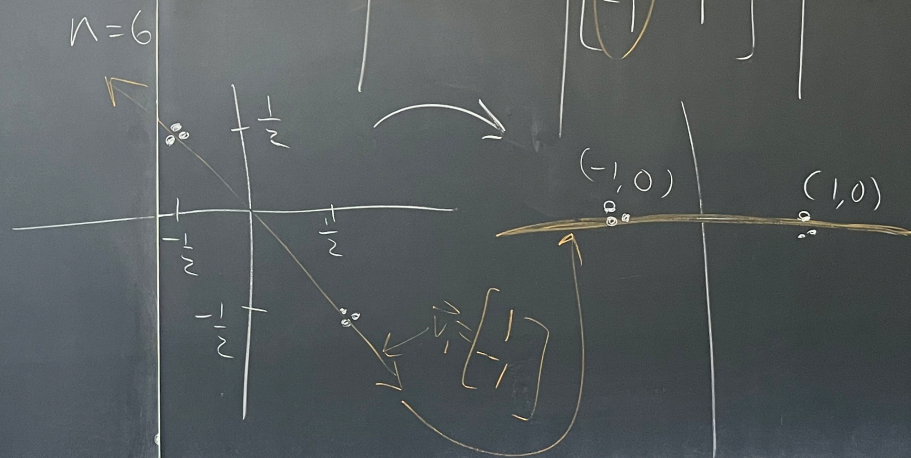
$$\vec{V}_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\lambda_1 = \frac{2}{5}$$

$$\vec{V}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\lambda_2 = 0$$

$$\begin{aligned} T &= XW \\ n \times 2 &= \begin{matrix} n \times p & p \times 2 \\ f_1 & f_2 \end{matrix} \\ &= \begin{bmatrix} -\frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} \vec{V}_1 \\ \vec{V}_2 \end{bmatrix} = \begin{bmatrix} 1 \\ ? \end{bmatrix} \end{aligned}$$



Outline for November 7

- Recap PCA and Handout 16
- **Begin: statistics and hypothesis testing**
- Central limit theorem
- Bootstrapping

Motivation for studying statistics

- 1) I have a new method that achieves 95% accuracy on a dataset. The previous best method achieved 92% accuracy. Is my method significantly better?
- 2) I have created a new treatment for high blood pressure. Did it significantly lower the blood pressure of the treatment group over the control group?
- 3) Which variants in the genome are statistically correlated with a specific disease?

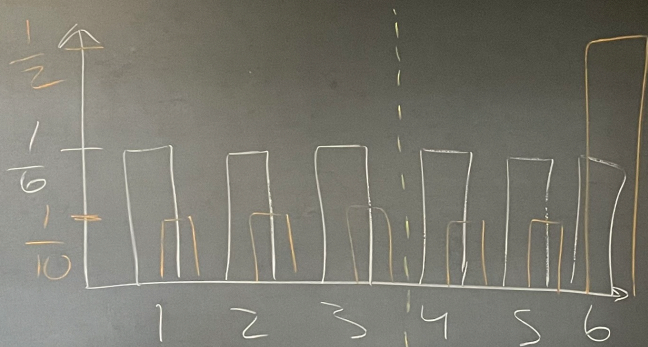
Motivation for studying statistics

- In general there are many questions that can only be answered properly with statistics.
- This one week on statistics is not a substitute for a full stats course, which I recommend everyone take!
- We're going to do a few key examples now to build up intuition, but this is a huge field and not my main area of research.

Distributions, expected value, variance

fair die

weighted die



Probability mass functions (pmf)

$$\sum_{x=1}^6 P(x) = 1, \quad \sum_{x \in \text{vals}(X)} P(\textcircled{x}) = 1$$

random variable

Expected Value mean "weighted average"
 $E[X] = \sum_{x \in \text{vals}(X)} x p(x)$ empirical based on data Theo.

fair $X_f \Rightarrow E[X_f] = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + \dots + 6 \cdot \frac{1}{6}$
 $= \boxed{3.5}$

Weighted $X_w \Rightarrow E[X_w] = \frac{1}{10}(1+2+\dots+5) + \frac{1}{2} \cdot 6$
 $= \boxed{4.5}$

Variance

$$\mu = E[X]$$

$$E[(X - \mu)^2] = \sum_{x \in \text{vals}(X)} (x - \mu)^2 p(x)$$

Spread

fair $\Rightarrow \text{Var}(X_f) = \frac{1}{6} \left[(1 - 3.5)^2 + \dots + (6 - 3.5)^2 \right]$

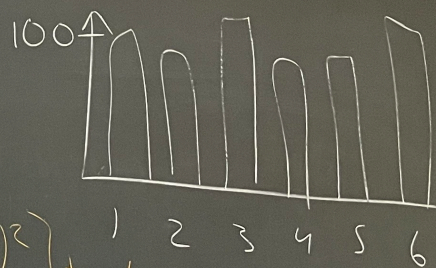
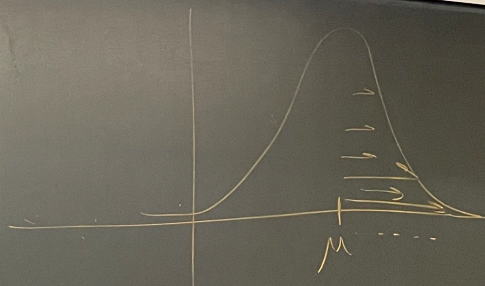
$$\approx 2.92$$

weighted $\Rightarrow \text{Var}(X_w) = \frac{1}{10} \left[(1 - 4.5)^2 + \dots + (5 - 4.5)^2 \right] + \frac{1}{2} (6 - 4.5)^2 = 3.25$

(empirical)

Sample variance $\Rightarrow \text{Var}(X) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \mu)^2$

weighting all examples equally.



$$\frac{1}{2} (6 - 4.5)^2 = 3.25$$

Outline for November 7

- Recap PCA and Handout 16
- Begin: statistics and hypothesis testing
- **Central limit theorem**
- Bootstrapping

Central Limit Theorem (CLT)

X_1, X_2, \dots, X_n are samples

from a population with mean μ & finite variance σ^2 ,

and \bar{X}_n is the sample mean.

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

THEN

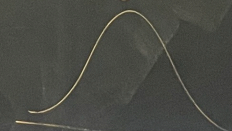
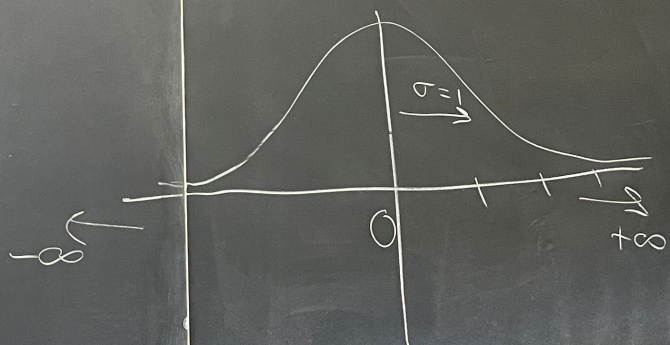
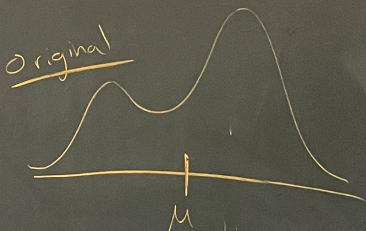
$$Z = \lim_{n \rightarrow \infty}$$

Normalizing $\sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$

$$\left(\frac{\bar{X}_n - \mu}{\left(\frac{\sigma}{\sqrt{n}} \right)} \right)$$

is a standard normal distribution

$N(0, 1)$
mean variance



probability density function (pdf)

$$\text{pdf} = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$
$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \boxed{1}$$

Hypothesis Testing

H_0 = null hypothesis (ex: die is fair)

H_1 = alternate hypothesis
(ex: die is weighted toward higher values)

apply CLT

z-score \nearrow

$$Z = \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}}$$
$$= \frac{4.2 - 3.5}{\sqrt{\frac{2.92}{20}}}$$

$$= \boxed{1.83}$$

0.7

← test statistic

$n=10$
 $\bar{X}_n = 4.2$

$n=20$

$\bar{X}_n = 4.2$ (sample mean)

one observation

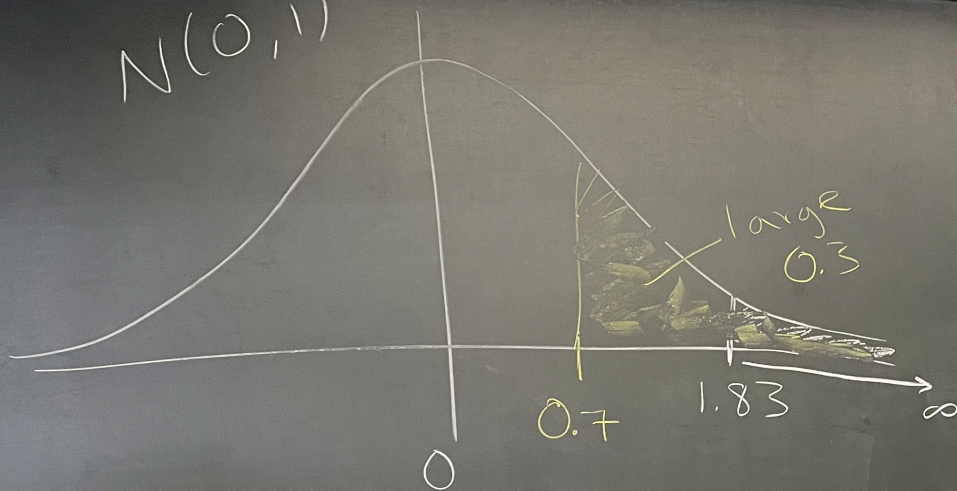
variance
 $\sigma^2 = 2.92$

std dev

$\sigma = \sqrt{2.92}$

sample
mean)
variation
ance

$N(0,1)$



p-value: probability of observing
a result as or more extreme
than yours under the null hypothesis

$$p\text{-value} = \int_{1.83}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

$$\approx \boxed{0.033}$$

usually compare with $0.05 = \alpha$

\Rightarrow reject null hypothesis Significance
level

Handout 17

Handout 17

$$(1) E[X] = 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = \frac{1}{2}$$

$$(2) \text{Var}(X) = \frac{1}{2} \left(0 - \frac{1}{2}\right)^2 + \frac{1}{2} \left(1 - \frac{1}{2}\right)^2 = \frac{1}{4}$$

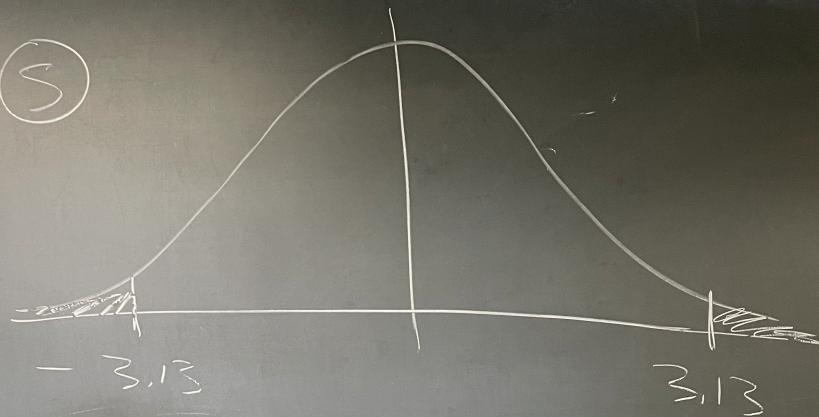
$$(3) \bar{X}_n = \frac{54}{80} = 0.675$$

$$(4) Z = \frac{0.675 - 0.5}{\sqrt{\frac{0.25}{80}}}$$

$$\approx 3.13$$

Z-score
test stat

⑤



$$p\text{-value} = 0.00175 \leq 0.05$$

$$\sqrt{\frac{92}{4}}$$

Outline for November 7

- Recap PCA and Handout 16
- Begin: statistics and hypothesis testing
- Central limit theorem
- Bootstrapping

Next time!