

Welcome !!  
😊

Plan

Logistics

Probability Review

Problem Solving

Set Size Estimation

go/rads 2024/

MTWR → 10-12 lecture

→ 2-3 problem solving

3-4 office hours

twitter @midd

Please come!

Please no computers!

Recommend: Read notes, come to class,  
Read notes again

## Grades

Participation (14)

1 pt per day

Problem Set (56)

one problem per class

3 pts for solution

1 pt self-grade

## Project (30)

Linear algebra  $\Rightarrow$  probability  
 $\Rightarrow$  algorithms

## Randomized Algorithms

So much data. Every day...

↳ NASA creates 20 terabytes

↳ 8 billion google searches

↳ 300 million terabytes internet

Moore's Law says compute doubles

But we're hitting physical limit

Processing data requires  
speed

Randomized lets us  
work in sublinear time

↳ Estimate # unique items?

↳ Similarity search

↳ Process matrices

↳ Reconstruct measurements

## Probability

$X$  is a random variable

If we flip a coin

$$X = \begin{cases} 1 & \text{wp } 1/2 \\ 0 & \text{wp } 1/2 \end{cases}$$

If we roll a die

$$X = \begin{cases} 5 \cdot 1 & \text{wp } 1/6 \\ 5 \cdot 2 & \text{wp } 1/6 \\ 5 \cdot 3 & \\ 5 \cdot 4 & \vdots \\ 5 \cdot 5 & \\ 5 \cdot 6 & \end{cases}$$

$\Pr(X=x)$  is the probability  
that  $X$  takes value  $x$

$$\Pr(X=4) = 1/6$$

$\uparrow$  die

$\mathbb{E}[X]$  is the expected value  
of  $X$

$$\mathbb{E}[X] = \sum_x x \cdot \Pr(X=x)$$

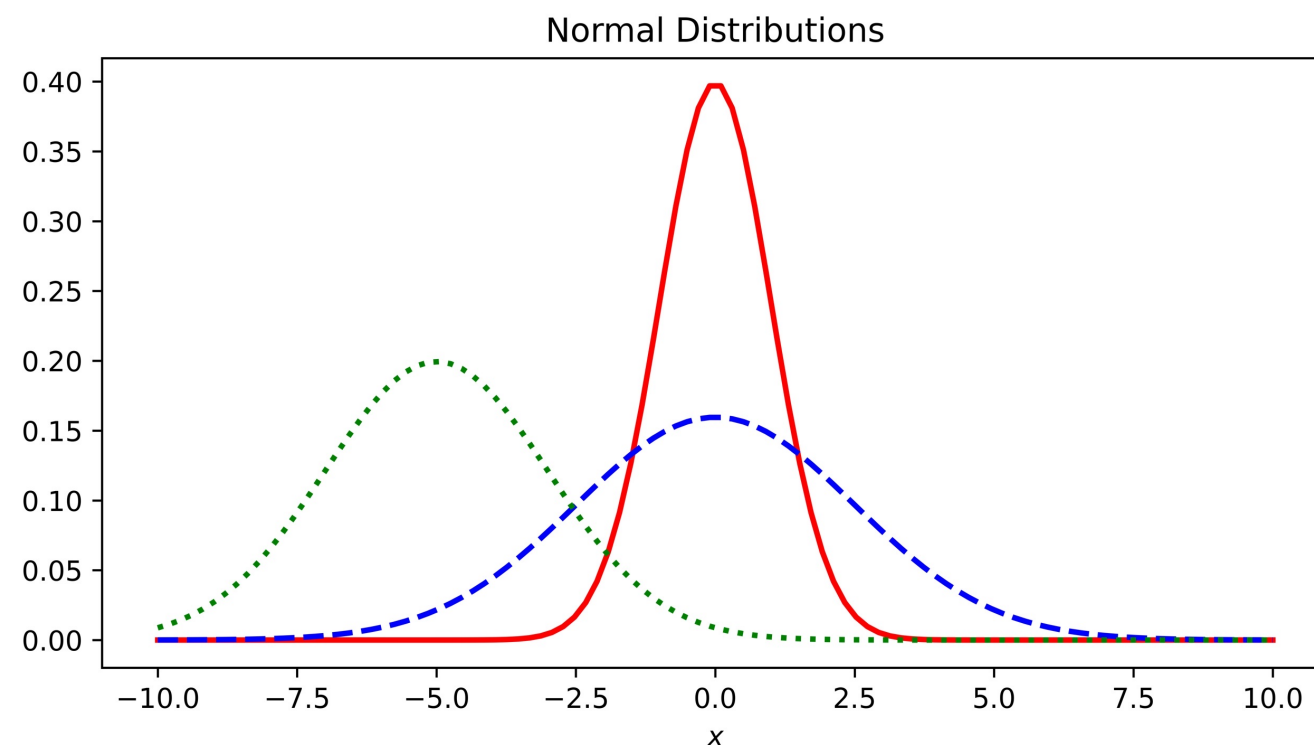
$$= 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 3 \cdot \frac{1}{6} + 4 \cdot \frac{1}{6} + 5 \cdot \frac{1}{6} + 6 \cdot \frac{1}{6}$$

$$= \frac{1}{6} (21) = 3.5$$

$\text{Var}(X)$  is the variance  
of a random variable:  
how much it varies  
from its expectation

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

↙ positive



$$\begin{aligned}\mathbb{E}[cX] &= \sum_x x \cdot c \cdot \text{Pr}(X=x) \\ &= c \sum_x x \cdot \text{Pr}(X=x) \\ &= c \cdot \mathbb{E}[X]\end{aligned}$$

$$\begin{aligned}\text{Var}(c \cdot X) &= \mathbb{E}[(cX - \mathbb{E}[cX])^2] \\ &= \mathbb{E}[(cX - c\mathbb{E}[X])^2] \\ &= \mathbb{E}[c^2 (X - \mathbb{E}[X])^2] \\ &= c^2 \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= c^2 \text{Var}(X)\end{aligned}$$

Events defined on rvs  
and variables

A = event that die is 1 or 2

B = event that die is odd

Does A give information  
about B?

$$\Pr(A \cap B) = \Pr(A) \cdot \Pr(B|A)$$

$$\frac{1}{6} = \frac{2}{6} \cdot \frac{1}{2}$$

$$\Pr(B|A) \triangleq \frac{\Pr(A \cap B)}{\Pr(A)}$$

A, B independent iff

$$\Pr(B|A) = \Pr(B)$$

$$\Leftrightarrow \Pr(A \cap B) = \Pr(A) \cdot \Pr(B)$$

X, Y are independent iff

$$\Pr(X=x | Y=y) = \Pr(X=x)$$

for all x, y

## Linearity of Expectation

$$\mathbb{E}[X + Y]$$

always?  
sometimes?  
never?

$$= \sum_x \sum_y (x+y) \Pr(X=x \cap Y=y)$$

$$= \sum_x x \sum_y \Pr(X=x \cap Y=y) + \sum_y y \sum_x \Pr(X=x \cap Y=y)$$

$$= \sum_x x \Pr(X=x) \underbrace{\sum_y \Pr(Y=y | X=x)}_1 + \sum_y y \cdot \Pr(Y=y) \cdot 1$$

$$= \mathbb{E}[X] + \mathbb{E}[Y]$$

$$E[XY] \stackrel{?}{=} E[X]E[Y]$$

always?  
 Sometimes?  
 never?

$$E[XY] = \sum_{x,y} xy \Pr(X=x \wedge Y=y)$$

$$\stackrel{\text{indep}}{=} \sum_x \sum_y xy \Pr(X=x) \Pr(Y=y)$$

$$= \sum_x x \Pr(X=x) \sum_y y \Pr(Y=y)$$

$$= E[X]E[Y]$$

$$X = \begin{cases} +1 & \text{wp } 1/2 \\ -1 & \text{wp } 1/2 \end{cases} \quad Y = X$$

$$E[X] = 1 \cdot 1/2 + (-1) \cdot 1/2 = 0$$

$$E[Y] = 0$$

$$E[X \cdot Y] = 1/2 \cdot 1 \cdot 1 + 1/2 \cdot (-1) \cdot (-1) = 1$$



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

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

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) \quad \text{when } X, Y \text{ are independent}$$

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

$$= E[(X - E[X]) + (Y - E[Y])]^2$$

$$= E[(X - E[X])^2 + 2(X - E[X])(Y - E[Y]) + (Y - E[Y])^2]$$

$$= E[(X - E[X])^2] + E[(Y - E[Y])^2] + 0$$

$$\begin{aligned} & E[XY - E[X] \cdot Y + E[X] \cdot E[Y] - X E[Y]] \\ &= E[XY] - E[X] \cdot E[Y] + E[X] \cdot E[Y] - E[Y] \cdot E[X] \\ &\stackrel{\text{indep}}{=} E[X]E[Y] \end{aligned}$$

$$X = \begin{cases} 1 & \text{wp } 1/2 \\ -1 & \text{wp } 1/2 \end{cases}$$

$$Y = -X$$

$$X = \begin{cases} 1 & \text{wp } 1/2 \\ -1 & \text{wp } 1/2 \end{cases}$$

$$Y = -X$$

$$E[X] = 1 \cdot 1/2 + (-1) \cdot 1/2 = 0$$

$$E[Y] = E[-1 \cdot X] = -1 \cdot E[X] = 0$$

$$\begin{aligned} \text{Var}(X) &= E[(X - 0)^2] = 1/2 \cdot 1^2 + 1/2 \cdot (-1)^2 \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{Var}(Y) &= \text{Var}(-1 \cdot X) = (-1)^2 \text{Var}(X) \\ &= \text{Var}(X) \end{aligned}$$

$$X + Y = \begin{cases} 1 + (-1) & \text{wp } 1/2 \\ -1 + 1 & \text{wp } 1/2 \end{cases}$$

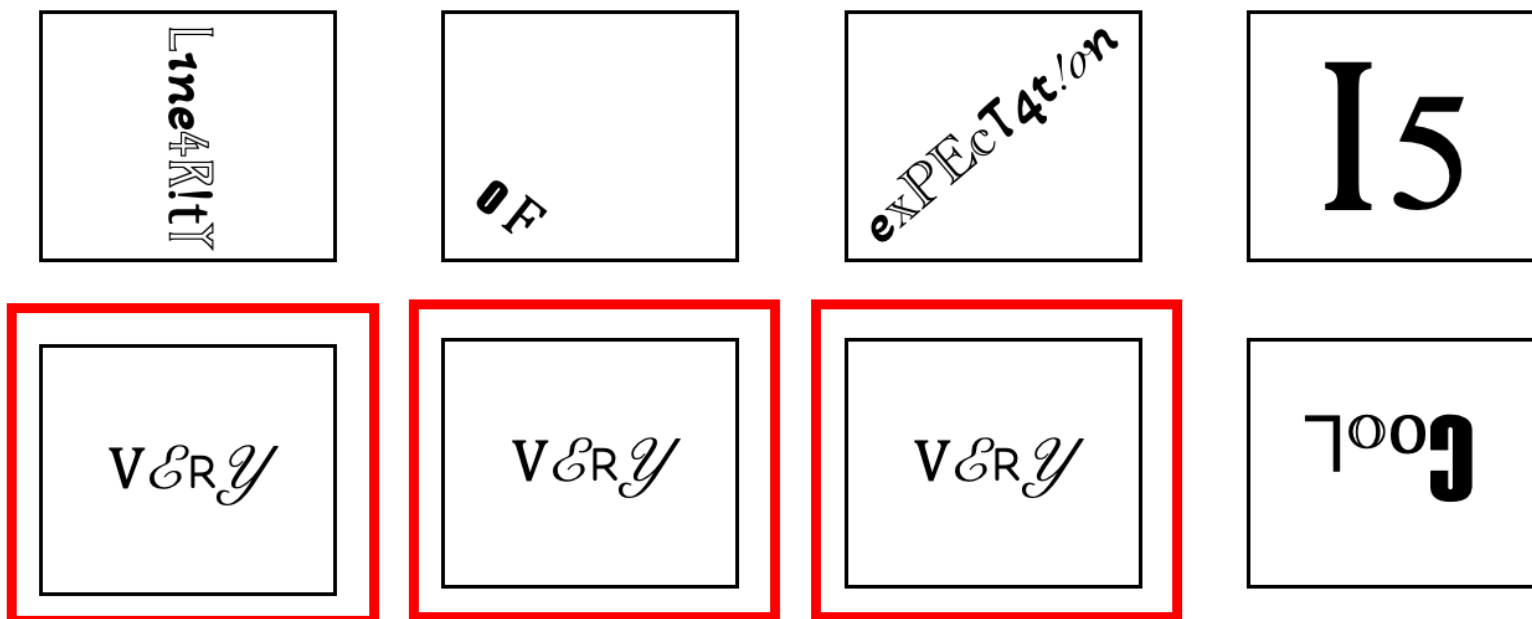
$$= 0 \quad \text{wp } 1$$

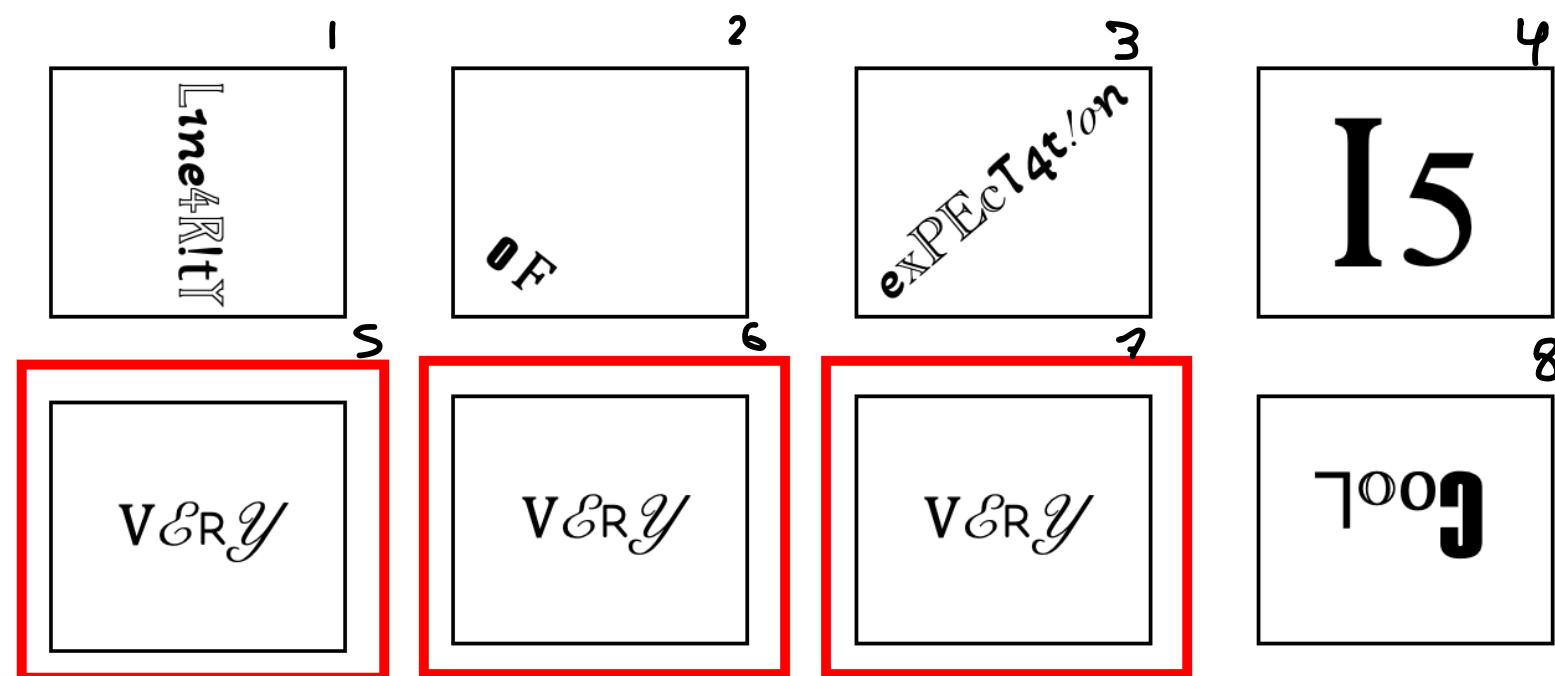
$$\text{Var}(X + Y) = 0 \neq 2 = \text{Var}(X) + \text{Var}(Y)$$

## Set Size Estimation

- ↳ internet traffic
- ↳ ecology
- ↳ social networks

naive: we keep making  
calls until we  
see 1 million CAPTCHAs





$$D_{i,j} = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ and } j^{\text{th}} \text{ item are same} \\ 0 & \text{else} \end{cases}$$

$$D_{1,2} = 0$$

$$D_{5,7} = 1$$

# samples

$$D = \sum_{i=1}^m \sum_{j=i+1}^m D_{i,j}$$

$$\begin{aligned} E[D_{i,j}] &= 1 \cdot \Pr(i^{\text{th}}, j^{\text{th}} \text{ same}) \\ &\quad + 0 \cdot \Pr(i^{\text{th}}, j^{\text{th}} \text{ not same}) \\ &= \Pr(i^{\text{th}}, j^{\text{th}} \text{ same}) \\ &= \frac{1}{n} \end{aligned}$$

$n = \# \text{ unique items}$

$$\begin{aligned} E[D] &= \sum_{i=1}^m \sum_{j=i+1}^m E[D_{i,j}] \\ &= \frac{1}{n} \sum_{i=1}^m \sum_{j=i+1}^m 1 \\ &= \frac{m(m-1)}{2 \cdot n} \end{aligned}$$

$$\mathbb{E}[D] = \frac{m(m-1)}{2 \cdot n}$$

Suppose we see  $D=10$   
when  $m=1000$ ,

$$n = 1,000,000$$

$$\mathbb{E}[D] = \frac{1000 \cdot 999}{2 \cdot 1,000,000} = .4995$$

### Markov's Inequality

$X$  non-negative r.v.

$$t > 0$$

$$\Pr(X \geq t) \leq \frac{\mathbb{E}[X]}{t}$$

$$\Pr(D \geq 10) \leq \frac{\mathbb{E}[D]}{10} = \frac{.4995}{10} = .04995$$

$$\begin{aligned}
 E[X] &= \sum_x x \Pr(X=x) \\
 &= \sum_{\substack{x \\ x \geq t}} x \Pr(X=x) + \sum_{\substack{x \\ x < t}} x \Pr(X=x)
 \end{aligned}$$

$$\geq \sum_{\substack{x \\ x \geq t}} t \Pr(X=x) + 0$$

$$\geq t \sum_{\substack{x \\ x \geq t}} \Pr(X=x) = t \Pr(X \geq t)$$

$$\Pr(X \geq t) \leq \frac{E[X]}{t}$$