



Algorithmic Analysis

CS109, Stanford University

Pset #4 is Due Tomorrow!

The screenshot shows a web browser window with the following elements:

- Browser Tab:** Pset 4 - Probabilistic Models
- Address Bar:** cs109psets.netlify.app/sum24/pset4/splash
- Page Title:** Pset 4 - Probabilistic Models For Joel Ramirez
- Navigation:** A vertical sidebar on the left contains a home icon and numbered links from 1 to 15.
- Content:**
 - A blue button labeled "Get Started".
 - A light blue box containing:
 - Due Date:** Tuesday, Aug 6, 5:00 PM Eastern Daylight Time (in 3 days).
 - Grace Period Date:** Tuesday, Aug 6, 6:00 PM Eastern Daylight Time (in 3 days).
 - Solutions Posted:** Friday, Aug 9, 6:00 PM Eastern Daylight Time (in 6 days).
 - A button labeled "Extension Request Forms" with a dropdown arrow.

<review>

Expectation (just cool)

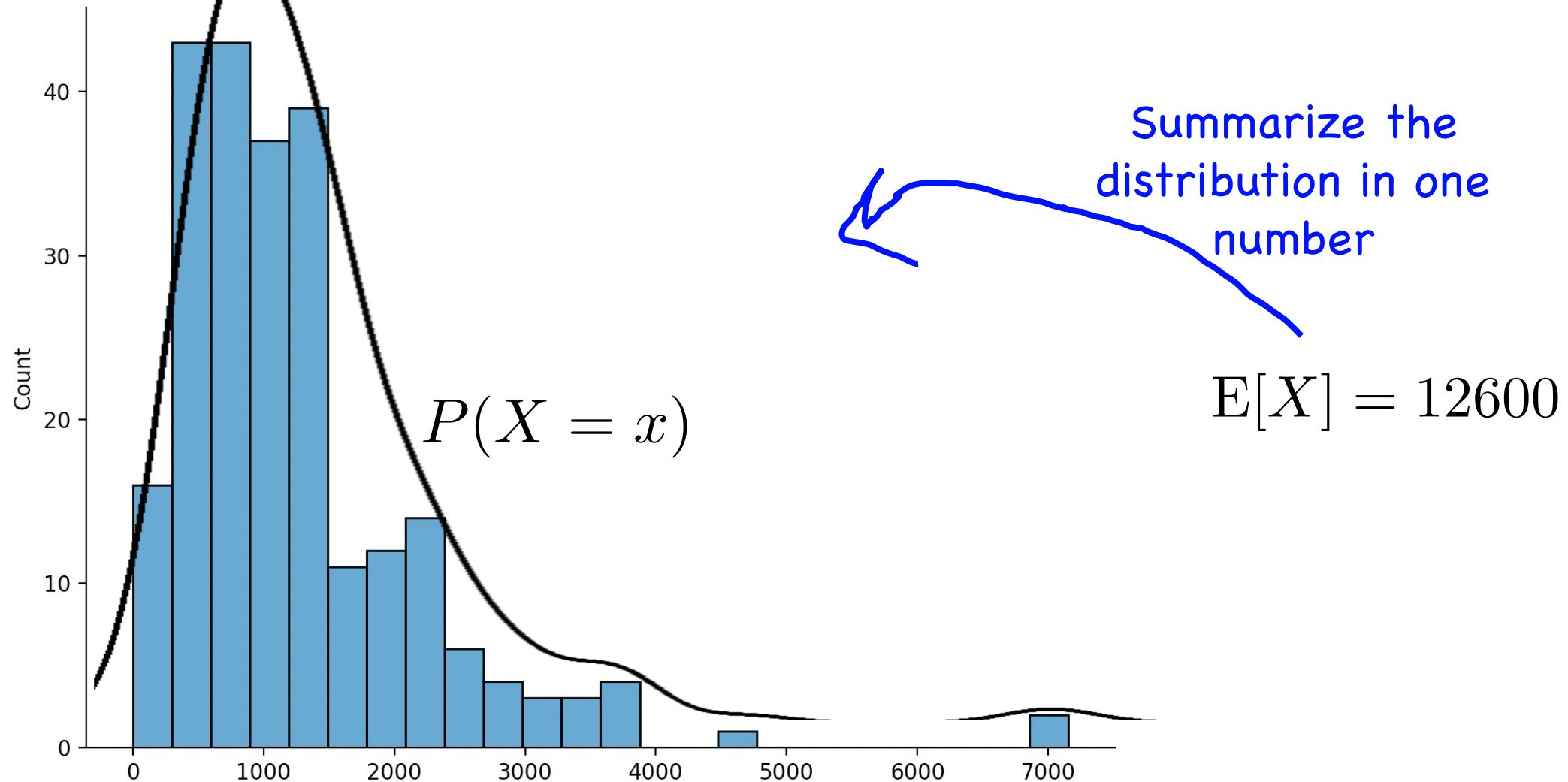
$$E[X] = \sum_x x \cdot P(X = x)$$

The probability that X takes on that value

All the values that X can take on

Limitation of Expectation (not cool?)

X = time to complete the medical diagnosis problem (in seconds)



Expectation of a Sum (much cool)

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

Generalized:
$$E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i]$$

Holds regardless of dependency between X_i 's

We shouldn't be surprised about this though...

Expectation of a Sum

What does this even mean?

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

Whenever you have two variables inside of the expectation you should think about how they interact *jointly*.

We enumerate every possible value of $x + y$.

Just like you've already been doing with individual variables).

$$\mathbb{E}[X + Y] = \sum_{x,y} (x + y)P(X = x, Y = y)$$

Let's exercise our math muscles to show that they're equivalent!
(Or maybe you can already tell!)

Expectation of a Function

Law of unconscious statistician

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

So for example...

$$\mathbb{E}[X^2] = \sum_x x^2 \cdot P(X = x)$$

Expectation of a Function

Law of the unconscious statistician

$$E[g(X)] = \sum_x g(x) \cdot P(X = x)$$

This one I won't prove...it's sorta lengthy...

Consider what we would do if a function outputs the same values for some x (e.g. x^2)...

This is why we call it the LOTUS

Bool was Cool



Boole was Cool

Let E_1, E_2, \dots, E_n be events with indicator RVs X_i

- Indicators are Bernoulli RVs that represent whether an event took place
- If event E_i occurs, then $X_i = 1$, else $X_i = 0$
- Hmmm: $E[X_i] = P(E_i)$

$$E[X_i] = 0 \cdot (1 - P(E_i)) + 1 \cdot P(E_i)$$

Bernoulli aka Indicator Random Variables were studied extensively by George Boole
Boole died of being too cool

Boole was Cool until he died

George Boole / Cause of death

Pneumonia

In 1864, Boole died due to **fever-induced pleural effusion after developing pneumonia**. Boole published around 50 articles and several separate publications in his lifetime.

An aside...things to care about



Differential Privacy

Aims to provide means to **maximize the accuracy** of probabilistic queries while minimizing the **probability** of identifying its records.



Cynthia Dwork's celebrity lookalike is Cynthia Dwork.

Differential Privacy

100 independent values $X_1 \dots X_{100}$ where $X_i \sim \text{Bern}(p)$

Y_i

What is
returned

```
# Maximize accuracy, while preserving privacy.  
def calculateYi(Xi):  
    obfuscate = random()  # random() returns True  
                           # or False with equal  
                           # likelihood  
    if obfuscate:  
        return indicator(random())  
    else:  
        return Xi
```

Differential Privacy

100 independent values $X_1 \dots X_{100}$ where $X_i \sim \text{Bern}(p)$

Y_i
What is
returned

```
# Maximize accuracy, while preserving privacy.  
def calculateYi(Xi):  
    obfuscate = random()           random() returns True  
    if obfuscate:                  or False with equal  
        return indicator(random()) likelihood  
    else:  
        return Xi
```

What is $E[Y_i]$?

$$E[Y_i] = P(Y_i = 1) = \frac{p}{2} + \frac{1}{4}$$

Differential Privacy

100 independent values $X_1 \dots X_{100}$ where $X_i \sim \text{Bern}(p)$

Y_i
What is
returned

```
# Maximize accuracy, while preserving privacy.  
def calculateYi(Xi):  
    obfuscate = random()           random() returns True  
    if obfuscate:                  or False with equal  
        return indicator(random()) likelihood  
    else:  
        return Xi
```

$$E[Y_i] = P(Y_i = 1) = \frac{p}{2} + \frac{1}{4}$$

$P(Y=1) = P(Y=1 | \text{obfuscate}=\text{True}) \cdot P(\text{obfuscate}=\text{True}) + P(Y=1 | \text{obfuscate}=\text{False}) \cdot P(\text{obfuscate}=\text{False})$

Let's derive this to convince you!

Differential Privacy

100 independent values $X_1 \dots X_{100}$ where $X_i \sim \text{Bern}(p)$

Y_i
What is
returned

```
# Maximize accuracy, while preserving privacy.  
def calculateYi(Xi):  
    obfuscate = random()           random() returns True  
    if obfuscate:                  or False with equal  
        return indicator(random()) likelihood  
    else:  
        return Xi
```

Let $Z = \sum_{i=1}^{100} Y_i$

What is $E[Z]$?

$$E[Z] = E\left[\sum_{i=1}^{100} Y_i\right] = \sum_{i=1}^{100} E[Y_i] = \sum_{i=1}^{100} \left(\frac{p}{2} + \frac{1}{4}\right) = 50p + 25$$

Differential Privacy

100 independent values $X_1 \dots X_{100}$ where $X_i \sim \text{Bern}(p)$

Y_i
What is
returned

```
# Maximize accuracy, while preserving privacy.  
def calculateYi(Xi):  
    obfuscate = random()           random() returns True  
    if obfuscate:                  or False with equal  
        return indicator(random()) likelihood  
    else:  
        return Xi
```

Let $Z = \sum_{i=1}^{100} Y_i$ $E[Z] = 50p + 25$ How do you estimate p ?

Challenge: What is the probability that our estimate is good?

End Review

And aside....

Computer Cluster Utilization

Computer cluster with k servers

- Requests independently go to server i with probability p_i
- Let event A_i = server i receives no requests
- X = # of events A_1, A_2, \dots, A_k that occur
- $E[X]$ after first n requests?

Since X is an expectation, can you express it as a sum?

-
- Let Bernoulli B_i be an **indicator** for A_i $X = \sum_{i=1}^k B_i$
 - Since requests independent: $P(A_i) = (1 - p_i)^n$

$$E[X] = E\left[\sum_{i=1}^k B_i\right] = \sum_{i=1}^k E[B_i] = \sum_{i=1}^k P(A_i) = \sum_{i=1}^k (1 - p_i)^n$$

Amazon Monetized This

amazon





* ~~52%~~ 74% of Amazons Profits

**More profitable than Amazon's North
America commerce operations



When stuck, brainstorm
about random variables



Toy Collection

You are trying to collect n distinct toys (and there are only n distinct toys)

- Each purchase, each toy is equally likely (buying with replacement)
- Let $X = \#$ purchases until you have ≥ 1 of each toy. What is $E[X]$?

Let $X_i = \#$ of **trials to get success after i -th success** where “success” is getting an unseen toy.

$$X = X_0 + X_1 + \dots + X_{n-1} \Rightarrow E[X] = E[X_0] + E[X_1] + \dots + E[X_{n-1}]$$

$$X_i \sim \text{Geo} \left(p = \frac{n-i}{n} \right)$$

After i successes, the probability of the next success is

$$p = (n-i) / n$$

$$E[X_i] = \frac{1}{p} = \frac{n}{n-i}$$

$$E[X] = \frac{n}{n} + \frac{n}{n-1} + \frac{n}{n-2} + \dots + \frac{n}{1} = n \left[\frac{1}{n} + \frac{1}{n-1} + \dots + 1 \right] = O(n \log n)$$

Break



Conditional Expectation

Conditional Expectation

X and Y are discrete random variables:

$$E[X|Y = y] = \sum_x xP(X = x|Y = y)$$



Analogously, continuous random variables:

$$E[X|Y = y] = \int_x xP(X = x|Y = y)$$

Conditional Expectation

$$E[X|Y = y] = \sum_x x \cdot P(X = x|Y = y)$$

Roll two 6-sided dice D_1 and D_2

- $X = \text{value of } D_1 + D_2$ $Y = \text{value of } D_2$
- What is $E[X | Y = 6]$?

$$\begin{aligned} E[X | Y = 6] &= \sum_x x P(X = x | Y = 6) \\ &= \left(\frac{1}{6}\right)(7 + 8 + 9 + 10 + 11 + 12) = \frac{57}{6} = 9.5 \end{aligned}$$

- Intuitively makes sense: $6 + E[\text{value of } D_1] = 6 + 3.5$

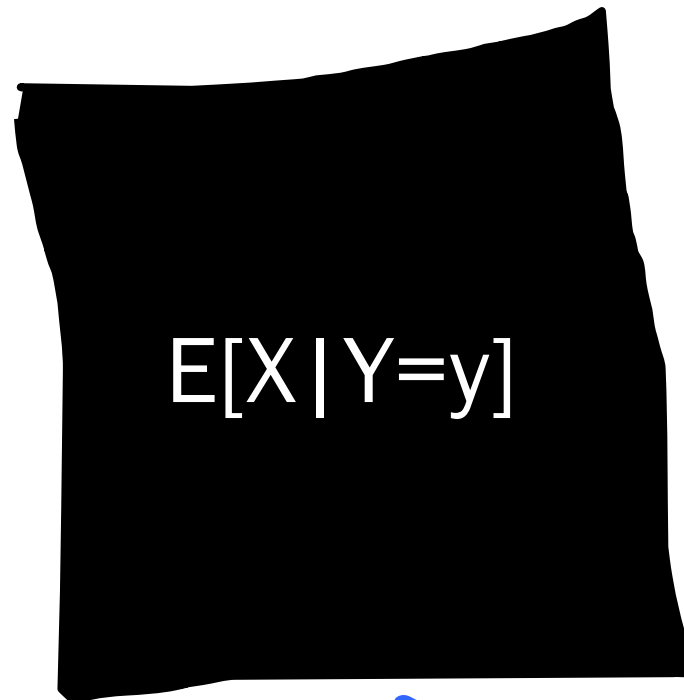
Law of Total Expectation

Conditional Expectation

$$E[X|Y = y] = \sum_x x \cdot P(X = x|Y = y)$$

Define $g(Y) = E[X | Y]$

This is a function of Y



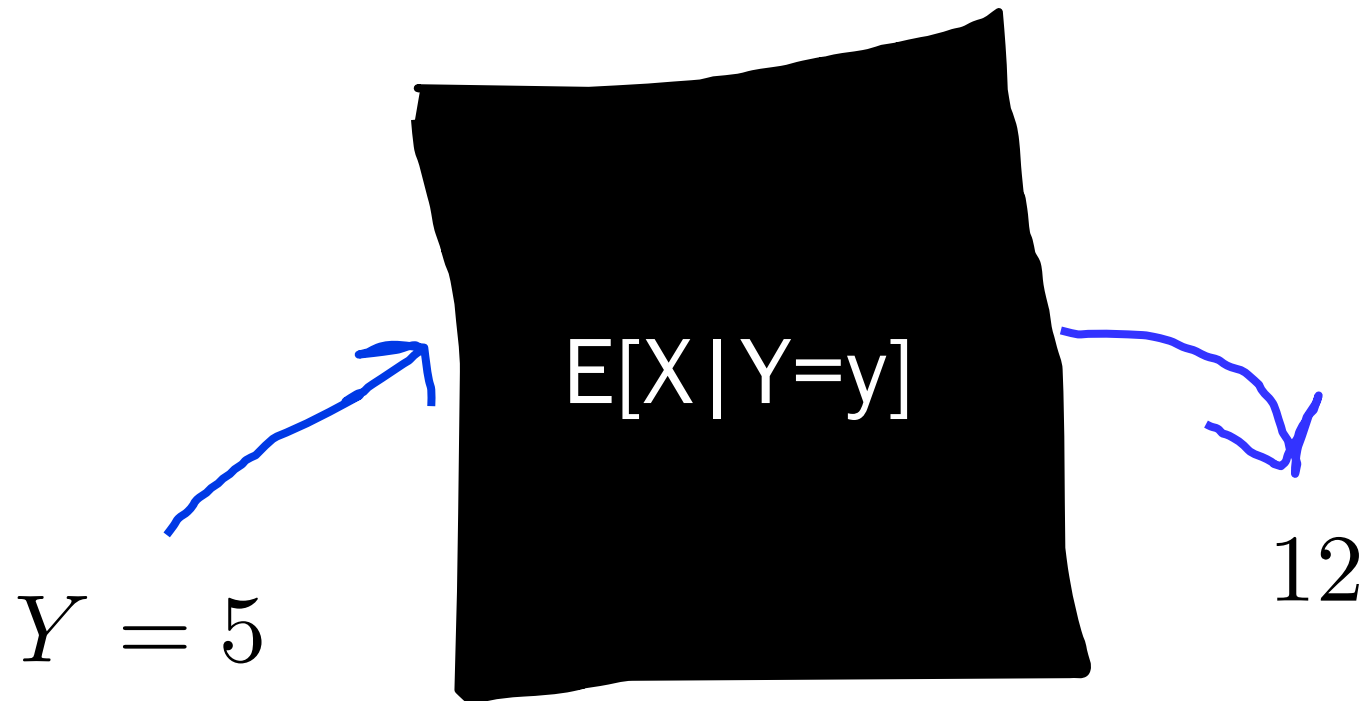
This is a function with Y as input

Conditional Expectation

$$E[X|Y = y] = \sum_x x \cdot P(X = x|Y = y)$$

Define $g(Y) = E[X | Y]$

This is a function of Y

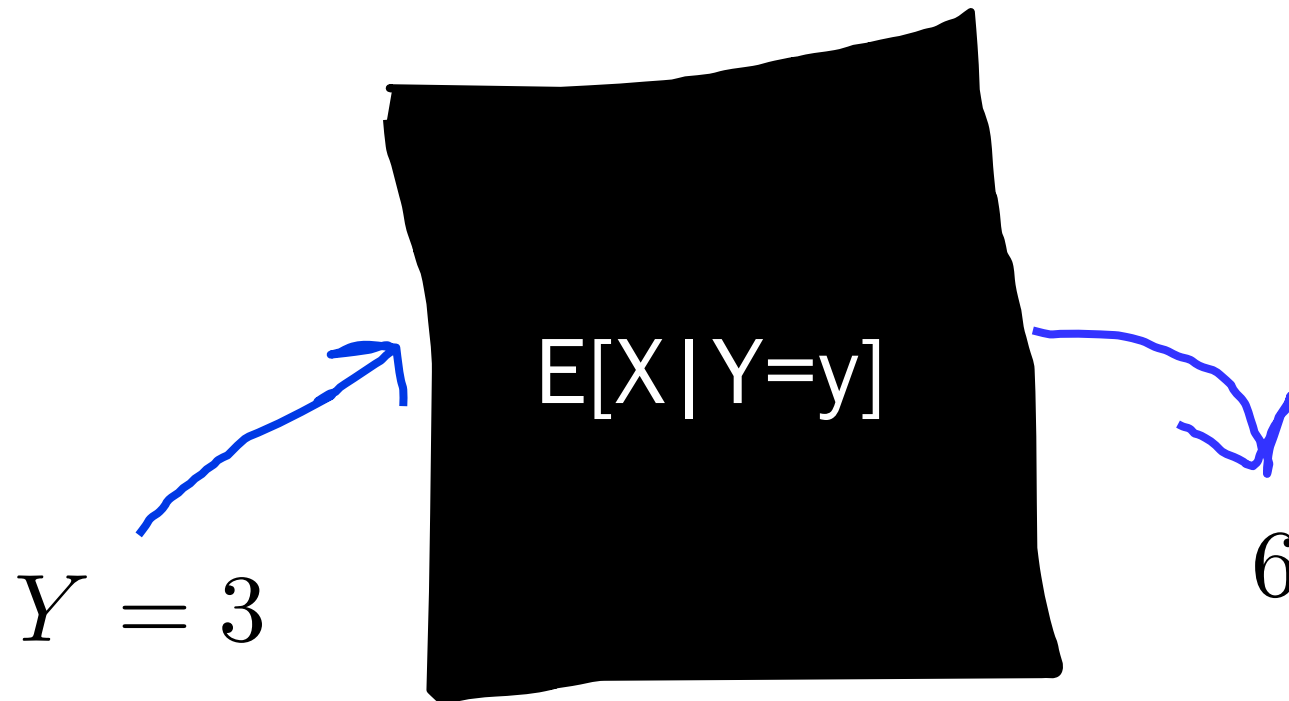


Conditional Expectation

$$E[X|Y = y] = \sum_x x \cdot P(X = x|Y = y)$$

Define $g(Y) = E[X | Y]$

This is a function of Y



Conditional Expectation

$$E[X|Y = y] = \sum_x x \cdot P(X = x|Y = y)$$

This is a number:

$$E[X]$$



This is a function of y :

$$E[X|Y = y]$$

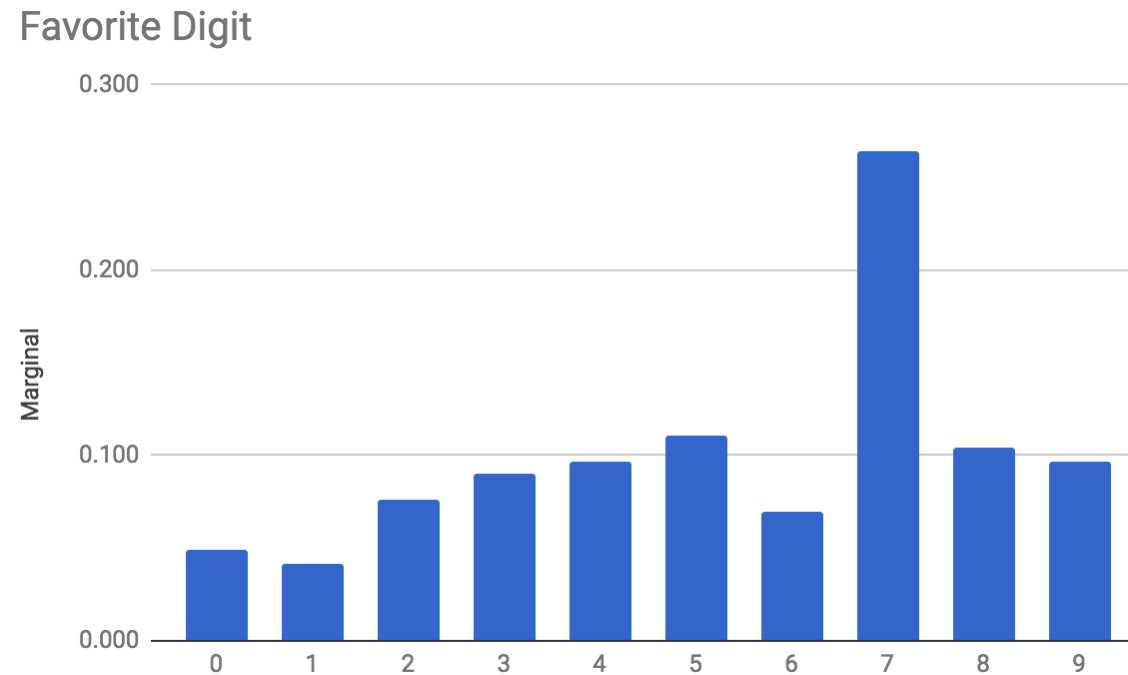
$$E[X = 5]$$

Doesn't make sense. Take expectation of random variables, not events

Expectation

$$E[X|Y = y] = \sum_x x \cdot P(X = x|Y = y)$$

X = favorite number
Y = year in school



$$E[X] = 0 * 0.05 + \dots + 9 * 0.10 = 5.38$$

Conditional Expectation

$$E[X|Y = y] = \sum_x x \cdot P(X = x|Y = y)$$

X = favorite number

Y = year in school

E[X | Y] ?

Year in school, Y = y	E[X Y = y]
2	5.5
3	5.8
4	6.0
5	4.7

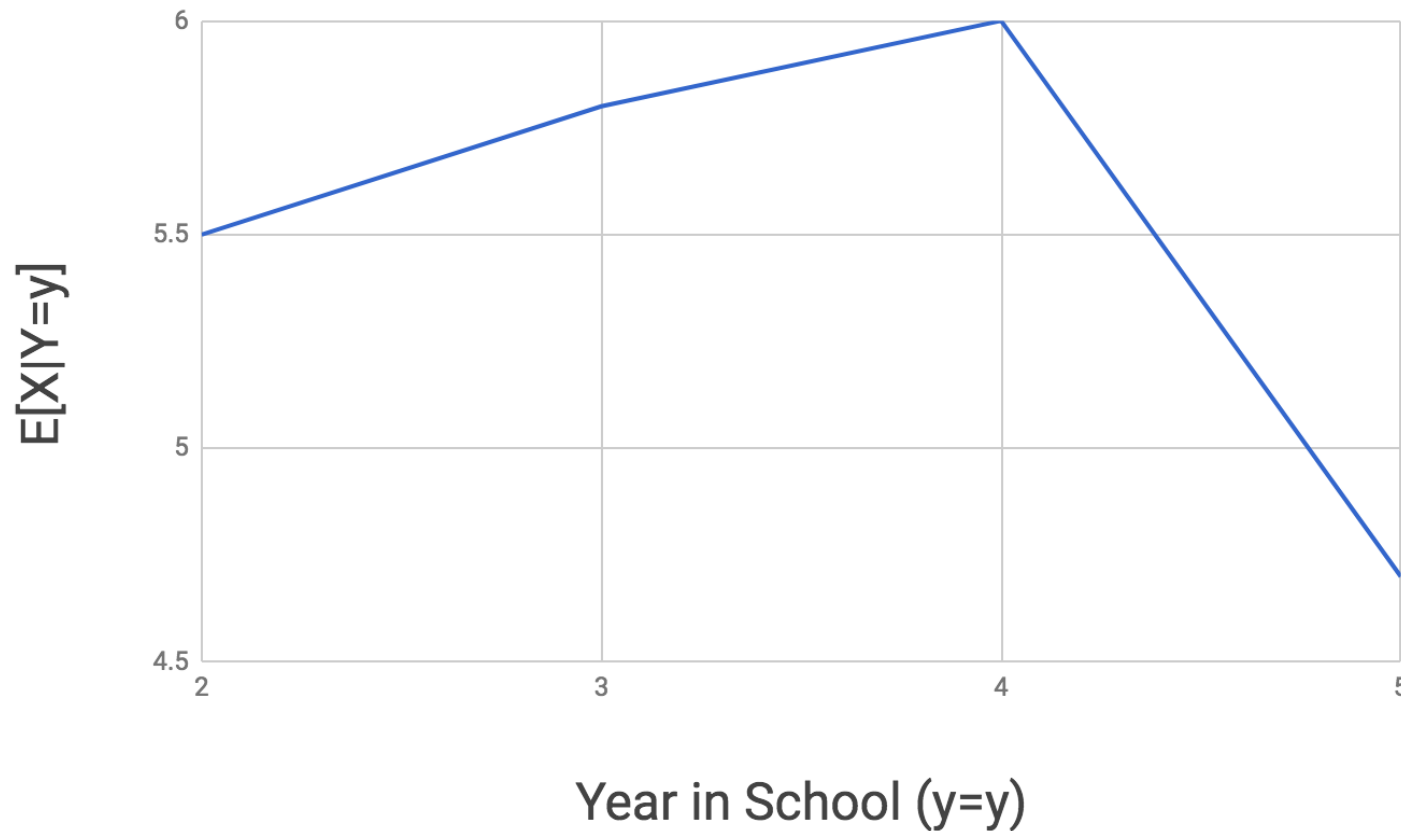
Conditional Expectation

$$E[X|Y = y] = \sum_x x \cdot P(X = x|Y = y)$$

X = favorite number

Y = year in school

$E[X | Y] ?$

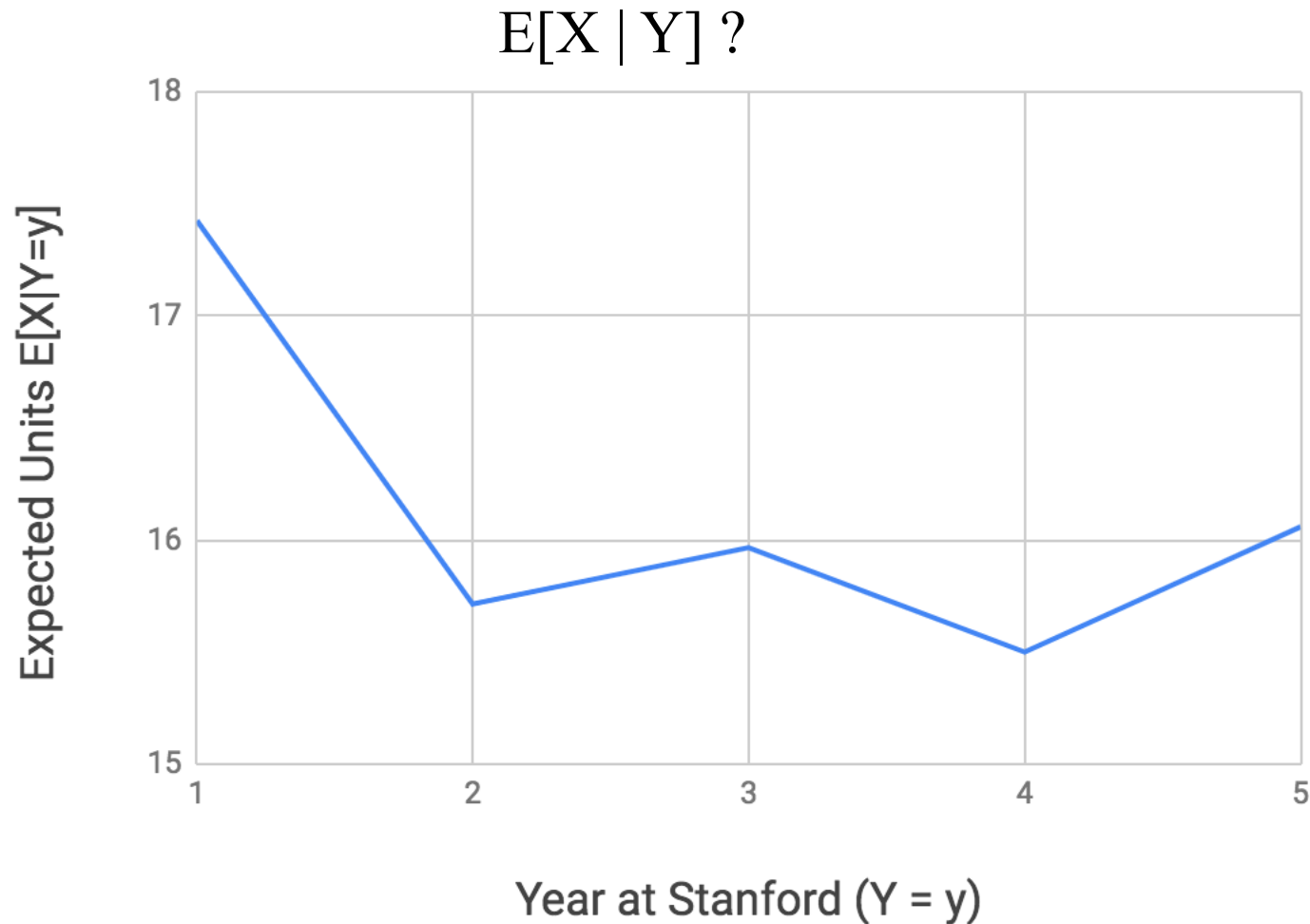


Conditional Expectation

$$E[X|Y = y] = \sum_x x \cdot P(X = x|Y = y)$$

X = units in fall quarter

Y = year in school



Want to see something cool?

What the heck does this give you?

$$\sum_y E[X|Y = y]P(Y = y)$$

Law of Total Expectation

$$E[X] = \sum_y E[X|Y = y]P(Y = y)?$$

$$\sum_y E[X|Y = y]P(Y = y) = \sum_y \sum_x xP(X = x|Y = y)P(Y = y)$$

Def of $E[X|Y]$

$$= \sum_y \sum_x xP(X = x, Y = y)$$

Chain rule!

$$= \sum_x \sum_y xP(X = x, Y = y)$$

I switch the order of the sums

$$= \sum_x x \sum_y P(X = x, Y = y)$$

Move that x outside the y sum

$$= \sum_x xP(X = x)$$

Marginalization

$$= E[X]$$

Def of $E[X]$

Law of Total Expectation



For any discrete random variable X
and any discrete random variable Y

$$E[X] = \sum_y E[X|Y = y]P(Y = y)$$

Recall the Law of Total Probability

$$P(X = x) = \sum_y P(X = x|Y = y)P(Y = y)$$

NETFLIX

(The Streaming Part)

How long does this code take to run?

Netflix streams millions of hours of videos per day. They REALLY care about the speed of the following code:

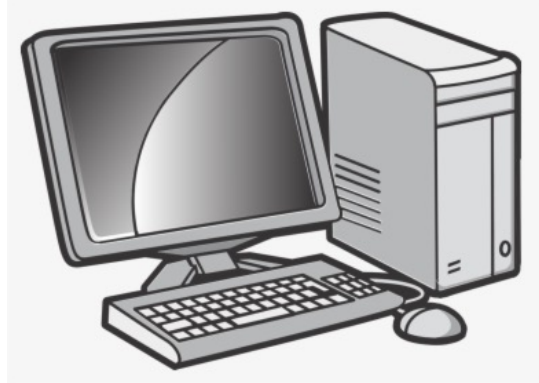
```
database.get_movie(movie_name)
```

How long does this line of code take? Say 512 MB movie.

1. 0.3s
2. 1.6 mins
3. 5 mins
4. 2 hours

All are correct! It is a RV!

How long does this code take to run?

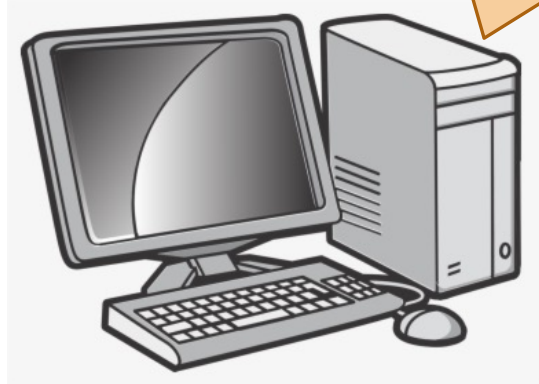


```
database.get_movie(movie_name)
```

1. 0.3s
2. 1.6 mins
3. 5 mins
4. 2 hours

Millisecond Latency

I have the file



```
database.get_movie(movie_name)
```

1. 0.3s
2. 1.6 mins
3. 5 mins
4. 2 hours

Minute Latency



1. 0.3s
2. **1.6 mins**
3. 5 mins
4. 2 hours

Many Minutes Latency

database.get_movie(movie_name)



私はファイルを持っています



1. 0.3s
2. 1.6 mins
3. 5 mins
4. 2 hours

Are we done?

```
database.get_movie(movie_name)
```



5mins across the world!!!!!!

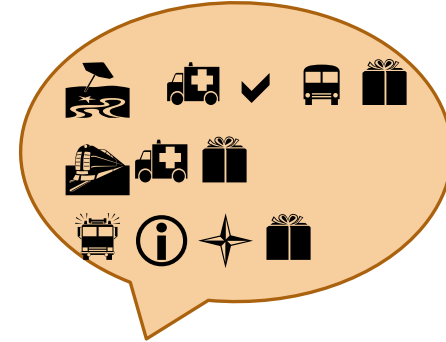


```
database.get_movie(movie_name)
```



Anyways

```
database.get_movie(movie_name)
```



1. 0.3s
2. 1.6 mins
3. 5 mins
4. 2 hours

Expected Run Time

Expected runtime (single file from database)

Assume the file location is distributed with the PMF:

- $P(\text{file on computer}) = 0.10$
- $P(\text{file in SoCal}) = 0.50$
- $P(\text{file in Japan}) = 0.37$
- $P(\text{file in Space}) = 0.03$

What is the expected runtime of `database.get_movie(movie_name)`

$$\begin{aligned}\mathbb{E}[\text{get_movie_database_time}] &= \mathbb{E}[\text{get_movie_database_time}|\text{Home}] \cdot \mathbb{P}(\text{Home}) \\ &+ \mathbb{E}[\text{get_movie_database_time}|\text{SoCal}] \cdot \mathbb{P}(\text{SoCal}) \\ &+ \mathbb{E}[\text{get_movie_database_time}|\text{Japan}] \cdot \mathbb{P}(\text{Japan}) \\ &+ \mathbb{E}[\text{get_movie_database_time}|\text{Space}] \cdot \mathbb{P}(\text{Space}) \\ &= 0.1 \cdot 0.3s + 0.5 \cdot 1.6min + 0.37 \cdot 5min + 0.03 \cdot 2hours \\ &\approx 6.25mins\end{aligned}$$

Times From Before:

1. Home: 0.3s
2. SoCal: 1.6 mins
3. Japan: 5 mins
4. Space: 2 hours

Analyze Recursive Code

Analyzing Recursive Code

```
int Recurse() {  
    int x = randomInt(1, 3); // Equally likely values  
  
    if (x == 1) return 3;  
    else if (x == 2) return (5 + Recurse());  
    else return (7 + Recurse());  
}
```

Let Y = value returned by `Recurse()`. What is $E[Y]$?

$$E[Y] = E[Y | X = 1]P(X = 1) + E[Y | X = 2]P(X = 2) + E[Y | X = 3]P(X = 3)$$

$$E[Y | X = 1] = 3$$

$$E[Y | X = 2] = E[5 + Y] = 5 + E[Y]$$

$$E[Y | X = 3] = E[7 + Y] = 7 + E[Y]$$

$$E[Y] = 3(1/3) + (5 + E[Y])(1/3) + (7 + E[Y])(1/3) = (1/3)(15 + 2E[Y])$$

$$E[Y] = 15$$

Uncertainty Theory

Beta
Distributions

Thompson
Sampling

Adding
Random Vars

Central Limit
Theorem


Sampling

Bootstrapping

Algorithmic
Analysis

Where are we in CS109?


On Wednesday...


Counting
Theory


Core
Probability

x_2
Random
Variables


Probabilistic
Models


Uncertainty
Theory


Machine
Learning

