



CS 109 Review

Noah Arthurs

CS 109



topics

machine learning

sampling, making conclusions from data

random variables / distributions

core probability fundamentals

skills

interpreting word problems into math

analyzing and producing code

methods

examples

demos

problem-solving

stories and memes!



CS 109

topics

machine learning

sampling, making conclusions from data

random variables / distributions

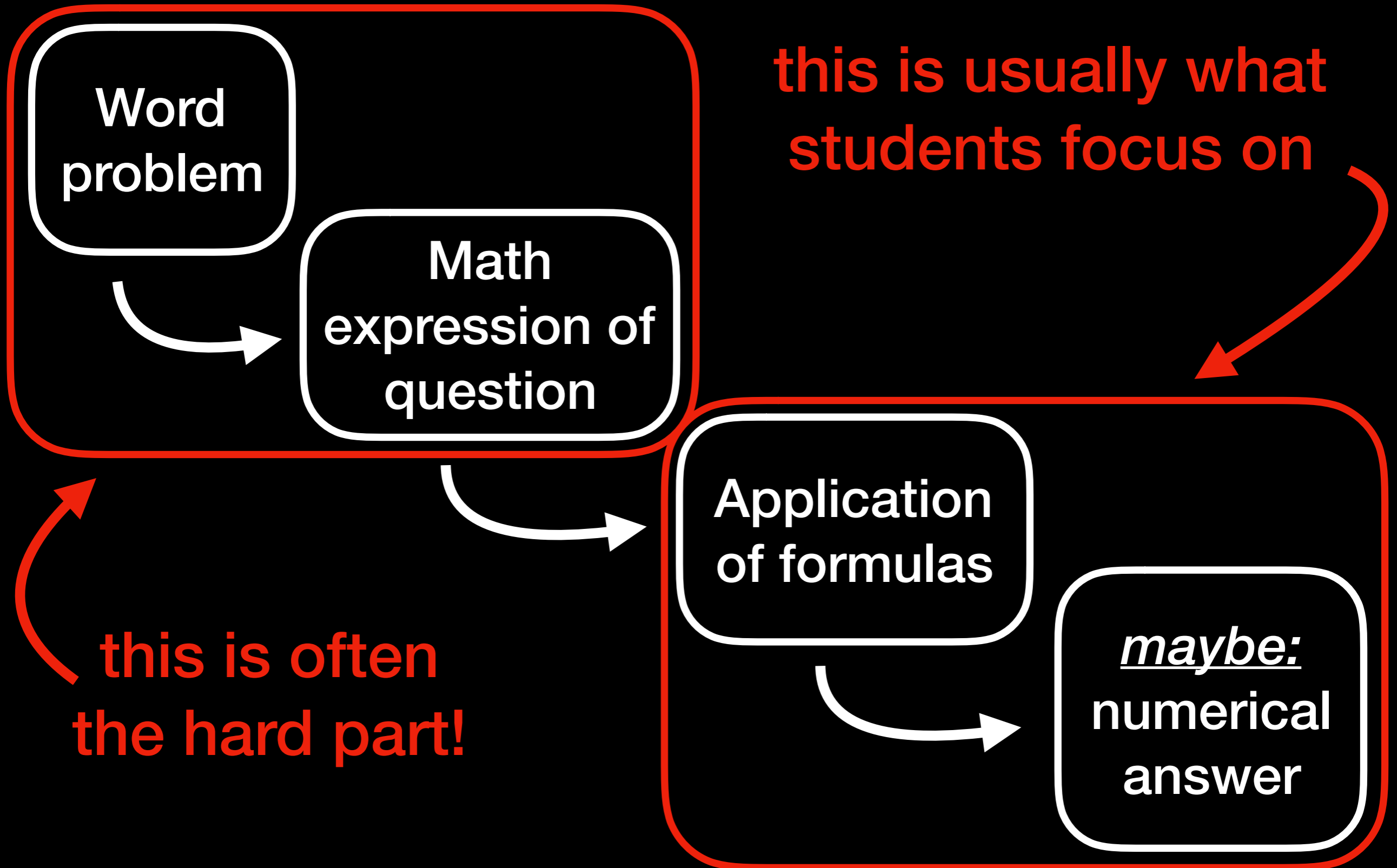
core probability fundamentals

skills

interpreting word problems into math

analyzing and producing code

Solving a CS109 problem



Step 1: Defining Your Terms

1. Let X represent _____
2. $X \sim$ _____
3. We want to know _____

Step 1: Defining Your Terms

Problem: If you roll 100 dice, what is the probability of getting less than 30 2's and 5's?

- 1. Let X represent the number of 2's and 5's we roll**
- 2. $X \sim \text{Binom}(100, 2/6)$**
- 3. We want to know $P(X < 30)$**

Translating English to Probability

<u>What the problem asks:</u>	<u>What you should immediately think:</u>
“What’s the probability of _____”	$P(\quad)$
“_____ given _____”, “_____ if _____”	$\quad \quad$
“at least _____”	<i>could we use what we know about everything less than ___?</i>
“approximate _____.”	<i>use an approximation!</i>
“How many ways...”	<i>combinatorics</i>

these are just a few, and these are why practice is the best way to prepare for the exam!



CS 109

topics

machine learning

sampling, making conclusions from data

random variables / distributions

core probability fundamentals

skills

interpreting word problems into math

analyzing and producing code

Code in CS 109

Analysis

Expectation of
binary tree depth

Bloom Filter Analysis

Expectation of
recursive die roll game

Implementation

Dithering

CO2 Levels

Biometric Keystrokes

Titanic

Peer Grading

Thompson Sampling



CS 109

topics

machine learning

sampling, making conclusions from data

random variables / distributions

core probability fundamentals

counting

conditional probability

probability principles

Counting

Sum Rule	Inclusion-Exclusion Principle
$outcomes = A + B $ if $ A \cap B = 0$	$ A + B - A \cap B $ for any $ A \cap B $
Product Rule	Pigeonhole Principle
$outcomes = A \times B $ if all outcomes of B are possible regardless of the outcome of A	If m objects are placed into n buckets, then at least one bucket has at least $\lceil m / n \rceil$ objects.

Combinatorics: Arranging Items

**Permutations
(ordered)**

**Combinations
(unordered)**

Distinct

$$n!$$

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Indistinct

$$\frac{n!}{k_1!k_2!\dots k_n!}$$

$$\binom{n+r-1}{r-1}$$

the divider method!

Probability basics

$$P(E) = \lim_{x \rightarrow \infty} \frac{n(E)}{n} \quad \text{in the general case}$$

$$\text{Probability} = \frac{\text{Event space}}{\text{Sample space}}$$

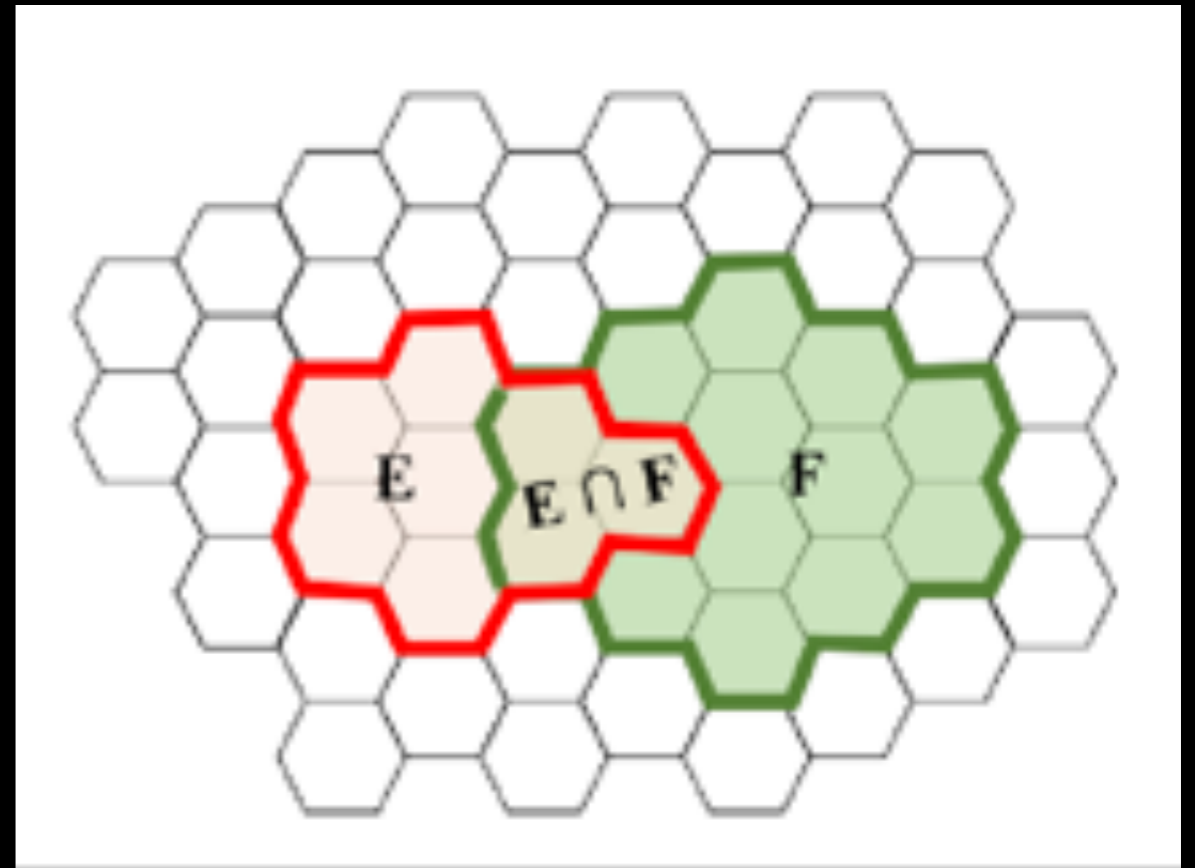
if all outcomes are equally likely!
(use counting with distinct objects)

Axioms: $0 \leq P(E) \leq 1$ $P(S) = 1$ $P(E^C) = 1 - P(E)$

Conditional Probability

definition:

$$P(E|F) = \frac{P(EF)}{P(F)}$$



Chain Rule:

$$* P(EF) = P(E \cap F)$$

$$P(EF) = P(E|F)P(F)$$

Law of Total Probability

Let's say we don't know $P(A)$, but we do know the probability of A given any value of B :

$$\begin{aligned}P(A) &= P(AB) + P(AB^C) \\ &= P(A|B)P(B) + P(A|B^C)P(B^C)\end{aligned}$$

If B can take on any value in S :

$$P(A) = \sum_{b \in S} P(A, B = b)$$

$$P(A) = \sum_{b \in S} P(A|B = b)P(B = b)$$

Bayes' Rule

$$P(E|F) = \frac{P(F|E)P(E)}{P(F)}$$

Bayes' Rule

posterior

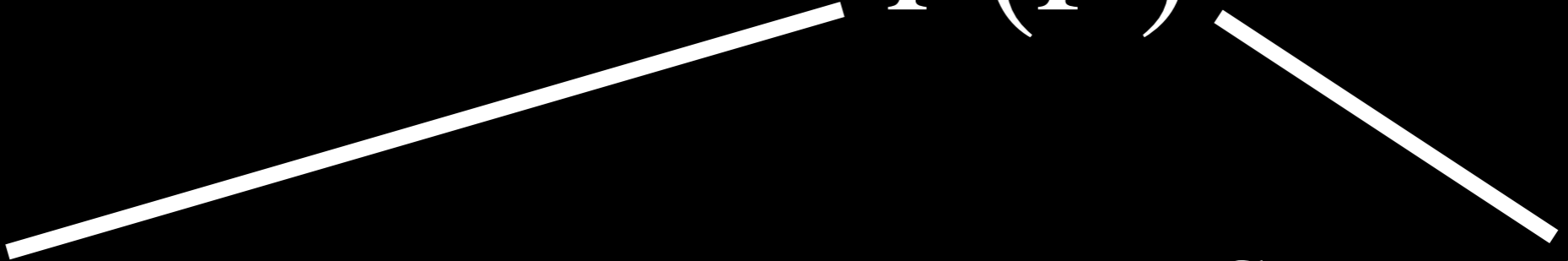
likelihood

prior

$$P(E|F) = \frac{P(F|E)P(E)}{P(F)}$$

normalization constant

Bayes' Rule

$$P(E|F) = \frac{P(F|E)P(E)}{P(F)}$$

$$P(F|E)P(E) + P(F|E^C)P(E^C)$$

divide the event F into all the possible ways it can happen; use LoTP

Which rule when?

$$P(EF) = P(E|F)P(F)$$

- Goes from an “and” to a conditional or vice versa
- Think about which event you want to condition on

$$P(A) = P(A|B)P(B) + P(A|B^C)P(B^C)$$

- We don't know about A but we do know about A|B
- Don't forget about the “and” version and “summation” version

$$P(E|F) = \frac{P(F|E)P(E)}{P(F)}$$

- Good for when E|F is hard but F|E is not so hard
- Common mistake: not trying chain rule first

Old Principles, New Tricks

Name of Rule	Original Rule	Conditional Rule
First axiom of probability	$0 \leq P(E) \leq 1$	$0 \leq P(E G) \leq 1$
Complement Rule	$P(E) = 1 - P(E^C)$	$P(E G) = 1 - P(E^C G)$
Chain Rule	$P(EF) = P(E F)P(F)$	$P(EF G) = P(E FG)P(F G)$
Bayes Theorem	$P(E F) = \frac{P(F E)P(E)}{P(F)}$	$P(E FG) = \frac{P(F EG)P(E G)}{P(F G)}$

Combining Events

$$P(ABC) = ?$$

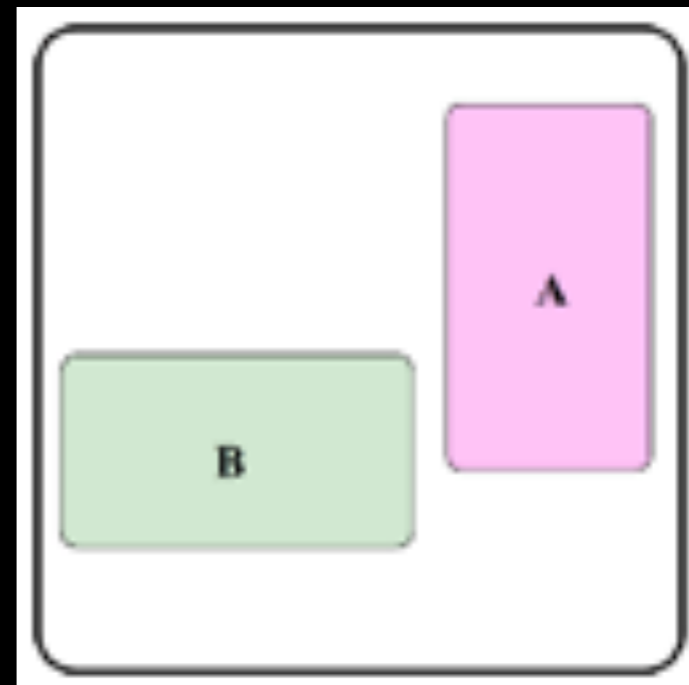
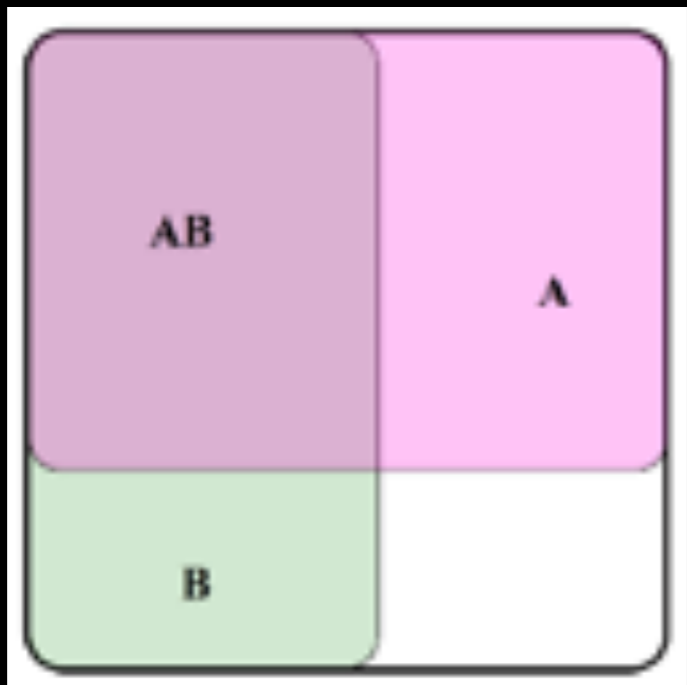
Let $X = AC$:

$$P(ABC) = P(BX) = P(B | X)P(X) = P(B | AC)P(AC)$$

There are three correct ways to apply chain rule to $P(ABC)$!

Independence

Independence	Mutual Exclusion
$P(EF) = P(E)P(F)$	$ E \cap F = 0$
“AND”	“OR”



Independence

Independence	Conditional Independence
$P(EF) = P(E)P(F)$	$P(EF G) = P(E G)P(F G)$ $P(E FG) = P(E G)$
“AND”	“AND [if]”

If E and F are independent.....

.....that does not mean they'll be independent if another event happens!

& vice versa

CS 109

topics

machine learning

sampling, making conclusions from data

multivariate distributions

random variables / distributions

discrete RVs

continuous RVs

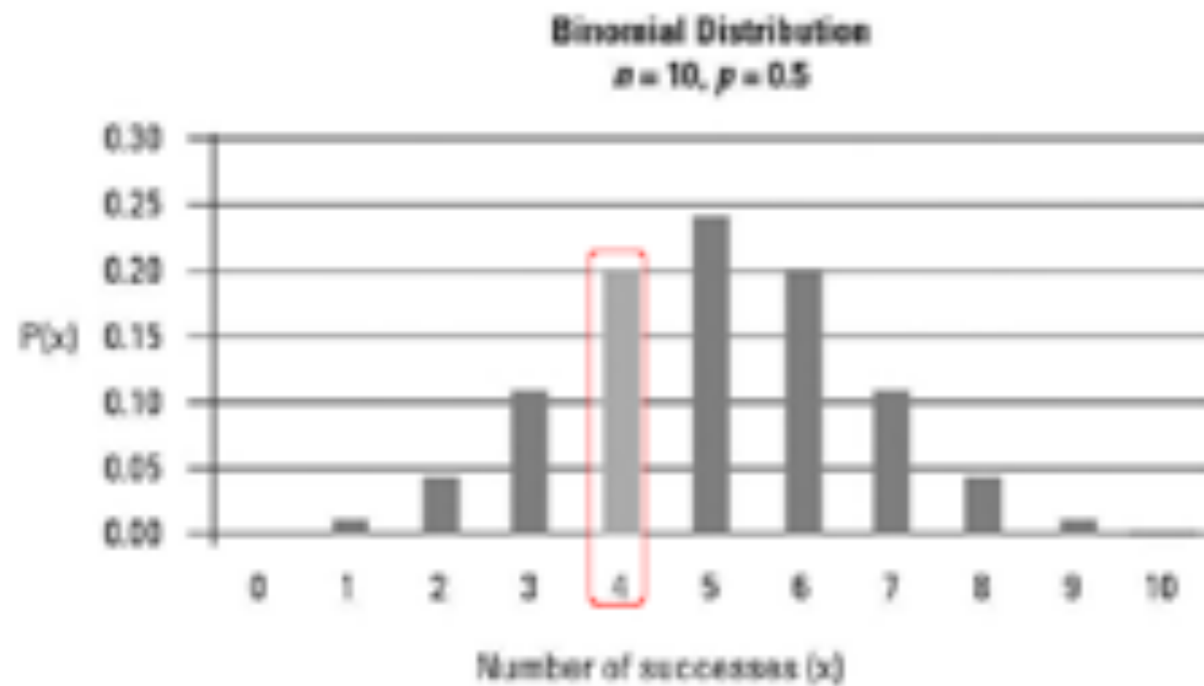
core probability fundamentals

properties of RVs

Probability Distributions

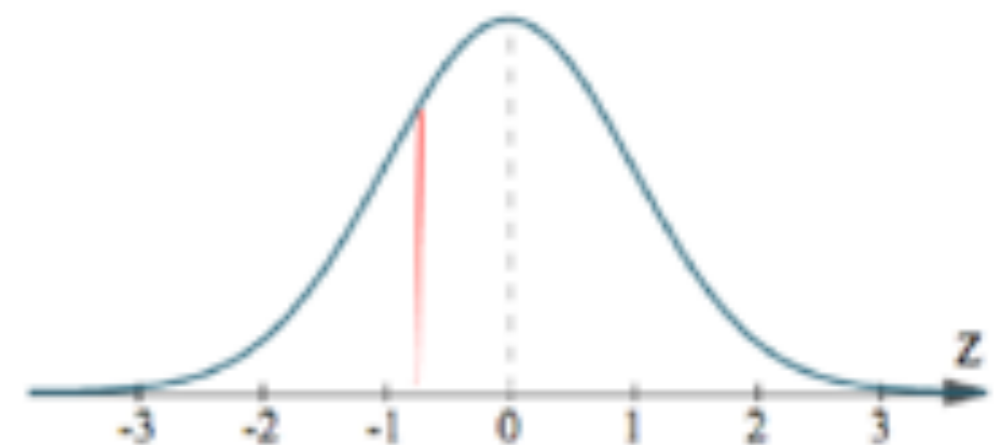
Discrete

PMF:



Continuous

PDF:

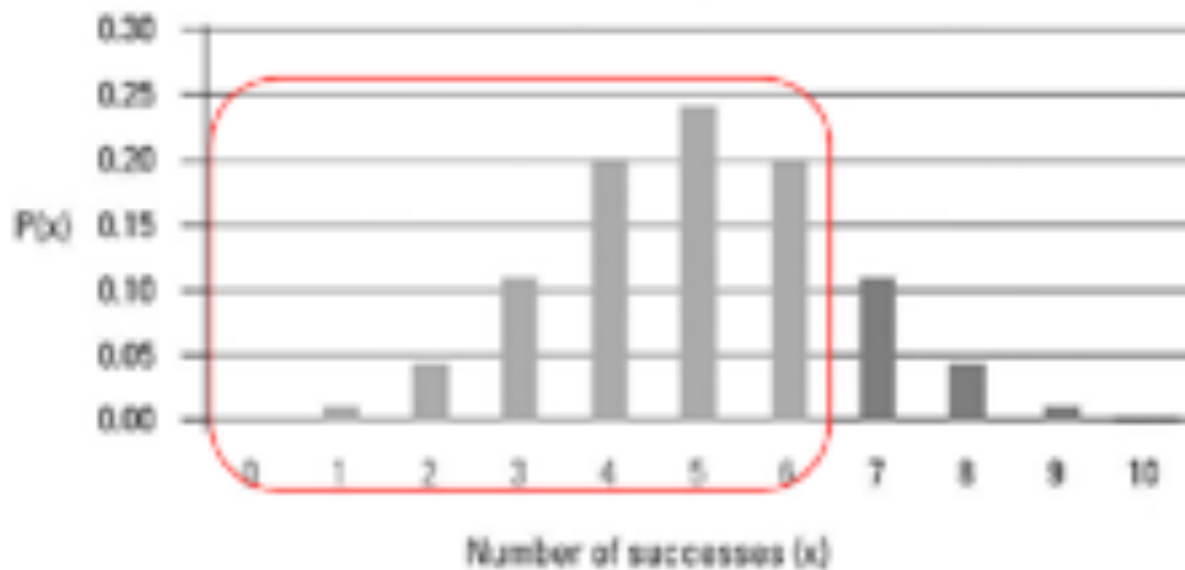


Probability Distributions

Discrete

CDF:

Binomial Distribution
 $n = 10, p = 0.5$



Continuous

CDF:



Expectation & Variance

Discrete definition

$$E[X] = \sum_{x:P(x)>0} x * P(x)$$

Continuous definition

$$E[X] = \int_x x * f_X(x) dx$$

Properties of Expectation

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

$$E[aX + b] = aE[X] + b$$

$$E[g(X)] = \sum_x g(x) * p_X(x)$$

Properties of Variance

$$Var(X) = E[X^2] - E[X]^2$$

$$Var(aX + b) = a^2 Var(X)$$

If X and Y are independent:

$$Var(X + Y) = Var(X) + Var(Y)$$

All our (discrete) friends

Ber(p)	Bin(n, p)	Poi(λ)	Geo(p)	NegBin (r, p)
$P(X) = p$	$\binom{n}{k} p^k (1-p)^{n-k}$	$\frac{\lambda^k e^{-\lambda}}{k!}$	$(1-p)^{k-1} p$	$\binom{k-1}{r-1} p^r (1-p)^{k-r}$
$E[X] = p$	$E[X] = np$	$E[X] = \lambda$	$E[X] = 1/p$	$E[X] = r/p$
$\text{Var}(X) = p(1-p)$	$\text{Var}(X) = np(1-p)$	$\text{Var}(X) = \lambda$	$\frac{1-p}{p^2}$	$\frac{r(1-p)}{p^2}$
Getting candy or not at a random house	# houses out of 20 that give out candy	# houses in an hour that give out candy	# houses to visit before getting candy	# houses to visit before getting candy 3 times

All our (continuous) friends

Uni(α, β)	Exp(λ)	N(μ, σ)
$f(x) = \frac{1}{\beta - \alpha}$	$f(x) = \lambda e^{-\lambda x}$	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
$P(a \leq X \leq b) = \frac{b - a}{\beta - \alpha}$	$F(x) = 1 - e^{-\lambda x}$	$F(x) = \Phi\left(\frac{x - \mu}{\sigma}\right)$
$E(x) = \frac{\alpha + \beta}{2}$	$E[x] = 1 / \lambda$	$E[x] = \mu$
$Var(x) = \frac{(\beta - \alpha)^2}{12}$	$Var(x) = \frac{1}{\lambda^2}$	$Var(x) = \sigma^2$
thickness of sidewalk pavement between houses	time until feet get too sore to trick or treat	weight of filled candy baskets

Discrete vs Continuous

Discrete

$$E[X] = \sum_{x:P(x)>0} x * P(x)$$

$$P(EF) = P(E|F)P(F)$$

$$P(A) = \sum_{b \in S} P(A|B=b)P(B=b)$$

$$P(E|F) = \frac{P(F|E)P(E)}{P(F)}$$

Continuous

$$E[X] = \int_x x * f_X(x) dx$$

$$f(E=e, F=f) = f(E=e|F=f)f(F=f)$$

$$P(A) = \int_b P(A|B=b)f_B(b) db$$

$$f(E=e|F) = \frac{P(F|E=e)f(E=e)}{P(F)}$$

Approximations

When can we approximate a binomial?

n is large

Binomial

```
graph TD; Binomial --> Normal; Binomial --> Poisson;
```

Normal

p is moderate

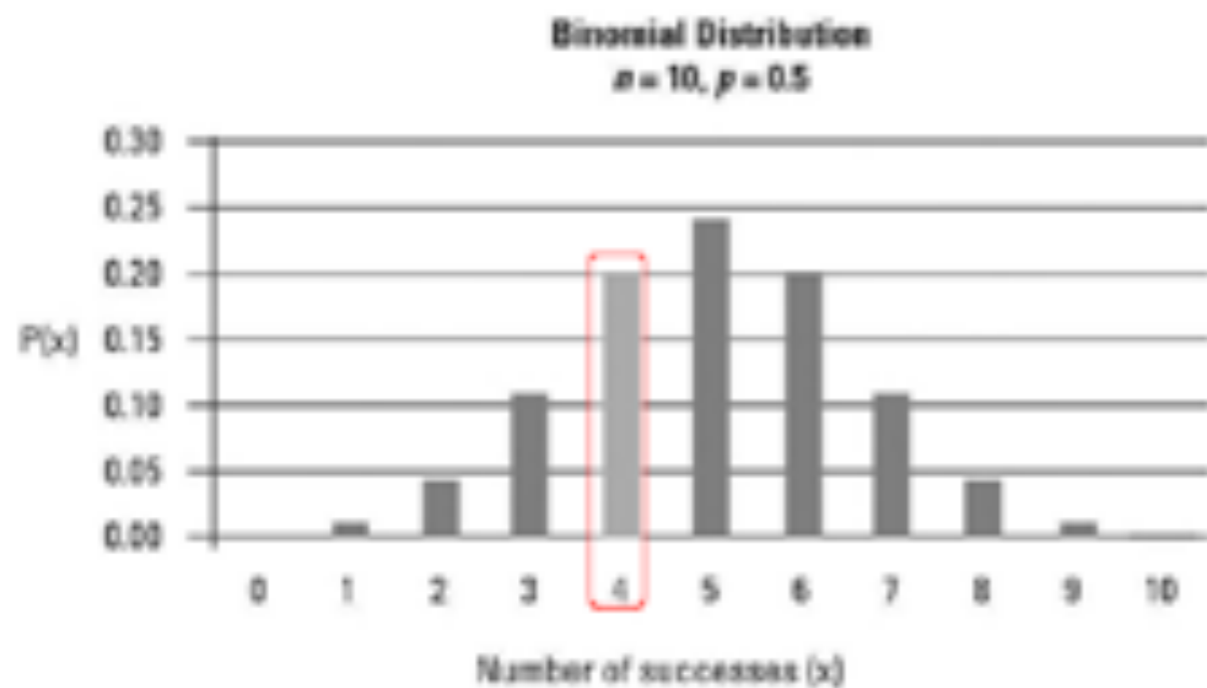
Poisson

p is small

Continuity Correction

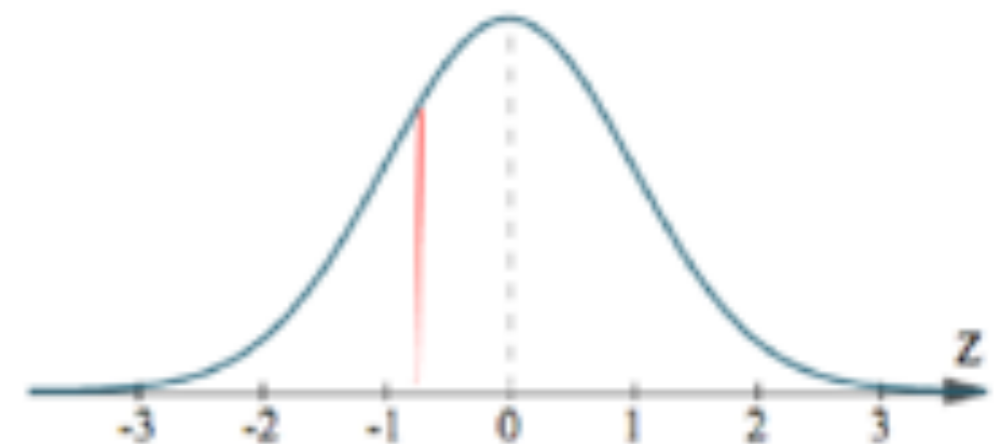
Discrete

PMF:



Continuous

PDF:



Only applies to PDF - why?

Joint Distributions

- Discrete case:

$$p_{x,y}(a, b) = P(X = a, Y = b)$$

$$P_x(a) = \sum_y P_{x,y}(a, y)$$

- Continuous case:

$$P(a_1 < x \leq a_2, b_1 < y \leq b_2) = \int_{a_1}^{a_2} \int_{b_1}^{b_2} f_{X,Y}(x, y) dy dx$$

$$f_X(a) = \int_{-\infty}^{\infty} f_{X,Y}(a, y) dy$$

Joint Distributions

- Discrete case:

$$p_{x,y}(a, b) = P(X = a, Y = b)$$

$$P_x(a) = \sum_y P_{x,y}(a, y)$$

- Continuous case:

$$P(a_1 < x \leq a_2, b_1 < y \leq b_2) = \int_{a_1}^{a_2} \int_{b_1}^{b_2} f_{X,Y}(x, y) dy dx$$

$$f_X(a) = \int_{-\infty}^{\infty} f_{X,Y}(a, y) dy$$

This is just marginalization!

Convolutions

$$X \sim \text{Bin}(n_1, p), Y \sim \text{Bin}(n_2, p) \Rightarrow X + Y \sim \text{Bin}(n_1 + n_2, p)$$

$$X \sim \text{Poi}(\lambda_1), Y \sim \text{Poi}(\lambda_2) \Rightarrow X + Y \sim \text{Poi}(\lambda_1 + \lambda_2)$$

$$X \sim N(\mu_1, \sigma_1^2), Y \sim N(\mu_2, \sigma_2^2) \Rightarrow X + Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

$$f_{X+Y}(a) = \int_{y=-\infty}^{\infty} f_X(a-y)f_Y(y)dy \quad \text{(if X, Y are indep.)}$$

Deriving Convolution

Discrete Case:

$$P(A + B = z) = \sum_b P(A + B = z | B = b)P(B = b)$$

Expanded Law of Total Probability!

Deriving Convolution

Discrete Case:

$$P(A + B = z) = \sum_b P(A + B = z | B = b)P(B = b)$$

$$P(A + B = z) = \sum_b P(A = b - z | B = b)P(B = b)$$

Deriving Convolution

Discrete Case:

$$P(A + B = z) = \sum_b P(A + B = z | B = b)P(B = b)$$

$$P(A + B = z) = \sum_b P(A = b - z | B = b)P(B = b)$$

If A and B are independent:

$$P(A + B = z) = \sum_b P(A = b - z)P(B = b)$$

Deriving Convolution

Discrete Case:

$$P(A + B = z) = \sum_b P(A + B = z | B = b)P(B = b)$$

$$P(A + B = z) = \sum_b P(A = b - z | B = b)P(B = b)$$

If A and B are independent:

$$P(A + B = z) = \sum_b P(A = b - z)P(B = b)$$

Continuous Case:

$$f(A + B = z) = \int_b f(A = b - z | B = b)f(B = b)db$$

Deriving Convolution

Discrete Case:

$$P(A + B = z) = \sum_b P(A + B = z | B = b)P(B = b)$$

$$P(A + B = z) = \sum_b P(A = b - z | B = b)P(B = b)$$

If A and B are independent:

$$P(A + B = z) = \sum_b P(A = b - z)P(B = b)$$

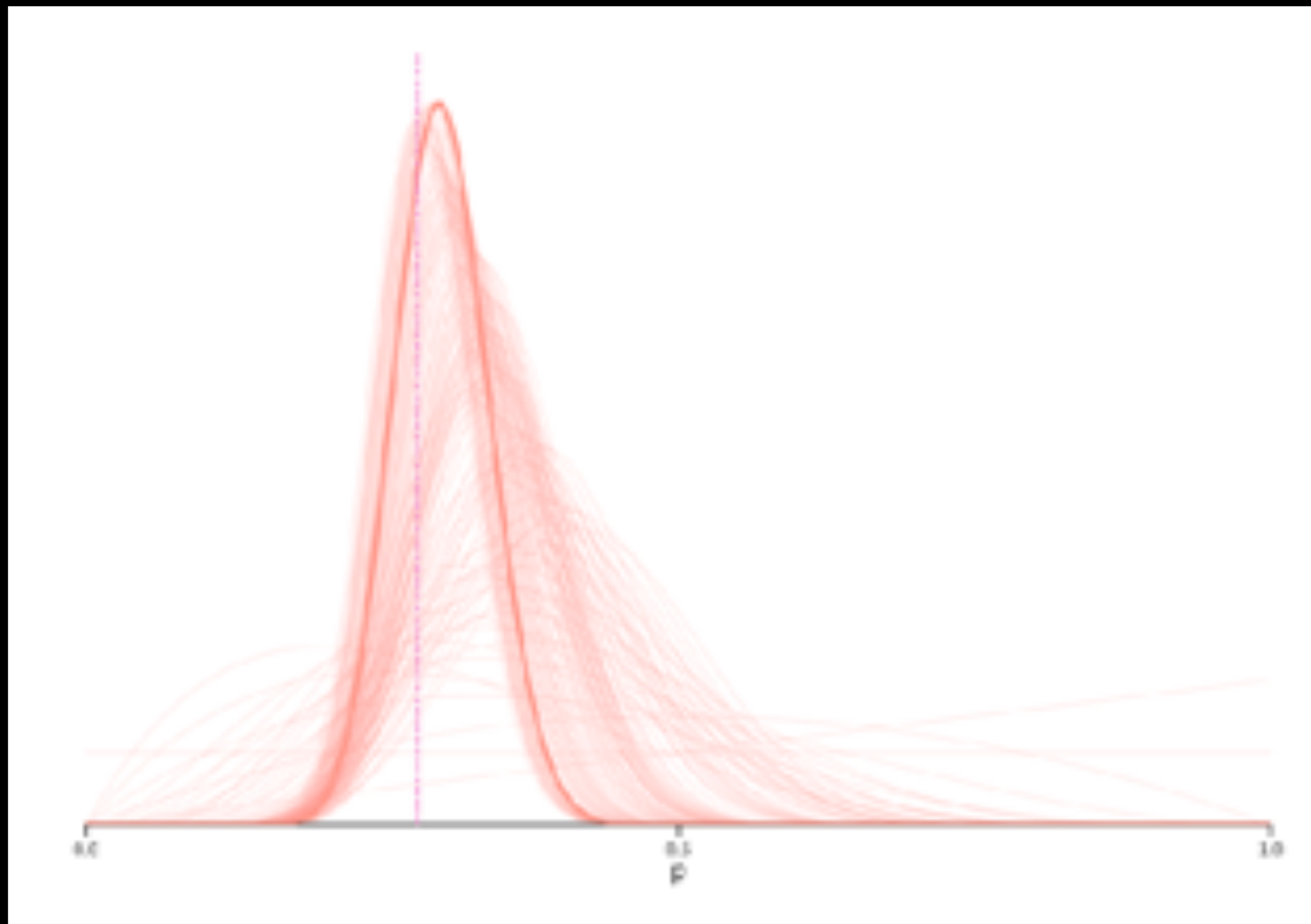
Continuous Case:

If A and B are independent

$$f(A + B = z) = \int_b f(A = b - z | \del{B = b})f(B = b)db$$

Beta

Our first look at the concept of estimating parameters by observing data!



<https://seeing-theory.brown.edu/bayesian-inference/index.html#section3>

CS 109

topics

machine learning

general inference

sampling, making conclusions from data

bootstrapping

CLT

random variables / distributions

unbiased estimators

core probability fundamentals

Sampling From Populations

Challenge: we want to know what the distribution of happiness looks like in Bhutan, but we have limited time and resources and the landscape looks like this:



climb every mountain!



Sampling From Populations



violating data collection norms so that it's unreasonable to assume that a sample is representative of the population

only asking people in Thimphu, e.g.



using statistical methods to draw reasonable conclusions about the population based on data from a random sample



understanding how your results might differ if you sample from the same population multiple times



being an omniscient entity who knows the true population distribution

Taking One Sample

Pick a random sample

if sample size is large enough and sampling methodology is good enough, you can consider it representative of the population!

Take measurements

we have handy equations for the sample mean and sample variance, which are unbiased estimators of the population mean and variance

Report estimate uncertainty

we can use the data from one sample to report our uncertainty about how our estimate of the mean might compare to the true mean (error bars!)

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

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

makes the estimate unbiased

$$\text{Std}(\bar{X}) \approx \sqrt{\left(\frac{S^2}{n}\right)}$$

Sample vs True

True mean and variance are properties of the underlying distribution. They are platonic ideals, completely unattainable!

$$\mu = E[X]$$

$$\sigma^2 = \text{Var}(X)$$

Sample mean and variance are unbiased estimates of true mean and variance based on a single IID sample.

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

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

Variance of sample mean tells us our uncertainty about how good of an estimate sample mean is.

$$\text{Var}(\bar{X}) = \text{Var}\left(\sum_{i=1}^n \frac{X_i}{n}\right) = \sum_{i=1}^n \text{var}\left(\frac{X_i}{n}\right) = \sum_{i=1}^n \frac{\text{Var}(X)}{n^2} \approx \frac{S^2}{n}$$

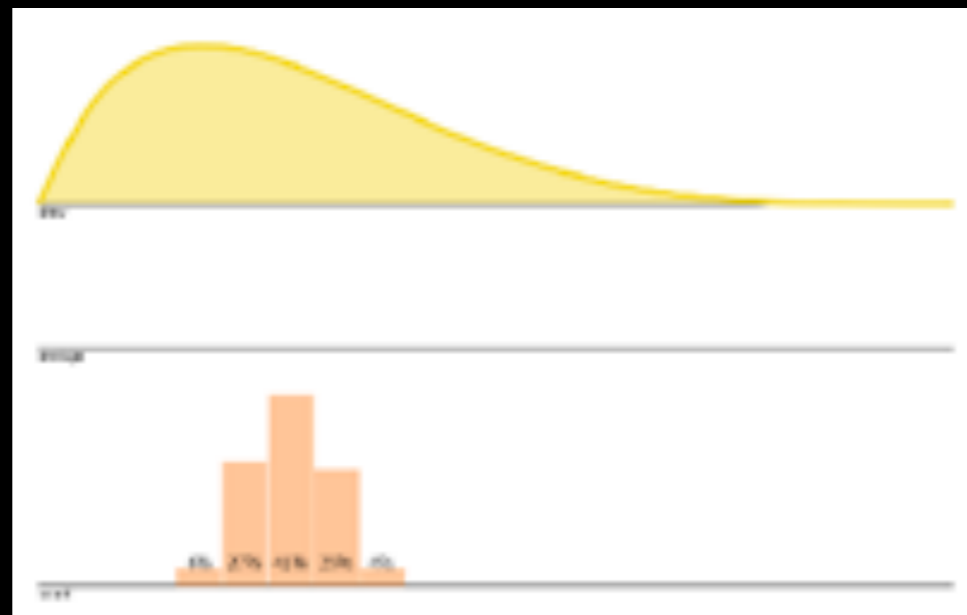
Taking Many Samples

Unbiased Estimators

the expected value of the estimated statistic is the value of the true population statistic (if many samples were to be taken)

Central Limit Theorem

if you sample from the same population a bunch of times, the mean and sum of all your samples (or any IID RVs) will be normally distributed no matter what your distribution looks like!



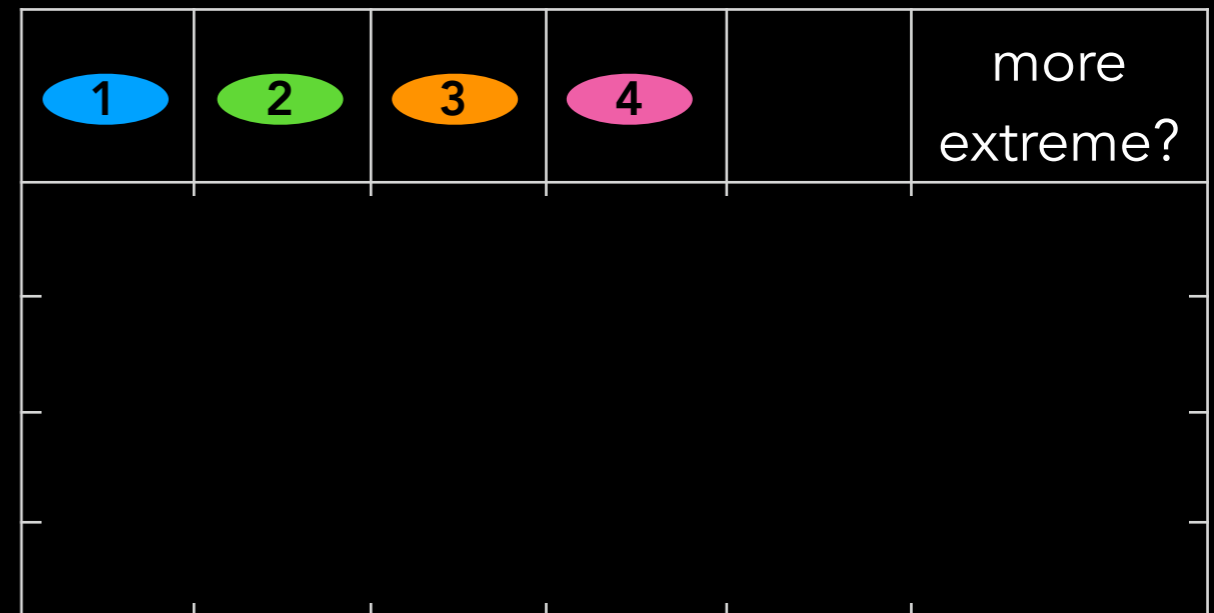
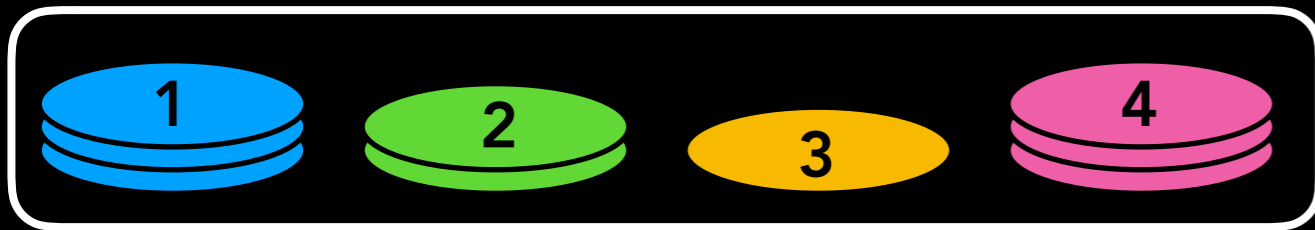
Bootstrapping: Simulating Many Samples From One

challenge

we want to find the probability that the data results we saw were due to chance, but we only have one sample of data

insight

since our sample represents our population, we can sample from the data we have and it's as if we had gone out and collected more



We sample with replacement from our data and calculate our statistic of interest each time, ending up with many estimates for our statistic of interest. We can even use this data to assess whether our observations are due to chance based on our p-value of choice.

General Inference: Sampling from a Bayesian Network to Find Joint Probability



Joint Sampling

generate many "particles" by tracing through the network, generating values for children based on their parents

Calculate Conditional Probability

we can calculate any conditional probability of specific variable assignments by simply counting the particles that match what we're looking for

$$P(\mathbf{X} = \mathbf{a} | \mathbf{Y} = \mathbf{b}) = \frac{N(\mathbf{X} = \mathbf{a}, \mathbf{Y} = \mathbf{b})}{N(\mathbf{Y} = \mathbf{b})}$$



CS 109

topics

machine learning

parameter estimation

classifiers

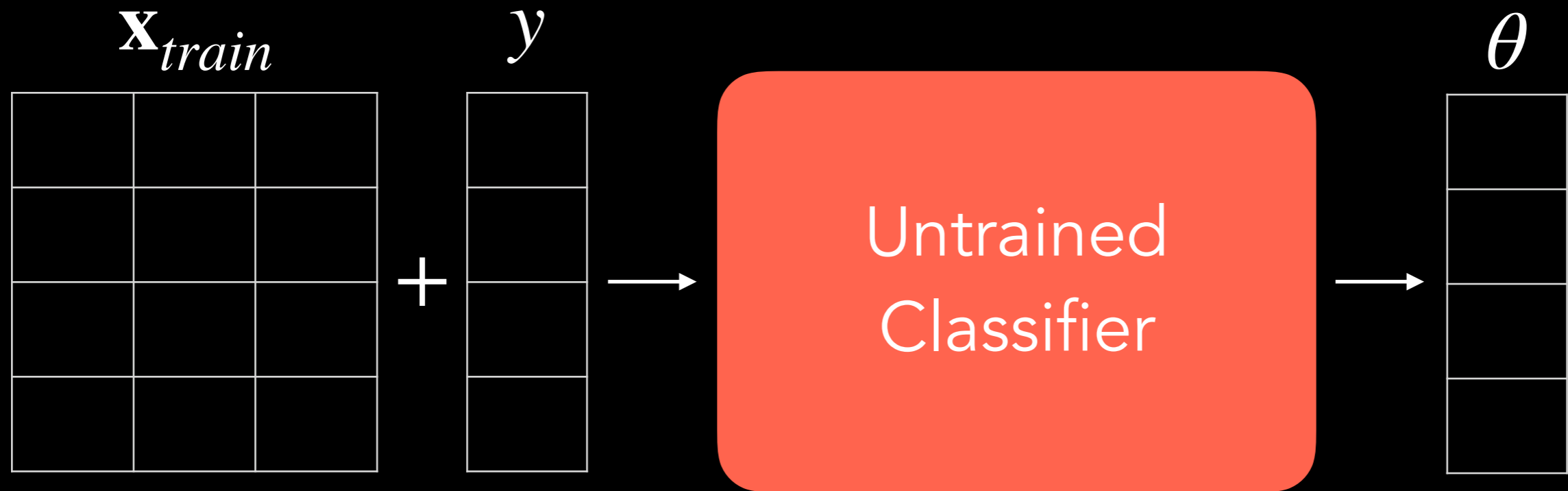
deep learning

sampling, making conclusions from data

random variables / distributions

core probability fundamentals

Classifiers



Parameter Estimation

Maximum Likelihood Estimation

1. Find likelihood: product of likelihoods of each sample/ datapoint given theta
2. Take the log of that expression
3. Take the derivative of that with respect to the parameters
4. Either set to 0 and solve
(if it's a simple case with closed form solution)
or plug into gradient ascent to find a value for theta that maximizes your likelihood

Maximum A Posteriori

1. Find likelihood: product of likelihoods of each sample/ datapoint given theta, times your prior likelihood of that theta
2. - 4. same as above

MLE vs. MAP

MLE:

$$\underset{\theta}{\operatorname{argmax}}(P(\text{data} | \theta)) = \underset{\theta}{\operatorname{argmax}} \left(\prod_{i=1}^n P(x^{(i)} | \theta) \right) = \underset{\theta}{\operatorname{argmax}} \left(\sum_{i=1}^n \log P(x^{(i)} | \theta) \right)$$

MLE vs. MAP

MLE:

$$\underset{\theta}{\operatorname{argmax}}(P(\text{data} | \theta)) = \underset{\theta}{\operatorname{argmax}} \left(\prod_{i=1}^n P(x^{(i)} | \theta) \right) = \underset{\theta}{\operatorname{argmax}} \left(\sum_{i=1}^n \log P(x^{(i)} | \theta) \right)$$

MAP:

$$\underset{\theta}{\operatorname{argmax}}(P(\theta | \text{data})) = \underset{\theta}{\operatorname{argmax}} \left(\frac{P(\text{data} | \theta)P(\theta)}{P(\text{data})} \right)$$

MLE vs. MAP

MLE:

$$\underset{\theta}{\operatorname{argmax}}(P(\text{data} | \theta)) = \underset{\theta}{\operatorname{argmax}} \left(\prod_{i=1}^n P(x^{(i)} | \theta) \right) = \underset{\theta}{\operatorname{argmax}} \left(\sum_{i=1}^n \log P(x^{(i)} | \theta) \right)$$

MAP:

$$\underset{\theta}{\operatorname{argmax}}(P(\theta | \text{data})) = \underset{\theta}{\operatorname{argmax}} \left(\frac{P(\text{data} | \theta)P(\theta)}{\cancel{P(\text{data})}} \right)$$

MLE vs. MAP

MLE:

$$\underset{\theta}{\operatorname{argmax}}(P(\text{data} | \theta)) = \underset{\theta}{\operatorname{argmax}} \left(\prod_{i=1}^n P(x^{(i)} | \theta) \right) = \underset{\theta}{\operatorname{argmax}} \left(\sum_{i=1}^n \log P(x^{(i)} | \theta) \right)$$

MAP:

$$\begin{aligned} \underset{\theta}{\operatorname{argmax}}(P(\theta | \text{data})) &= \underset{\theta}{\operatorname{argmax}} \left(\frac{P(\text{data} | \theta)P(\theta)}{P(\text{data})} \right) = \underset{\theta}{\operatorname{argmax}} \left(P(\theta) \prod_{i=1}^n P(x^{(i)} | \theta) \right) \\ &= \underset{\theta}{\operatorname{argmax}} \left(\log P(\theta) + \sum_{i=1}^n \log P(x^{(i)} | \theta) \right) \end{aligned}$$

Gradient Ascent



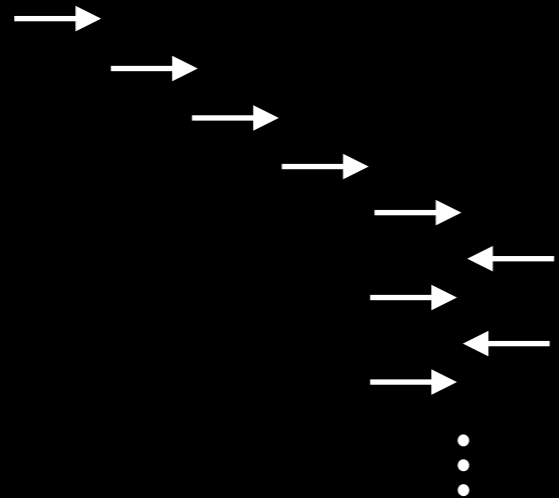
step size

$$\eta = 1$$

step direction

$$= \text{sign} \left[\frac{\partial \text{prob}}{\partial \theta} \right]$$

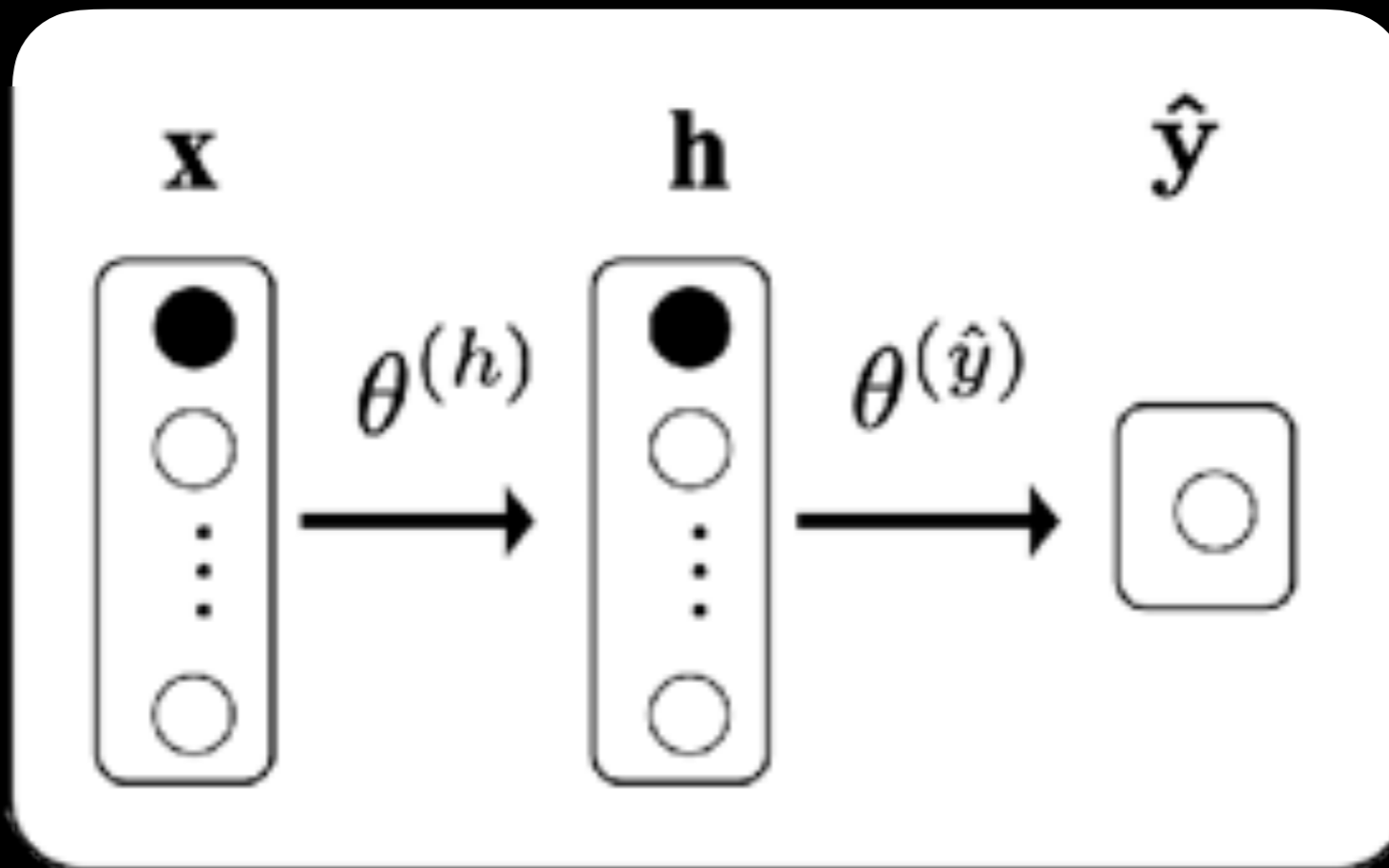
0 1 2 3 4 5 6 7 8 9 10 11 12 θ



Classifier Algorithms

<u>Naïve Bayes</u>	Algorithm	<u>Logistic Regression</u>
All features in \mathbf{x} are conditionally independent given classification	Assumption	Sigmoid gives us the probability of class 1
Whether $y = 0$ or $y = 1$ maximizes the probability of our data	What are we optimizing/figuring out?	The value(s) for θ such that the probability of our data is maximized
Learn (from data) estimates for $\hat{P}(Y = y), \hat{P}(X_i = x_i Y = y)$: $\hat{P}(x_i y) = \frac{(\text{ex. where } X_i = x_i \text{ and } Y = y) + 1}{(\text{ex. where } Y = y) + 2}$ $\hat{P}(Y = y) = \frac{\text{ex. where } Y = y}{\text{total examples}}$	How do we do that mathematically?	Probability of 1 datapoint $P(y \mathbf{x}) = \sigma(\theta^T \mathbf{x})^y \cdot [1 - \sigma(\theta^T \mathbf{x})]^{1-y}$ Use data & gradient ascent to improve thetas $LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log [1 - \sigma(\theta^T \mathbf{x}^{(i)})]$ $\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)}$

Neural Networks



1. Make deep learning assumption: $P(Y = y | \mathbf{X} = \mathbf{x}) = (\hat{y})^y(1 - \hat{y})^{1-y}$
2. Calculate log likelihood for all data: $LL(\theta) = \sum_{i=0}^n y^{(i)}(\log \hat{y}^{(i)}) + (1 - \hat{y}^{(i)}) \log [1 - \hat{y}^{(i)}]$
3. Find partial derivative of LL with respect to each theta:
 use the chain rule!

$$\frac{\partial LL(\theta)}{\partial \theta_j^{(\hat{y})}} = \frac{\partial LL(\theta)}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta_j^{(\hat{y})}} \quad \frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}} = \frac{\partial LL(\theta)}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \mathbf{h}_j} \cdot \frac{\partial \mathbf{h}_j}{\partial \theta_{i,j}^{(h)}}$$

Good luck on the final!

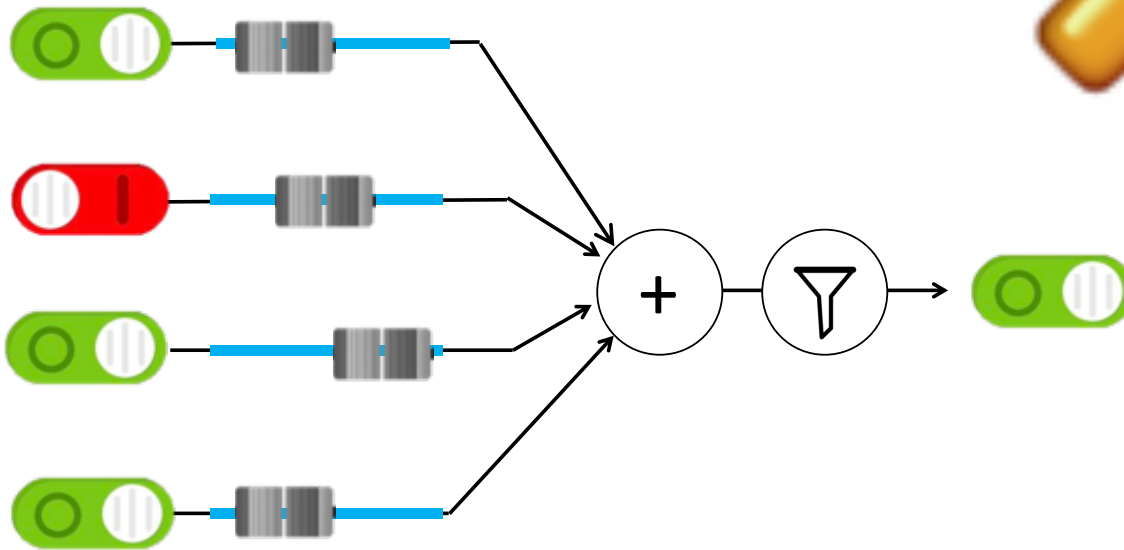
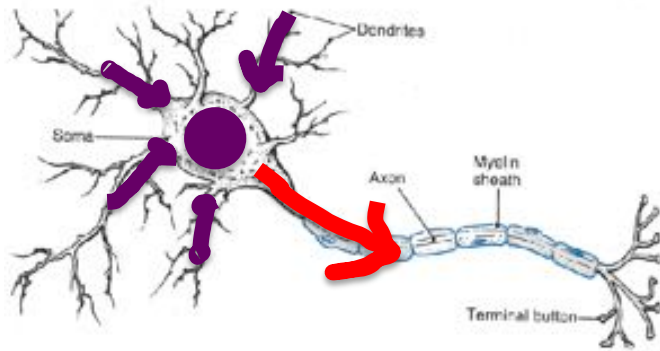




Stretch!

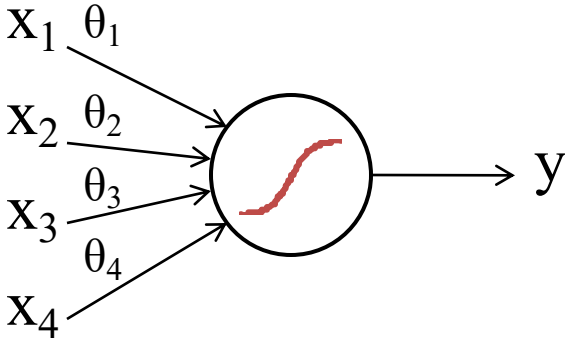
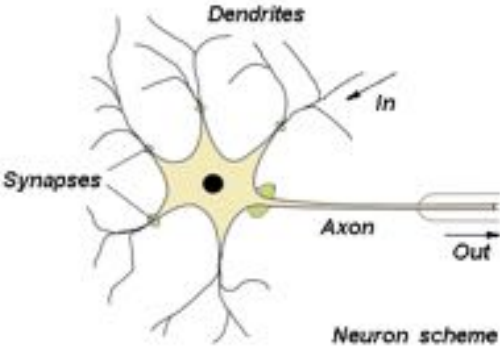
Review

Artificial Neurons

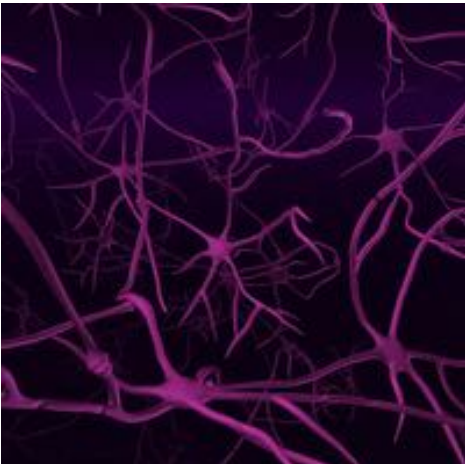


Biological Basis for Neural Networks

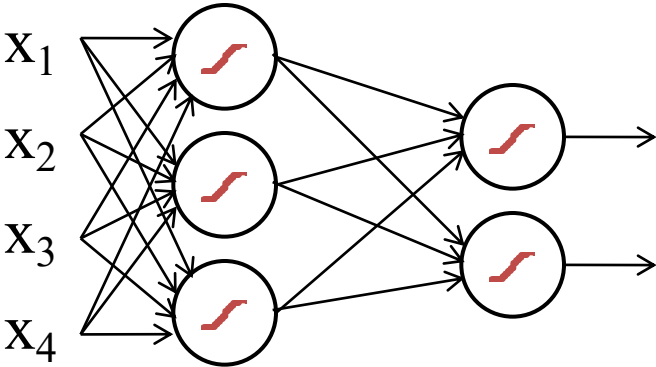
- A neuron



- Your brain



Actually, it's probably someone else's brain

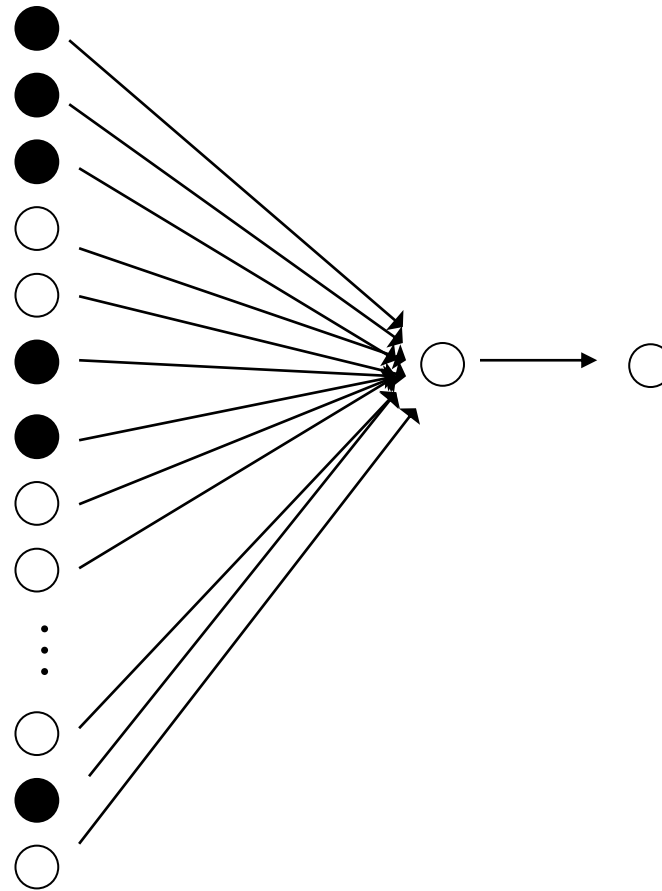
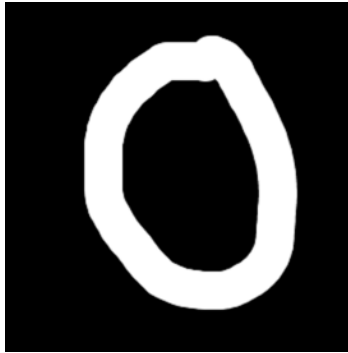


(aka Neural Networks)



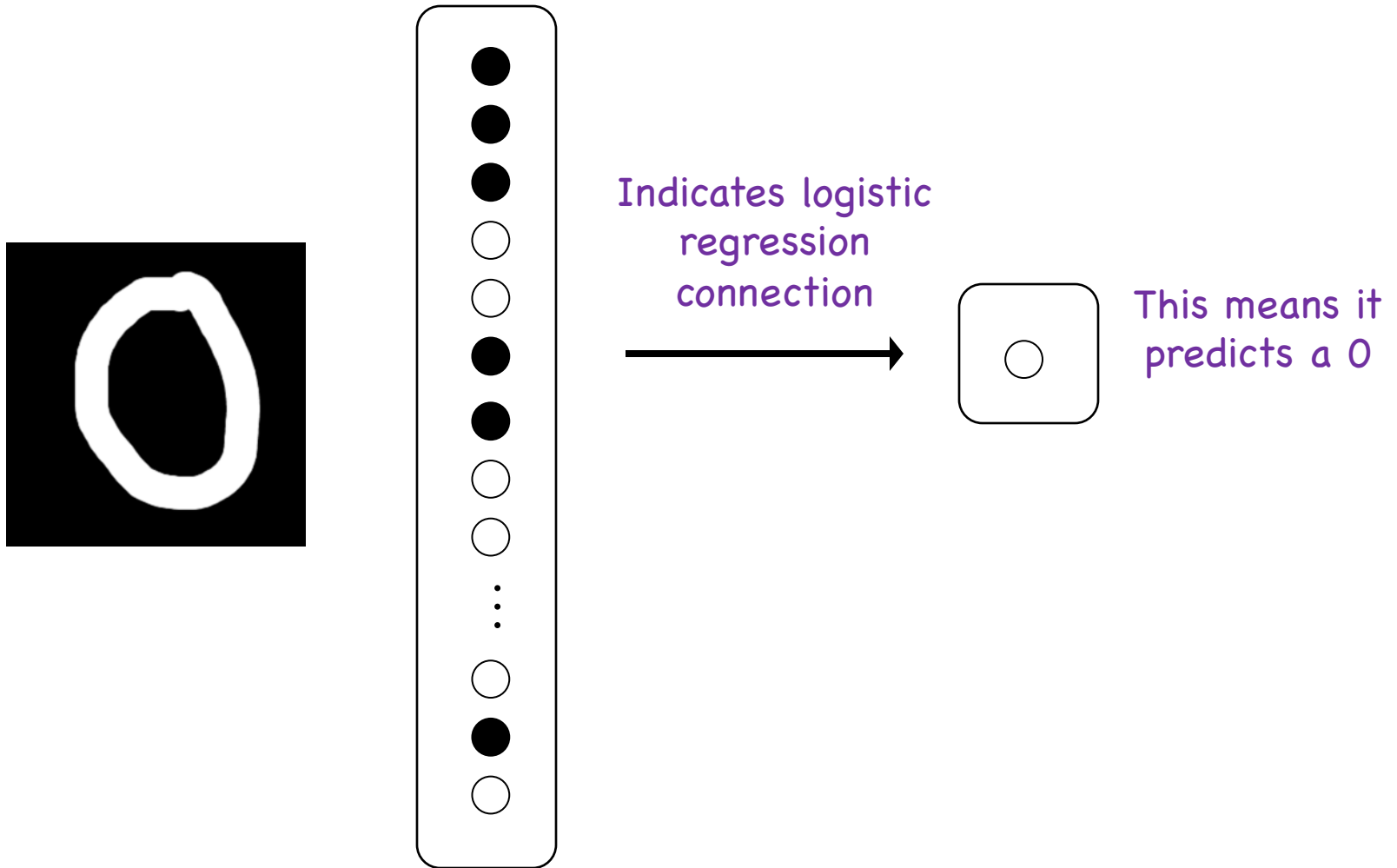
Deep learning is (at its core) many logistic regression pieces stacked on top of each other.

Logistic Regression

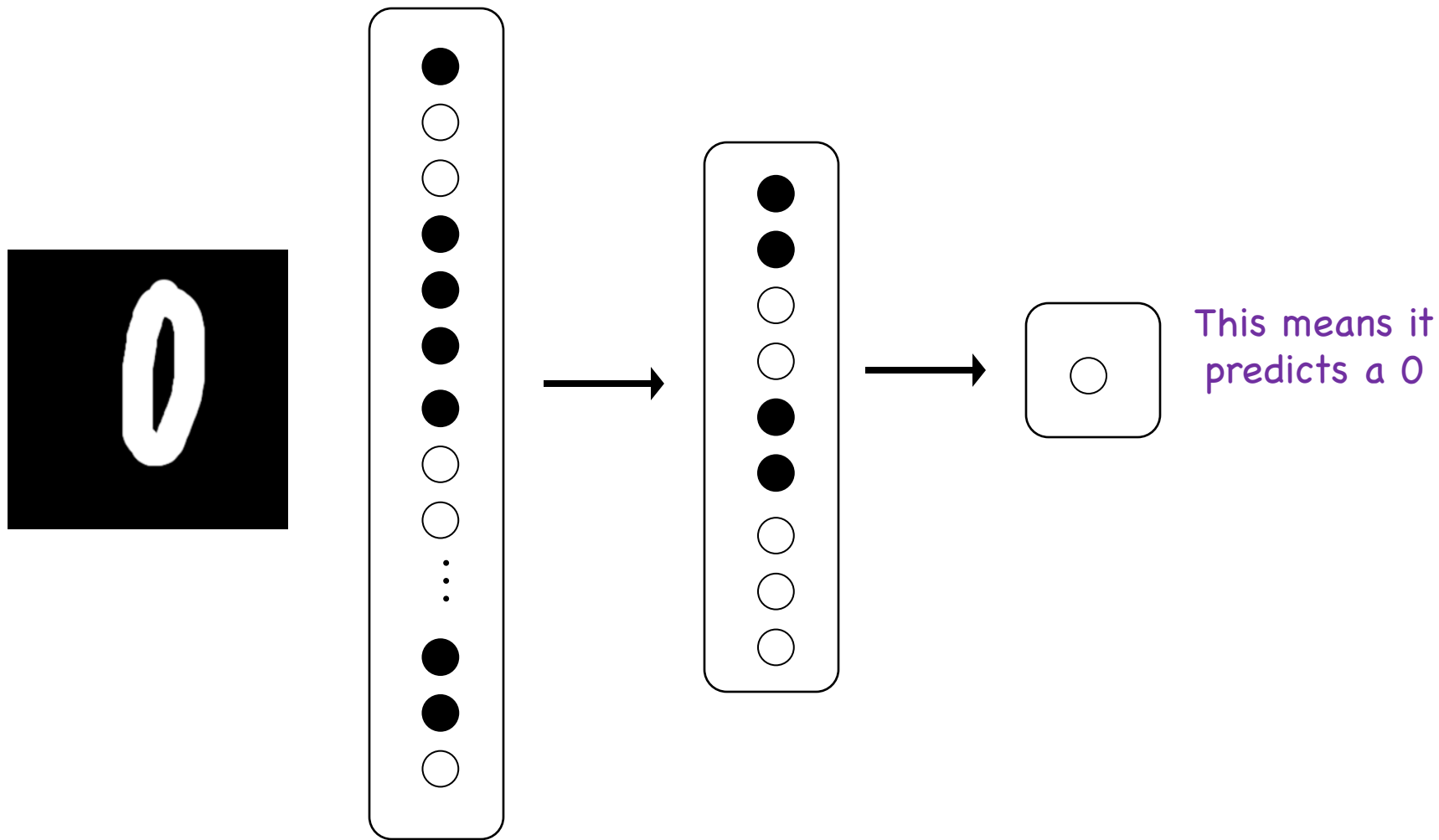


This means it predicts a 0

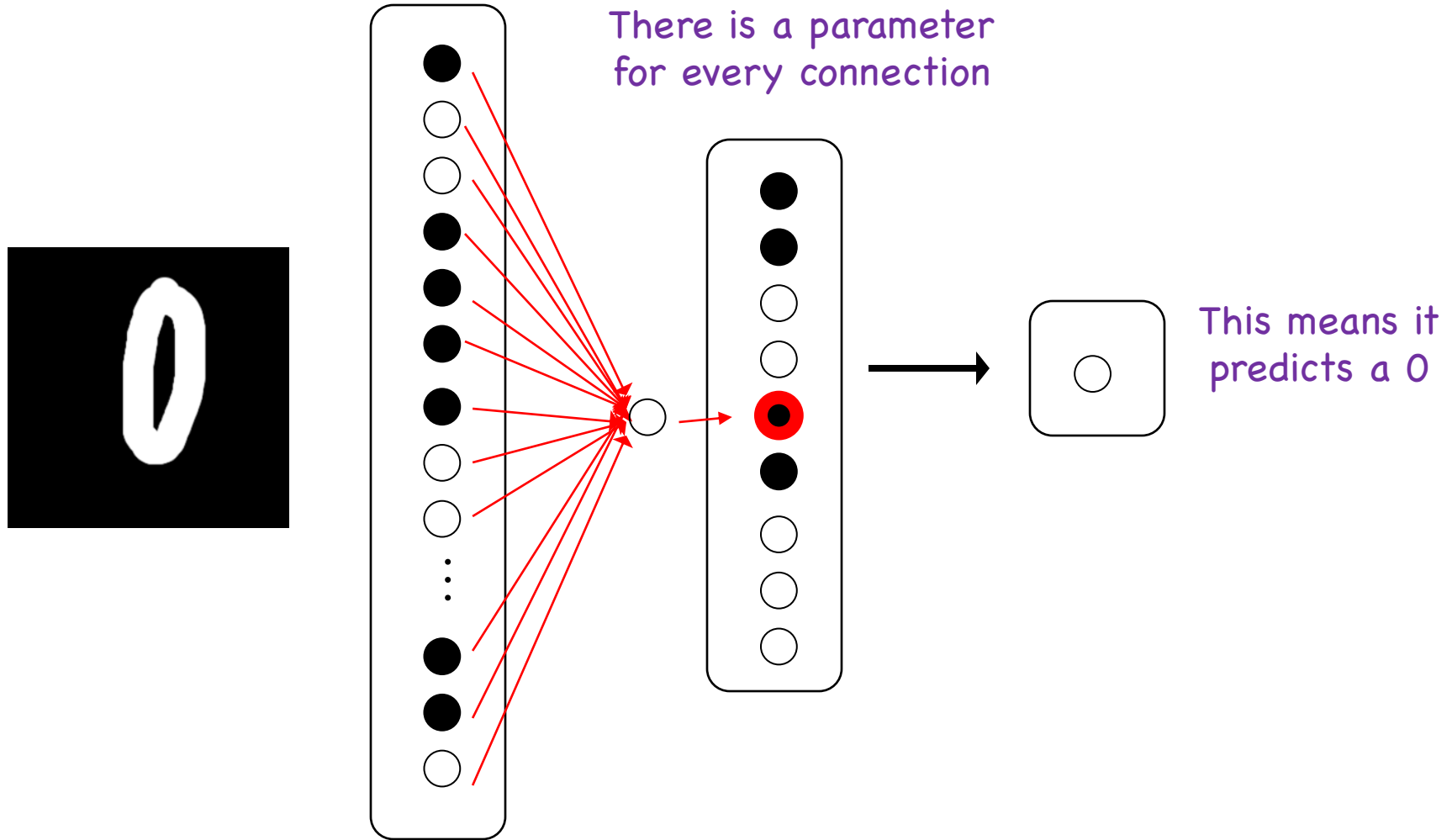
Logistic Regression



We Can Put Neurons Together

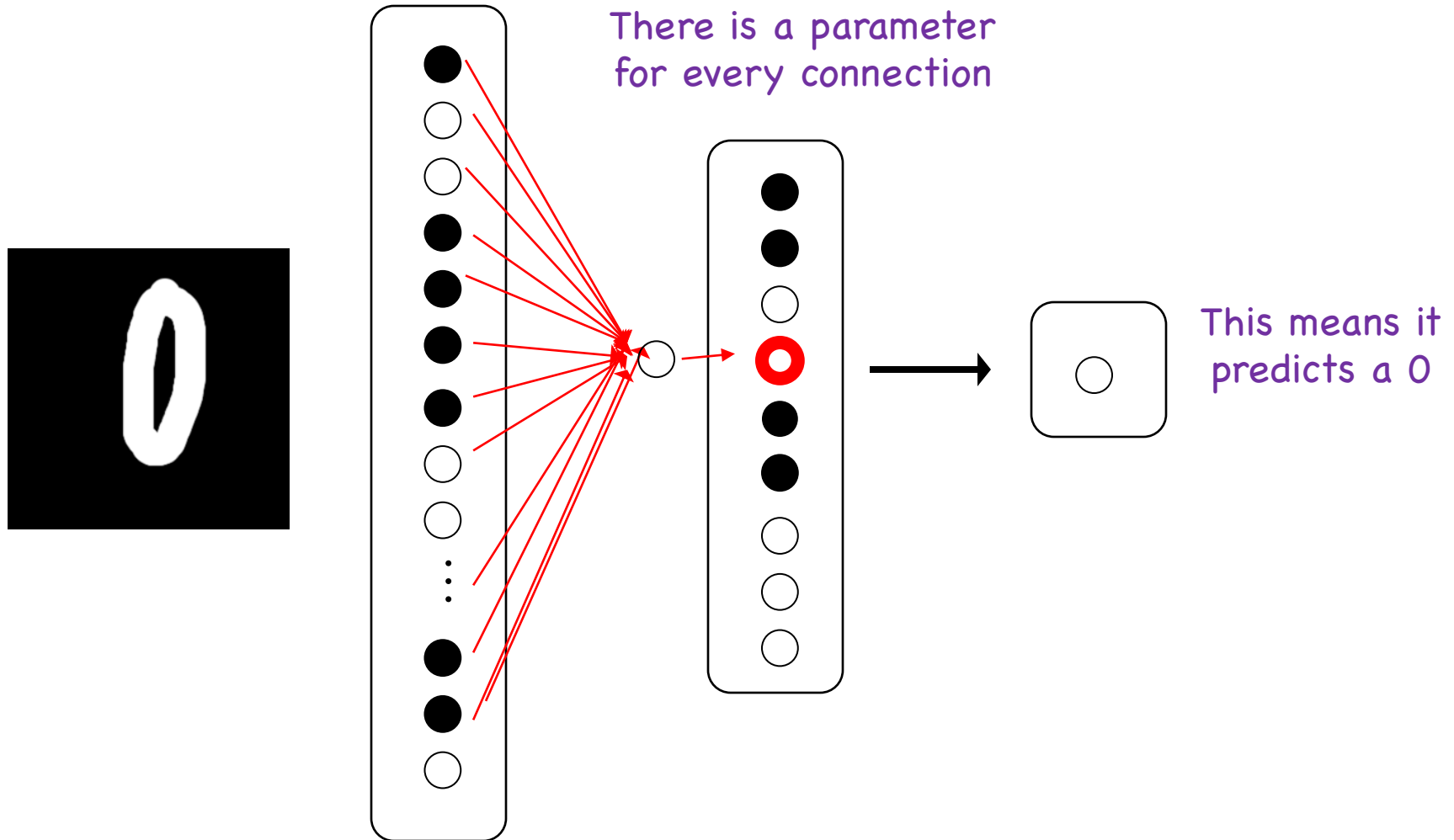


We Can Put Neurons Together



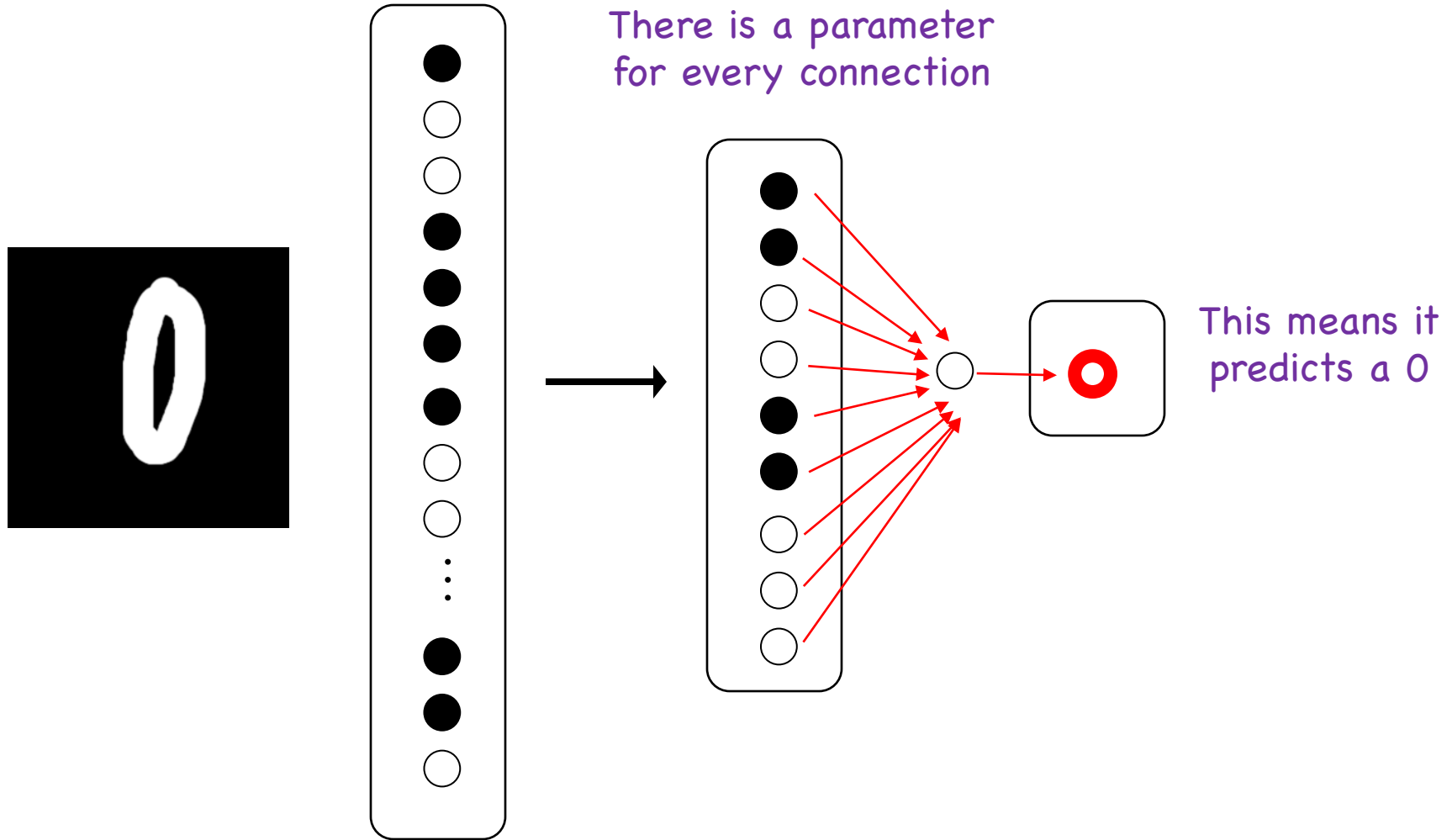
Look at a single “hidden” neuron

We Can Put Neurons Together



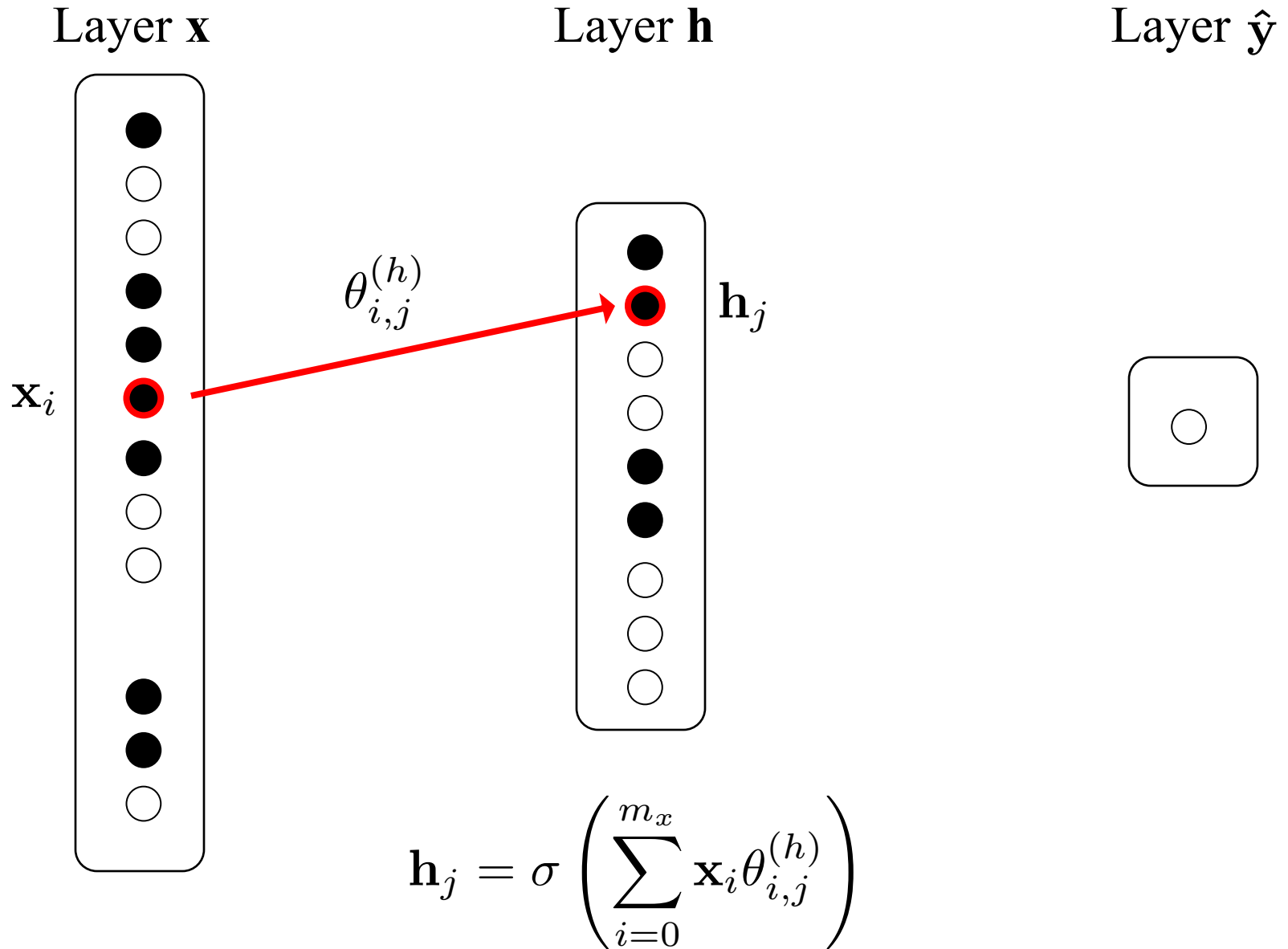
Look at another “hidden” neuron

We Can Put Neurons Together



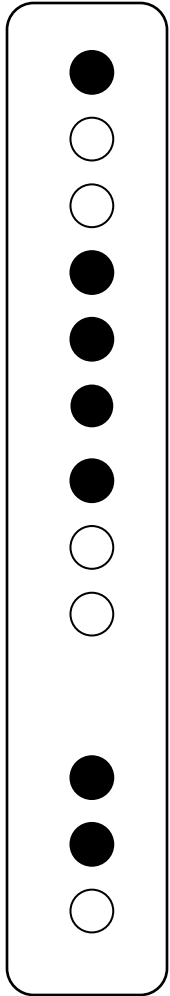
Look at another neuron

New Notation

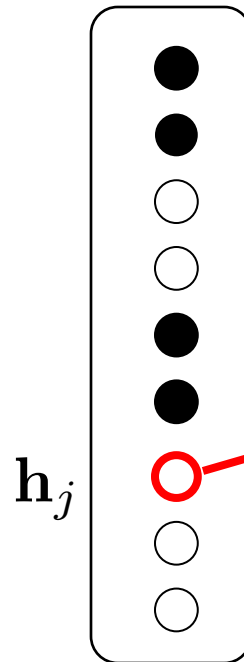


New Notation

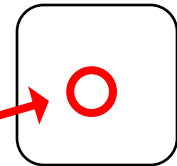
Layer x



Layer h



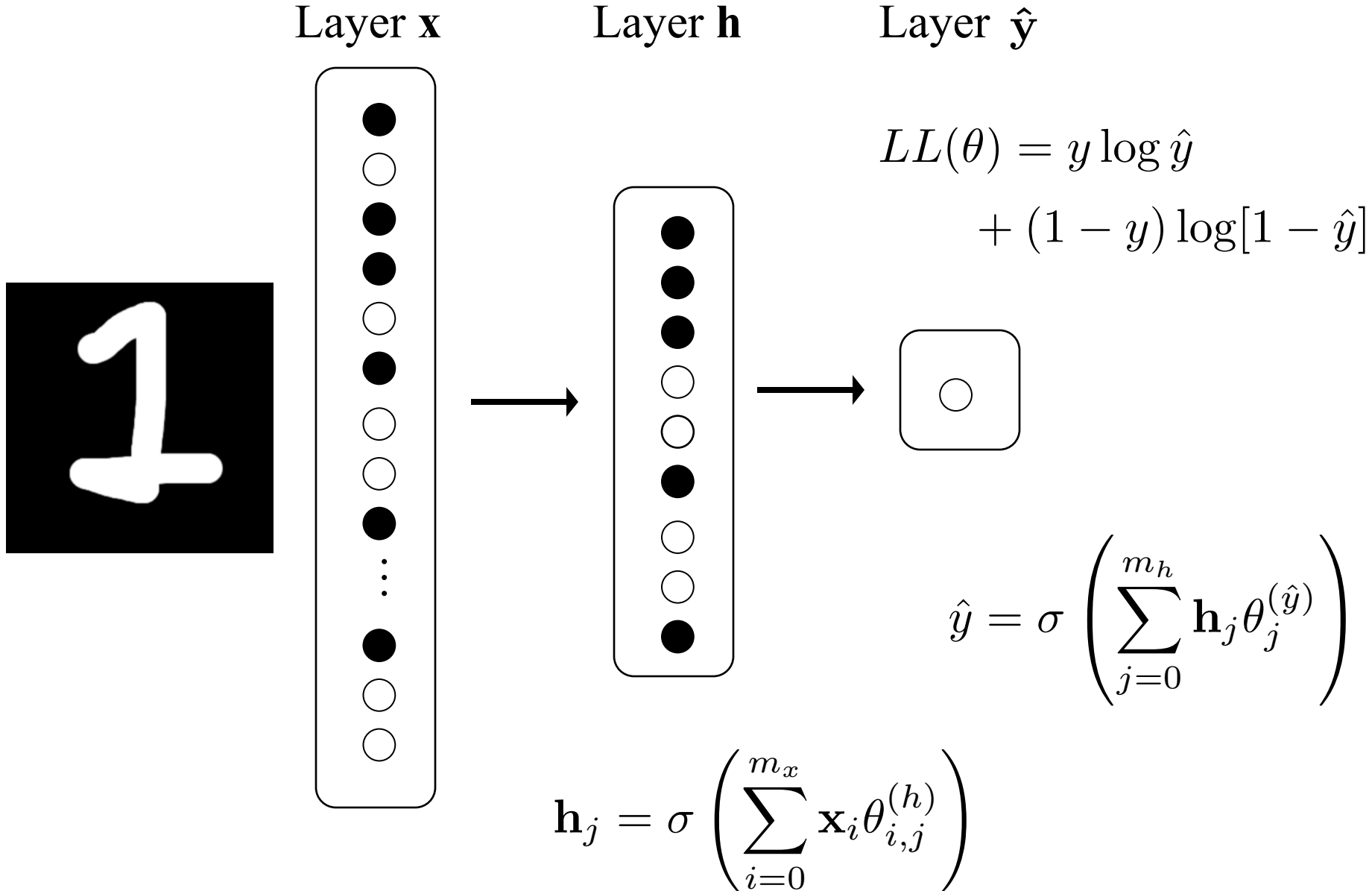
Layer \hat{y}



$\theta_j^{(\hat{y})}$

$$\hat{y} = \sigma \left(\sum_{j=0}^{m_h} h_j \theta_j^{(\hat{y})} \right)$$

Forward Pass



Same Assumption, Same LL

$$P(Y = 1|X = \mathbf{x}) = \hat{y} \quad \hat{y} = \sigma(\theta^T \mathbf{x})$$

For one datum

$$P(Y = y|\mathbf{X} = \mathbf{x}) = (\hat{y})^y (1 - \hat{y})^{1-y} \quad Y \sim \text{Bern}(\hat{y})$$

For IID data

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n P(Y = y^{(i)} | X = \mathbf{x}^{(i)}) \\ &= \prod_{i=1}^n (\hat{y}^{(i)})^{y^{(i)}} \cdot \left[1 - (\hat{y}^{(i)}) \right]^{(1-y^{(i)})} \end{aligned}$$

Take the log

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log[1 - \hat{y}^{(i)}]$$

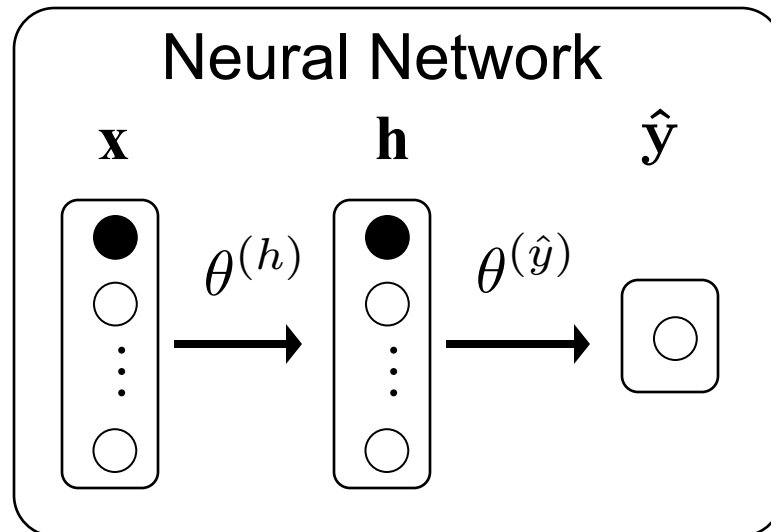
Derivative Goals

Loss with respect to
output layer params

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}}$$

Loss with respect to
hidden layer params

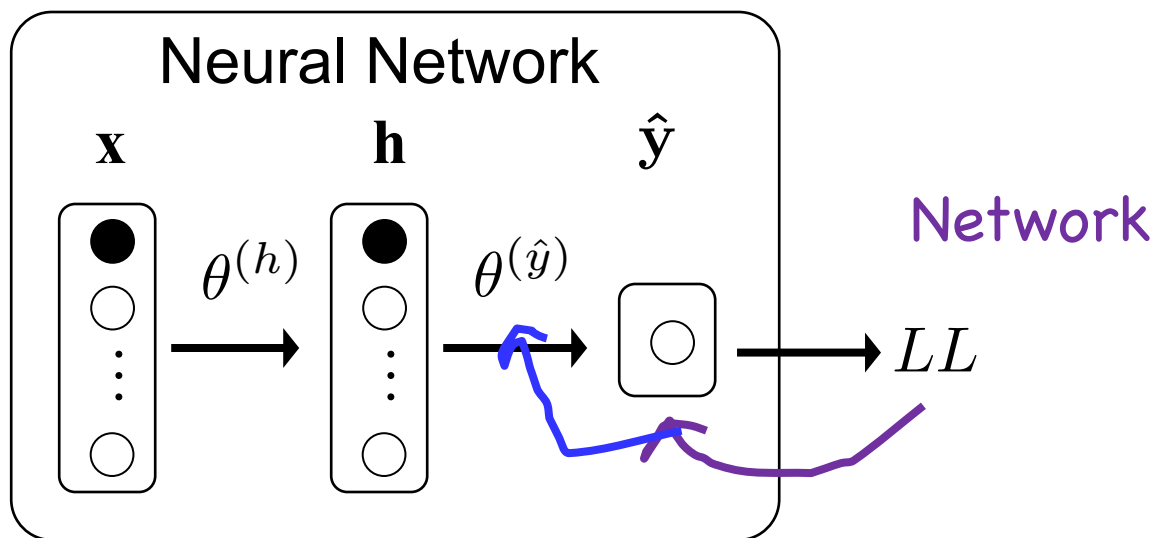
$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}}$$



Chain Rule Example 1

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}}$$

Goal



$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}} = \frac{\partial LL}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta_i^{(\hat{y})}}$$

Decomposition

Make it Simple

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}} = \text{[Yellow Box]} \cdot \text{[Turtle]}$$

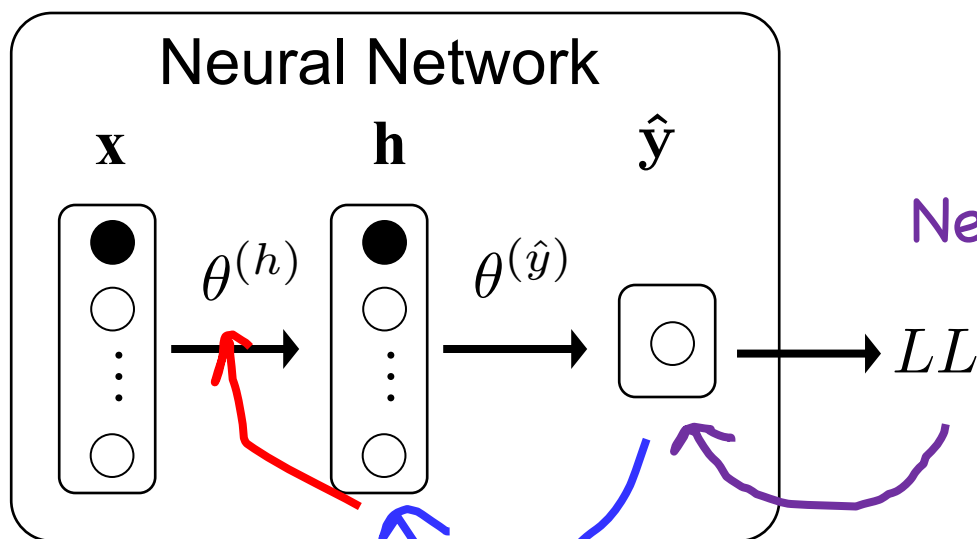
$$\text{[Yellow Box]} = \frac{y}{\hat{y}} - \frac{(1 - y)}{(1 - \hat{y})}$$

$$\text{[Turtle]} = \hat{y}[1 - \hat{y}] \cdot h_i$$

Chain Rule Example 2

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}}$$

Goal



Network

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}} = \frac{\partial LL}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \mathbf{h}_j} \cdot \frac{\partial \mathbf{h}_j}{\partial \theta_{i,j}^{(h)}}$$

Decomposition

Make it Simple

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}} = \begin{array}{|c|c|c|} \hline \img alt="Spike" data-bbox="425 171 533 316"/> & \img alt="Tortoise" data-bbox="535 171 643 316"/> & \img alt="Piranha Plant" data-bbox="645 171 753 316"/> \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline \img alt="Spike" data-bbox="292 358 403 505"/> \\ \hline \end{array} = \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})}$$

$$\begin{array}{|c|} \hline \img alt="Tortoise" data-bbox="292 569 403 716"/> \\ \hline \end{array} = \hat{y}[1-\hat{y}]\theta_j^{(\hat{y})}$$

$$\begin{array}{|c|} \hline \img alt="Piranha Plant" data-bbox="300 758 411 907"/> \\ \hline \end{array} = \mathbf{h}_j[1-\mathbf{h}_j]\mathbf{x}_j$$

End Review



*"True friendship comes when the silences
between two people is comfortable."*

Your random variables are correlated

Covariance and Correlation

Noah Arthurs

CS109, Stanford University

Recall our Ebola Bats



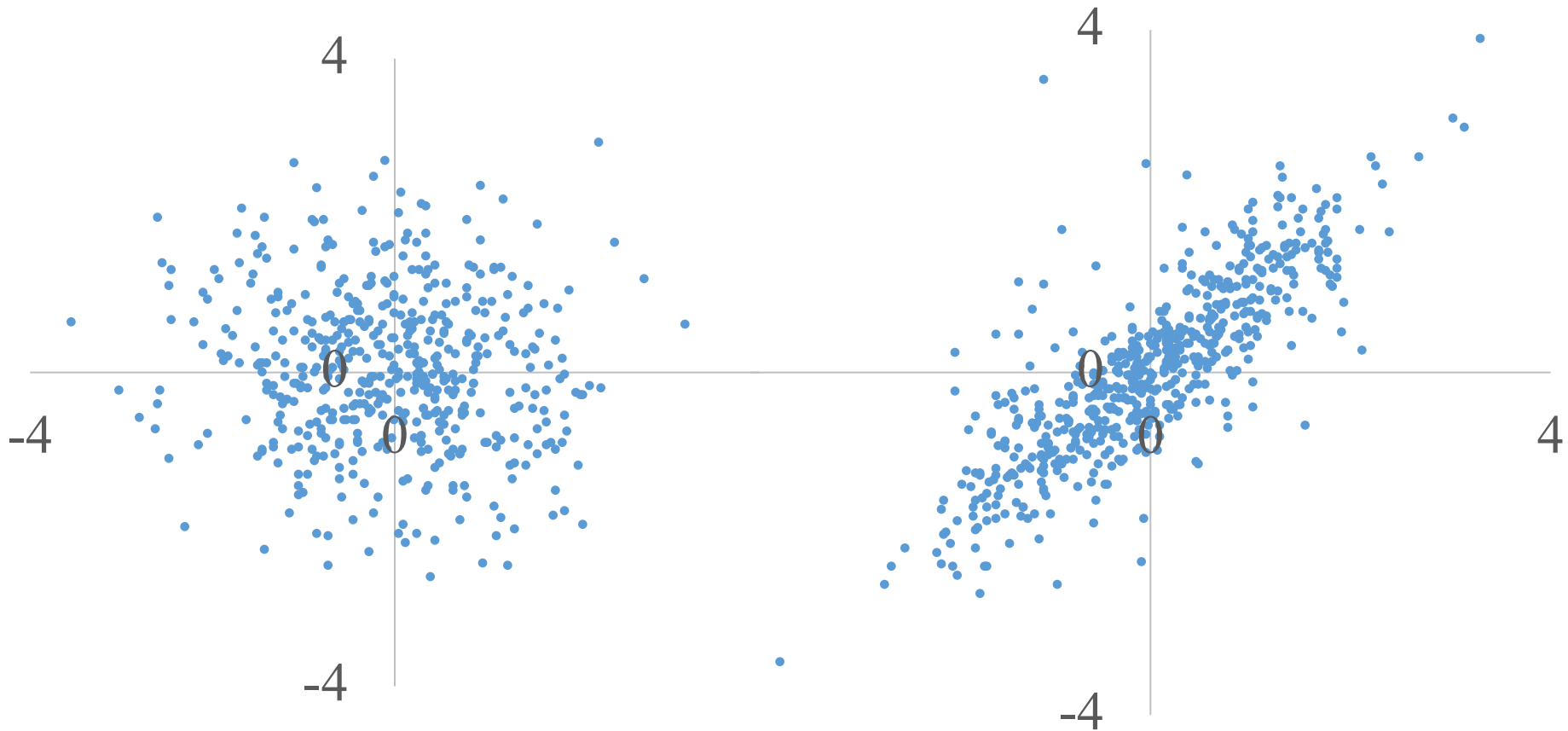
Bat Data

Gene1	Gene2	Gene3	Gene4	Gene5	Trait
TRUE	FALSE	TRUE	TRUE	FALSE	FALSE
FALSE	FALSE	TRUE	TRUE	TRUE	TRUE
TRUE	FALSE	TRUE	FALSE	FALSE	FALSE
TRUE	FALSE	TRUE	TRUE	TRUE	FALSE
FALSE	TRUE	TRUE	TRUE	TRUE	TRUE
FALSE	FALSE	FALSE	TRUE	FALSE	FALSE
TRUE	FALSE	FALSE	TRUE	FALSE	FALSE
TRUE	FALSE	FALSE	TRUE	FALSE	FALSE
TRUE	FALSE	TRUE	FALSE	FALSE	FALSE
FALSE	TRUE	FALSE	TRUE	FALSE	FALSE
TRUE	TRUE	FALSE	TRUE	FALSE	FALSE
TRUE	FALSE	FALSE	TRUE	FALSE	FALSE
TRUE	FALSE	TRUE	TRUE	TRUE	FALSE
FALSE	FALSE	TRUE	TRUE	FALSE	FALSE
TRUE	FALSE	FALSE	TRUE	FALSE	FALSE
TRUE	FALSE	FALSE	TRUE	FALSE	FALSE
			...		
TRUE	FALSE	FALSE	TRUE	FALSE	FALSE

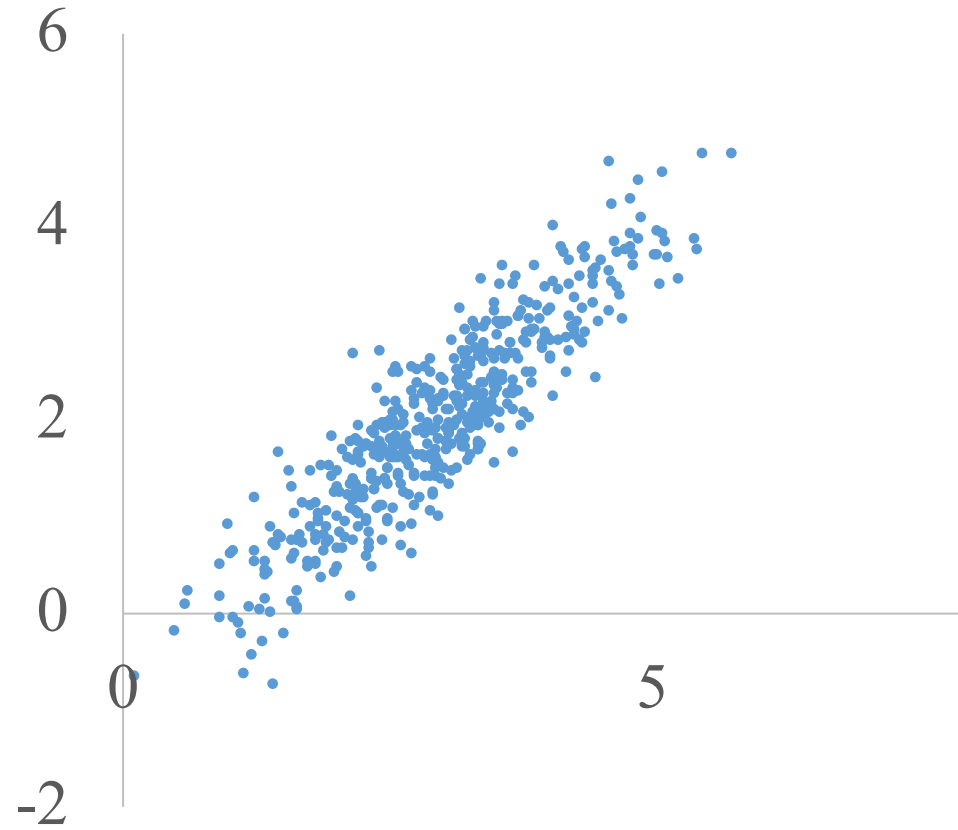
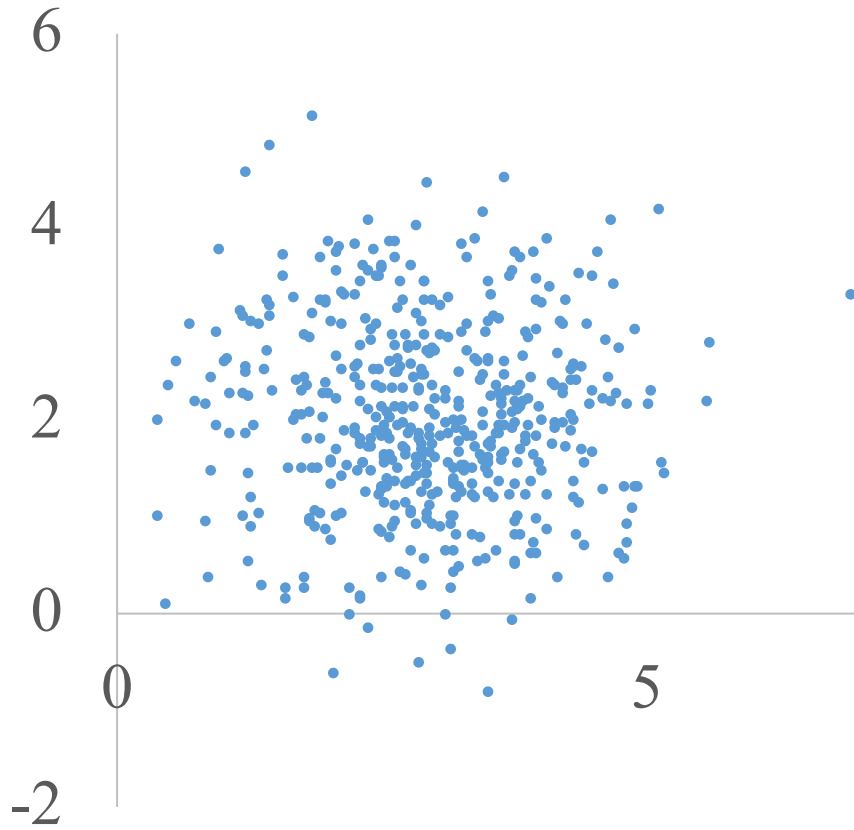
Expression Amount

Gene5	Trait
0.76	0.83
0.94	0.85
0.82	0.03
0.94	0.32
0.50	0.10
0.40	0.53
0.90	0.67
0.29	0.71
0.72	0.25
0.15	0.24
0.79	0.98
0.68	0.77
0.71	0.37
0.36	0.18
0.62	0.08
0.59	0.38
0.82	0.76

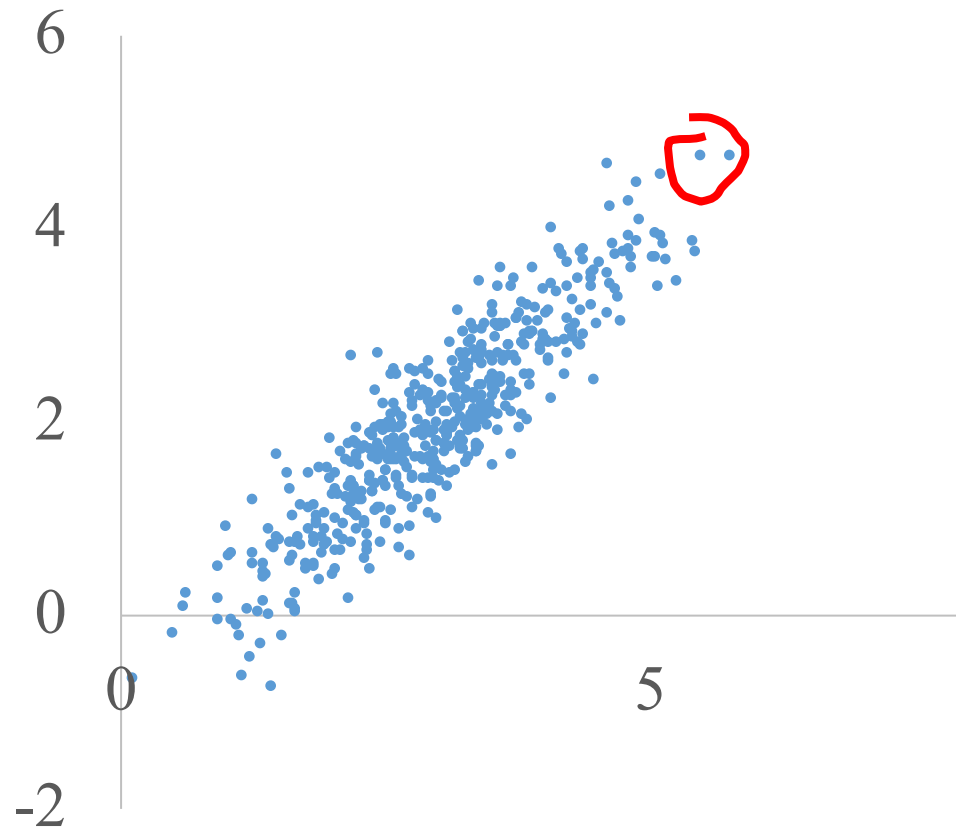
Spot The Difference



Spot The Difference



Vary Together

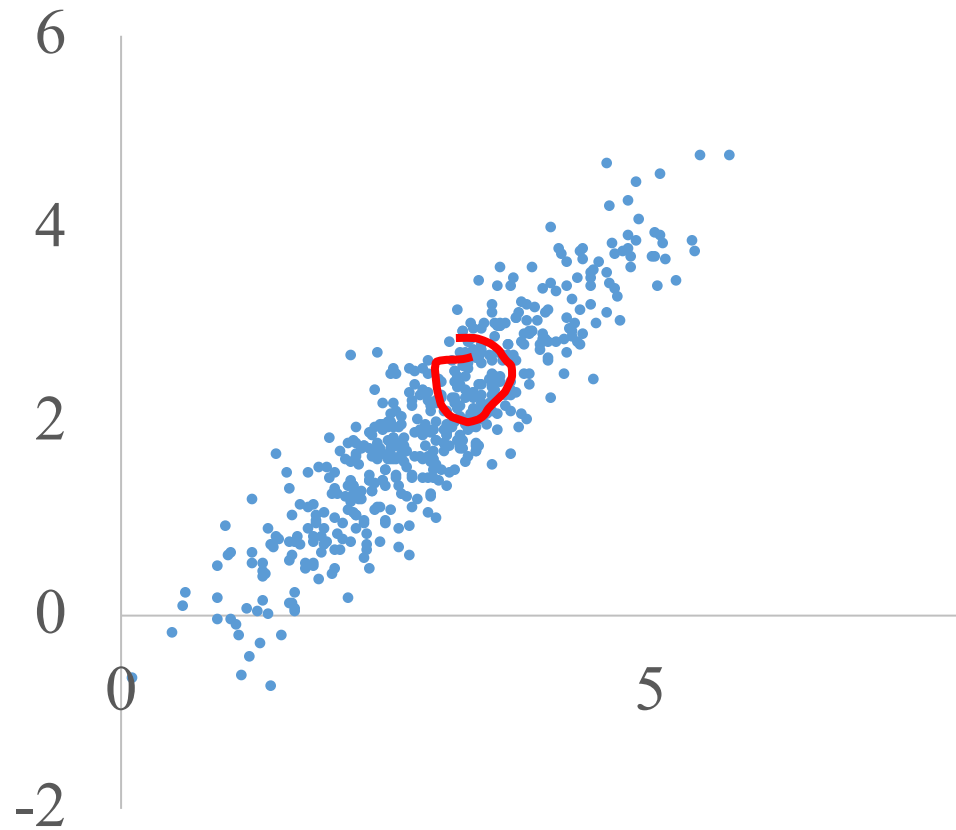


$$x - E[x] = 3$$

$$y - E[y] = 2.6$$

$$(x - E[x])(y - E[y]) = 7.8$$

Vary Together

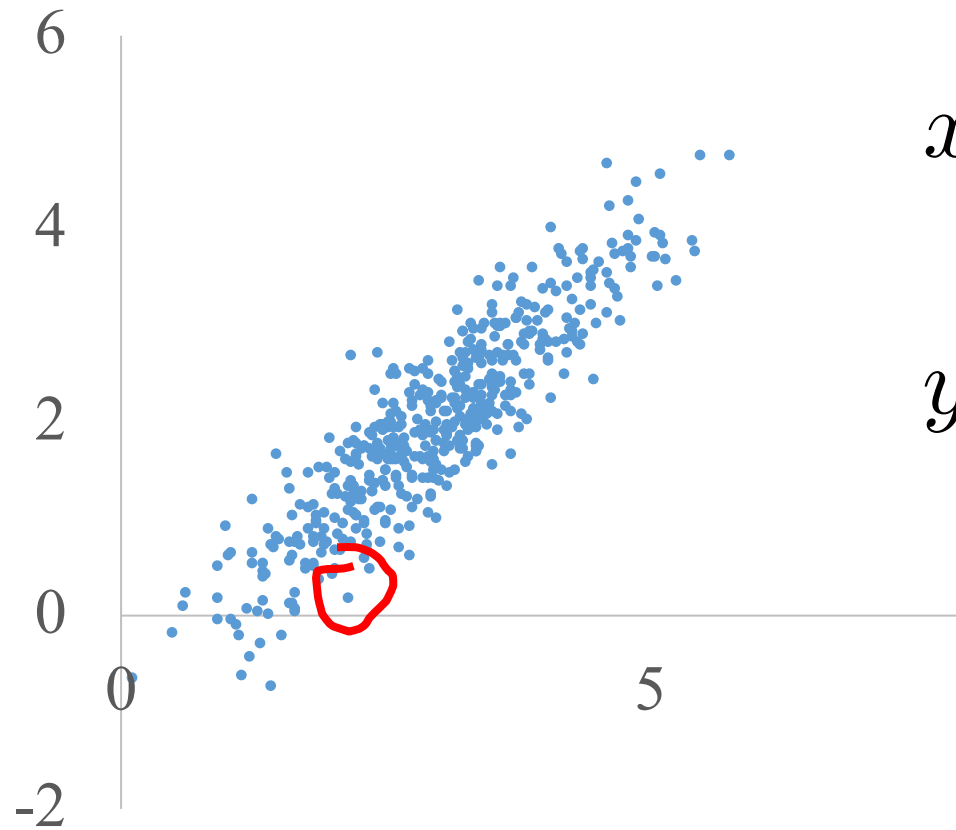


$$x - E[x] \approx 0$$

$$y - E[y] \approx 0$$

$$(x - E[x])(y - E[y]) = 0$$

Vary Together

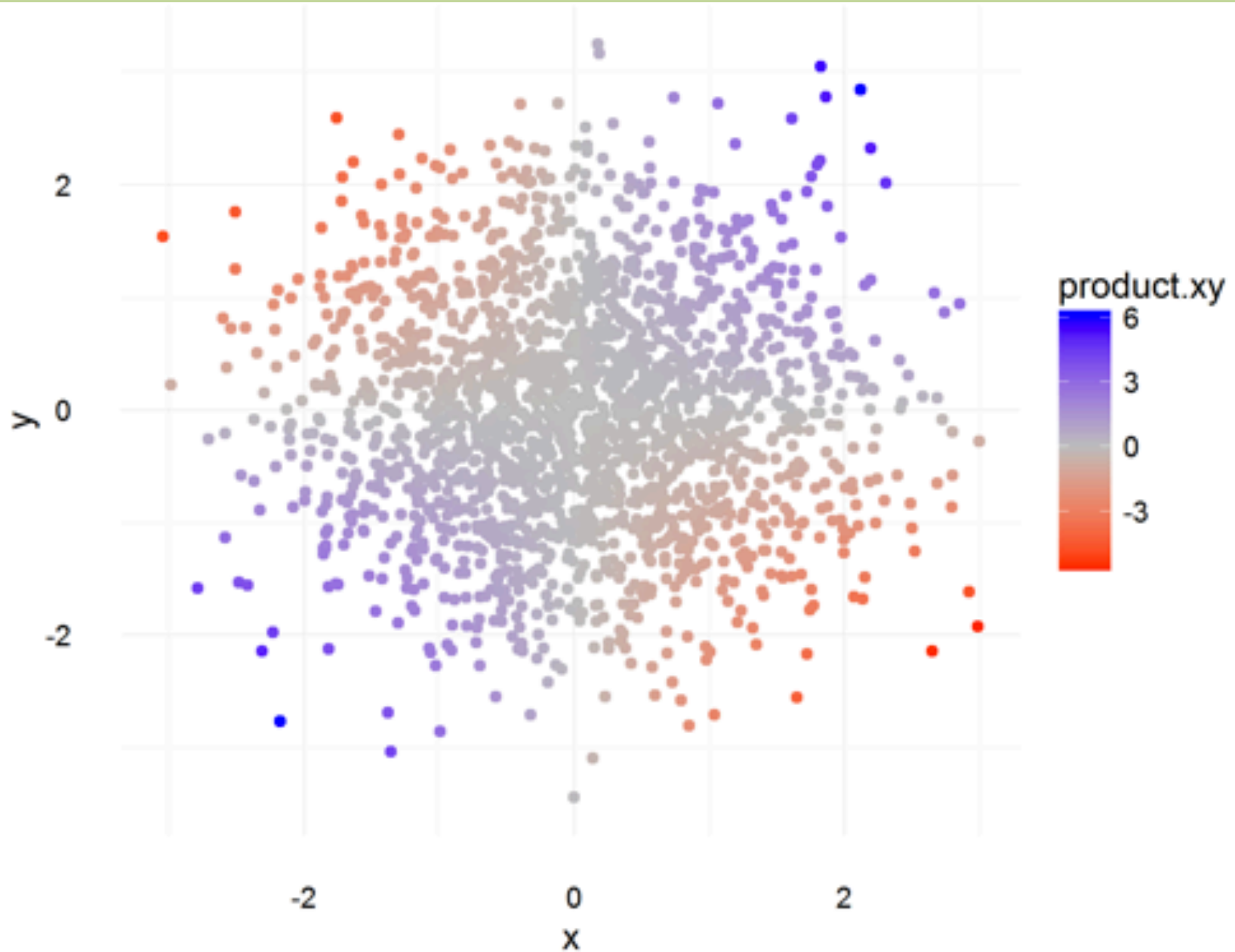


$$x - E[x] = -1.1$$

$$y - E[y] = -2.8$$

$$(x - E[x])(y - E[y]) \approx 3.1$$

Understanding Covariance



The Dance of the Covariance

- Say X and Y are arbitrary random variables
- Covariance of X and Y :

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

x	y	$(x - E[X])(y - E[Y])p(x,y)$
Above mean	Above mean	Positive
Bellow mean	Bellow mean	Positive
Bellow mean	Above mean	Negative
Above mean	Bellow mean	Negative

The Dance of the Covariance

- Say X and Y are arbitrary random variables
- Covariance of X and Y :

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

- Equivalently:

$$\begin{aligned}\text{Cov}(X, Y) &= E[XY - E[X]Y - XE[Y] + E[Y]E[X]] \\ &= E[XY] - E[X]E[Y] - E[X]E[Y] + E[X]E[Y] \\ &= E[XY] - E[X]E[Y]\end{aligned}$$

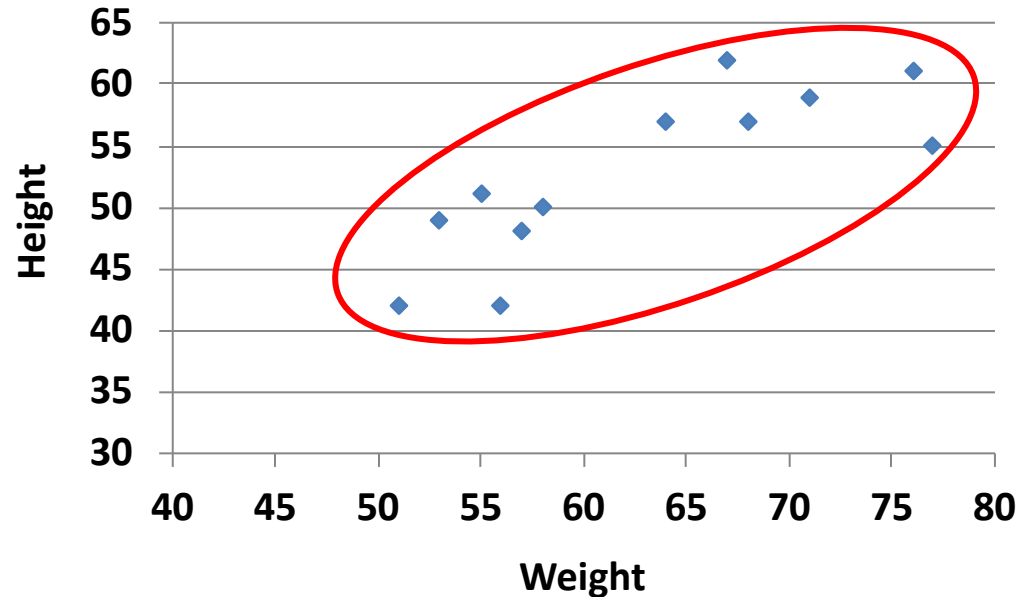
- X and Y independent, $E[XY] = E[X]E[Y] \rightarrow \text{Cov}(X, Y) = 0$
- But $\text{Cov}(X, Y) = 0$ does **not** imply X and Y independent!

Covariance and Data

- Consider the following data:

Weight	Height	Weight * Height
64	57	3648
71	59	4189
53	49	2597
67	62	4154
55	51	2805
58	50	2900
77	55	4235
57	48	2736
56	42	2352
51	42	2142
76	61	4636
68	57	3876

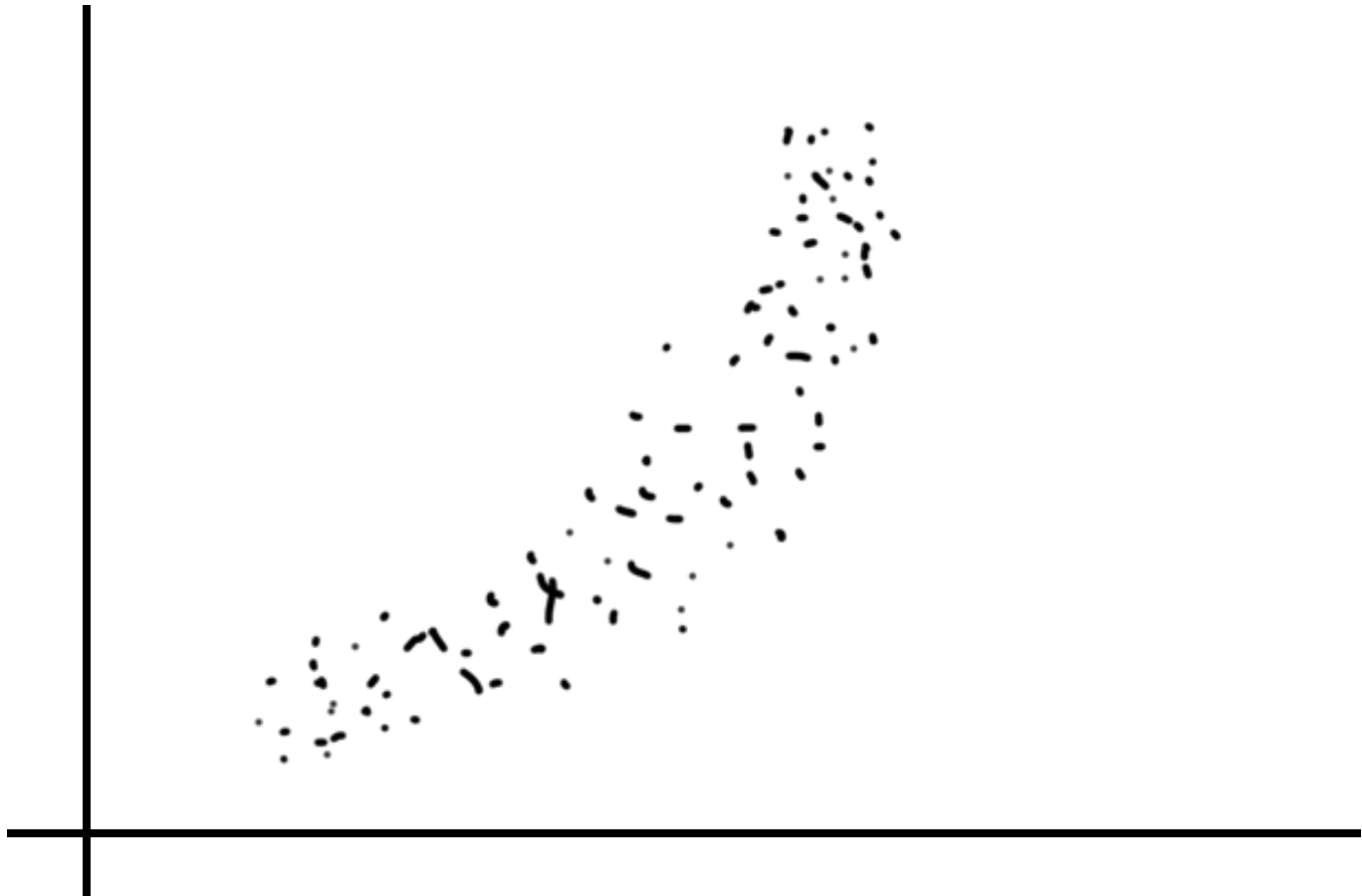
$$\begin{aligned} E[W] &= 62.75 & E[H] &= 52.75 & E[W*H] &= 3355.83 \end{aligned}$$



$$\begin{aligned} \text{Cov}(W, H) &= E[W*H] - E[W]E[H] \\ &= 3355.83 - (62.75)(52.75) \\ &= 45.77 \end{aligned}$$

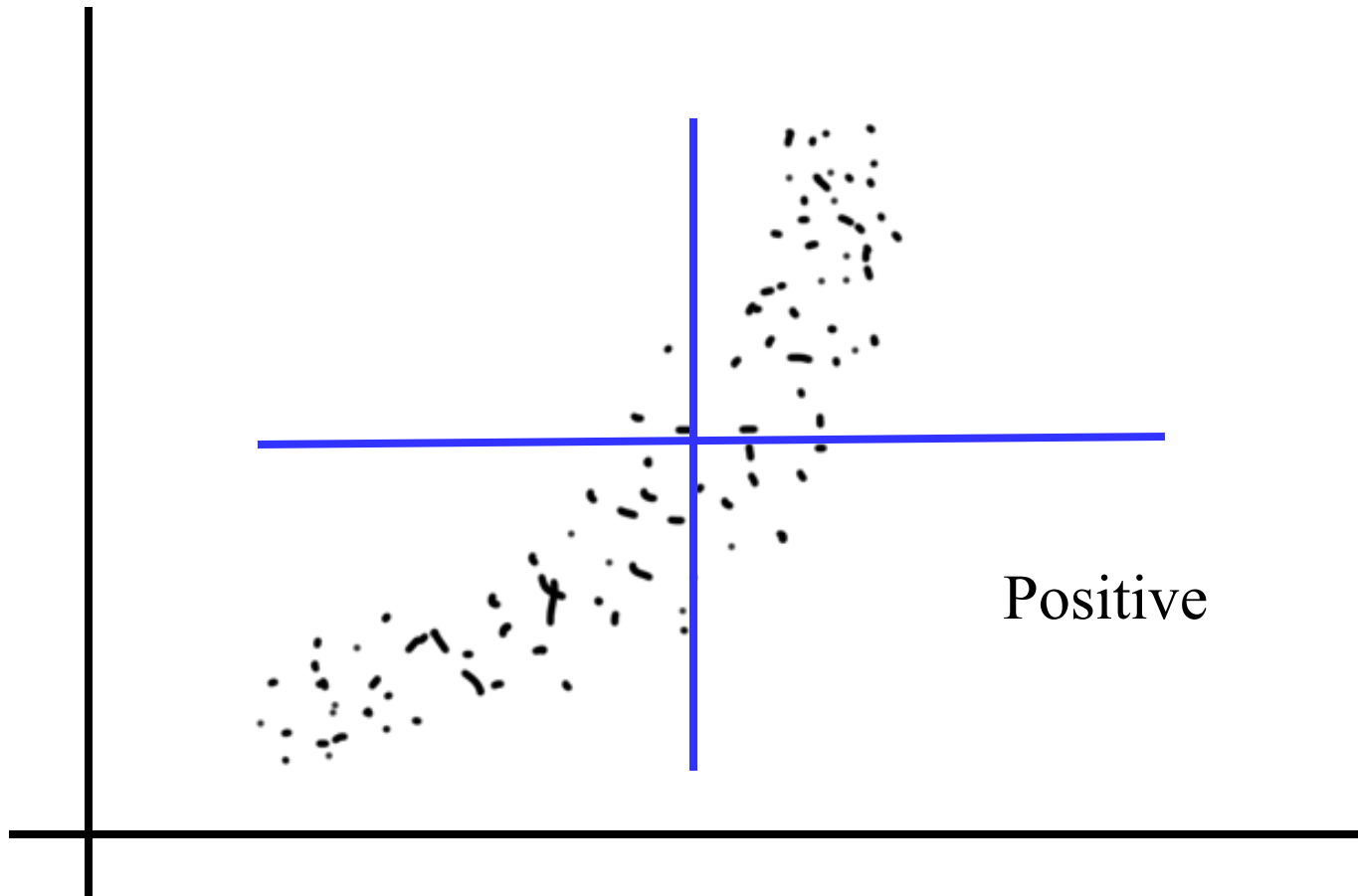
Covariance

Poll: (a) positive, (b) negative, (c) zero



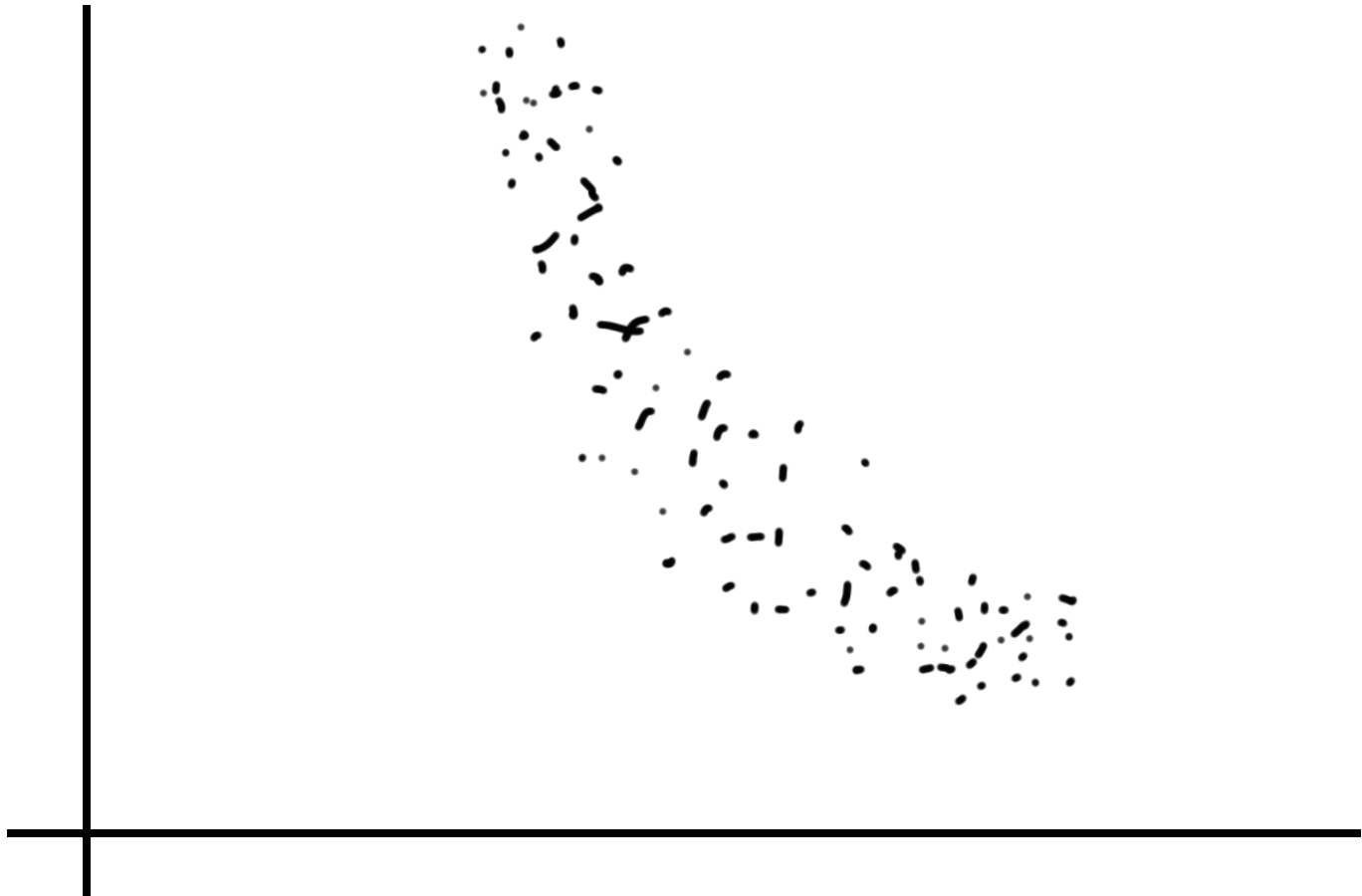
Covariance

Is the Covariance: (a) positive, (b) negative, (c) zero



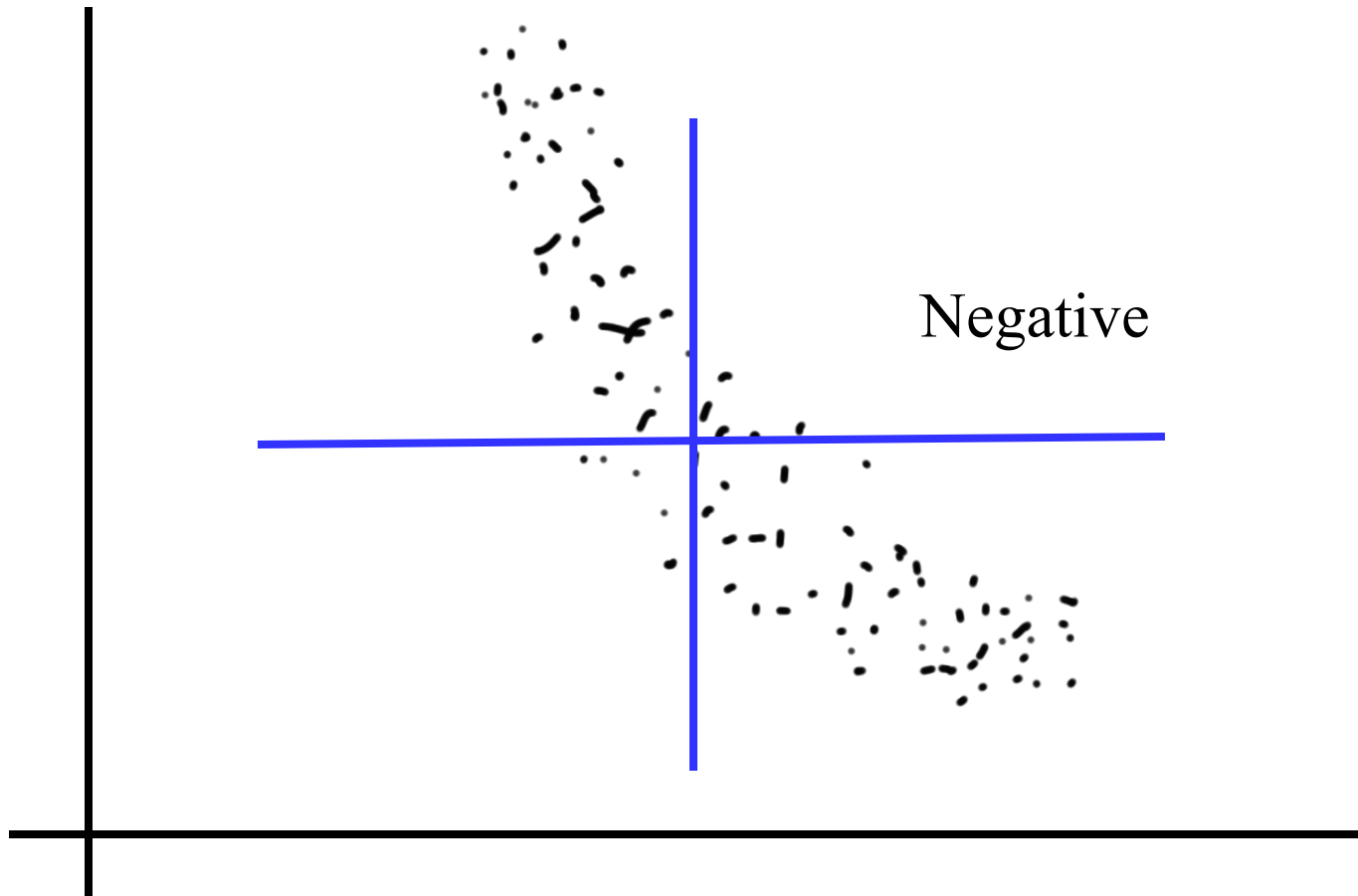
Covariance

Is the Covariance: (a) positive, (b) negative, (c) zero



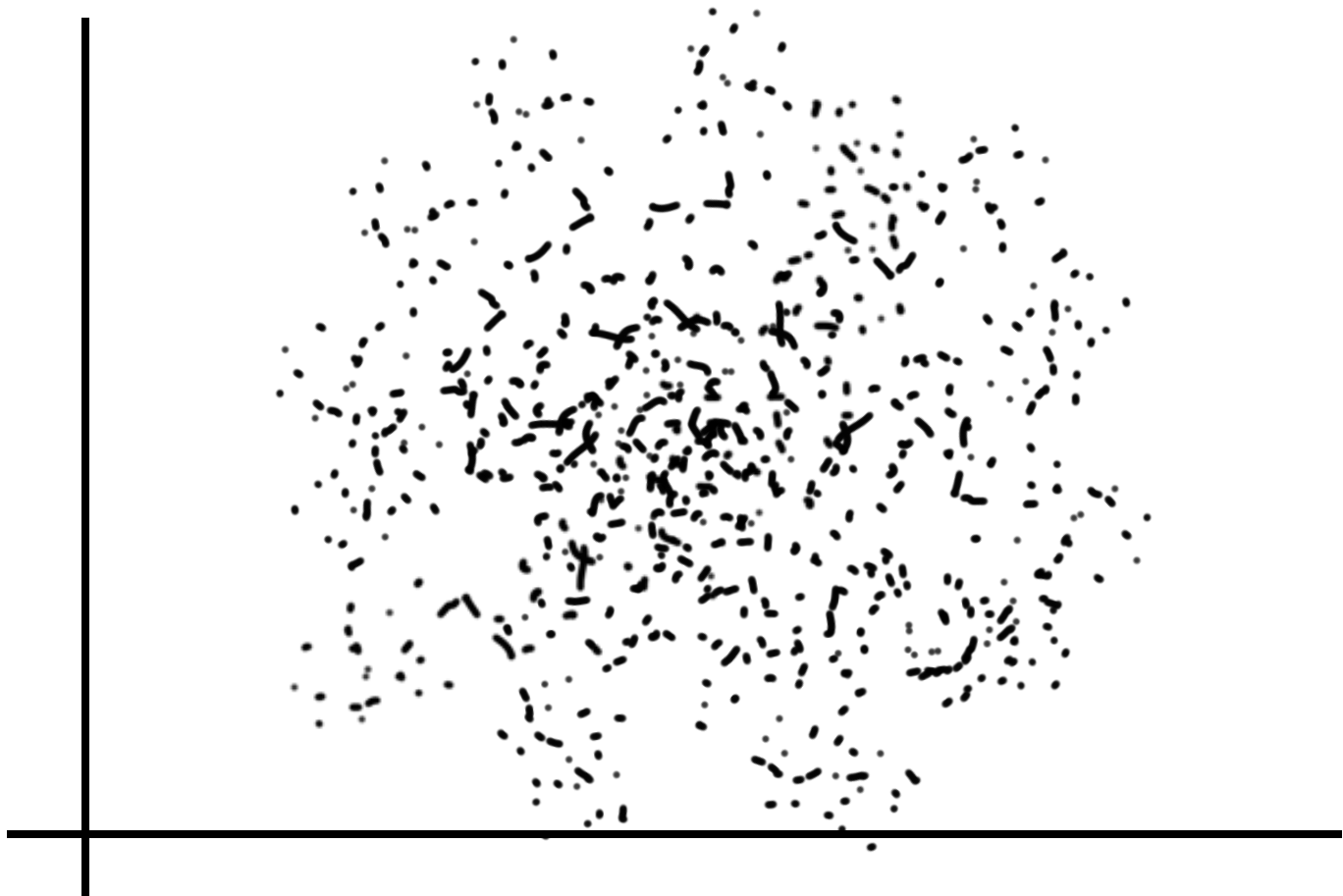
Covariance

Is the Covariance: (a) positive, (b) negative, (c) zero



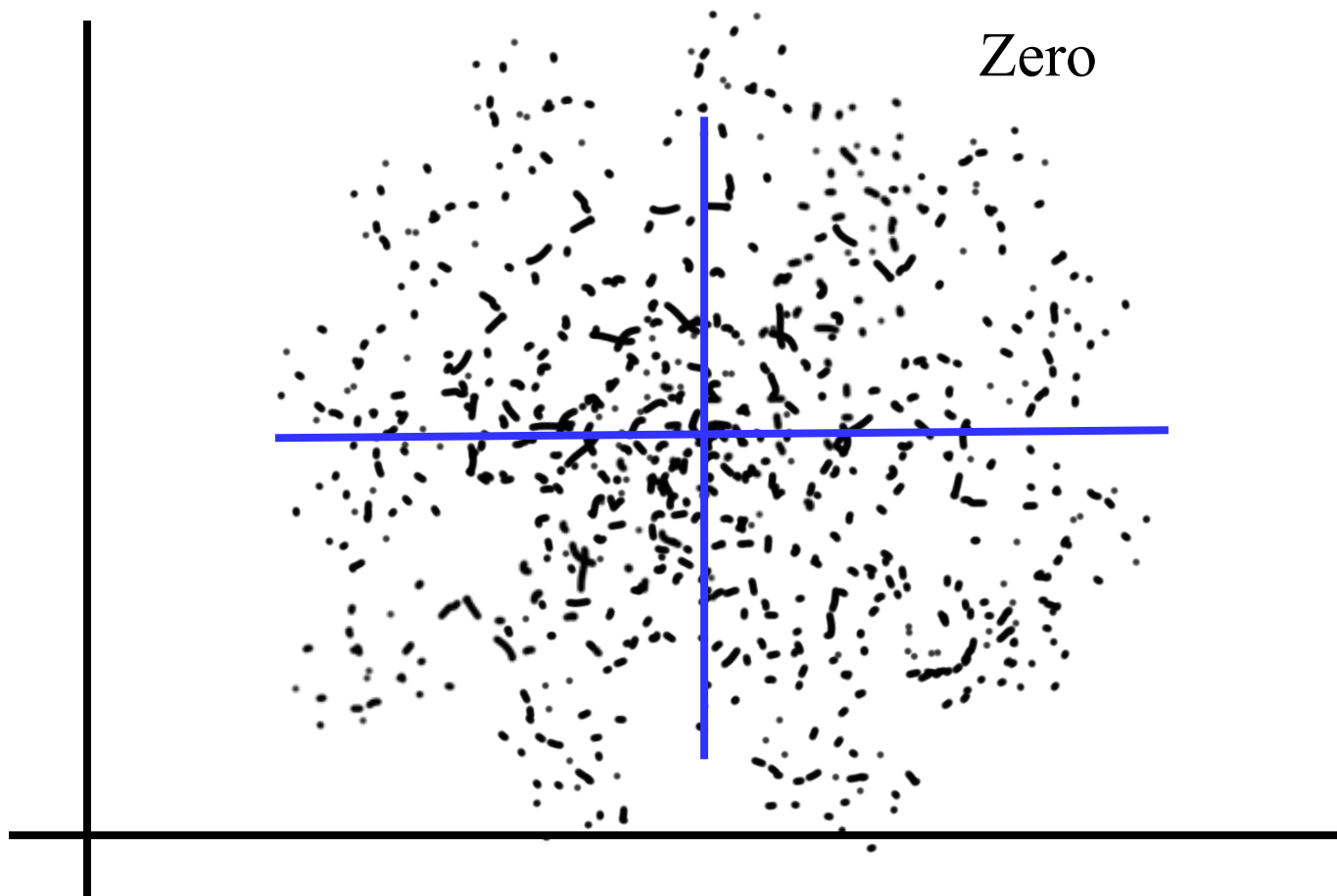
Covariance

Is the Covariance: (a) positive, (b) negative, (c) zero



Covariance

Is the Covariance: (a) positive, (b) negative, (c) zero



Independence and Covariance

- X and Y are random variables with PMF:

$Y \backslash X$	-1	0	1	$p_Y(y)$
0	1/3	0	1/3	2/3
1	0	1/3	0	1/3
$p_X(x)$	1/3	1/3	1/3	1

$$Y = \begin{cases} 0 & \text{if } X \neq 0 \\ 1 & \text{otherwise} \end{cases}$$

- $E[X] = -1(1/3) + 0(1/3) + 1(1/3) = 0$
 - $E[Y] = 0(2/3) + 1(1/3) = 1/3$
 - Since $XY = 0$, $E[XY] = 0$
 - $\text{Cov}(X, Y) = E[XY] - E[X]E[Y] = 0 - 0 = 0$
- But, X and Y are clearly dependent!

Properties of Covariance

- Say X and Y are arbitrary random variables
 - $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
 - $\text{Cov}(X, X) = E[X^2] - E[X]E[X] = \text{Var}(X)$
 - $\text{Cov}(aX + b, Y) = a\text{Cov}(X, Y)$
- Covariance of sums of random variables
 - X_1, X_2, \dots, X_n and Y_1, Y_2, \dots, Y_m are random variables
 - $$\text{Cov}\left(\sum_{i=1}^n X_i, \sum_{j=1}^m Y_j\right) = \sum_{i=1}^n \sum_{j=1}^m \text{Cov}(X_i, Y_j)$$

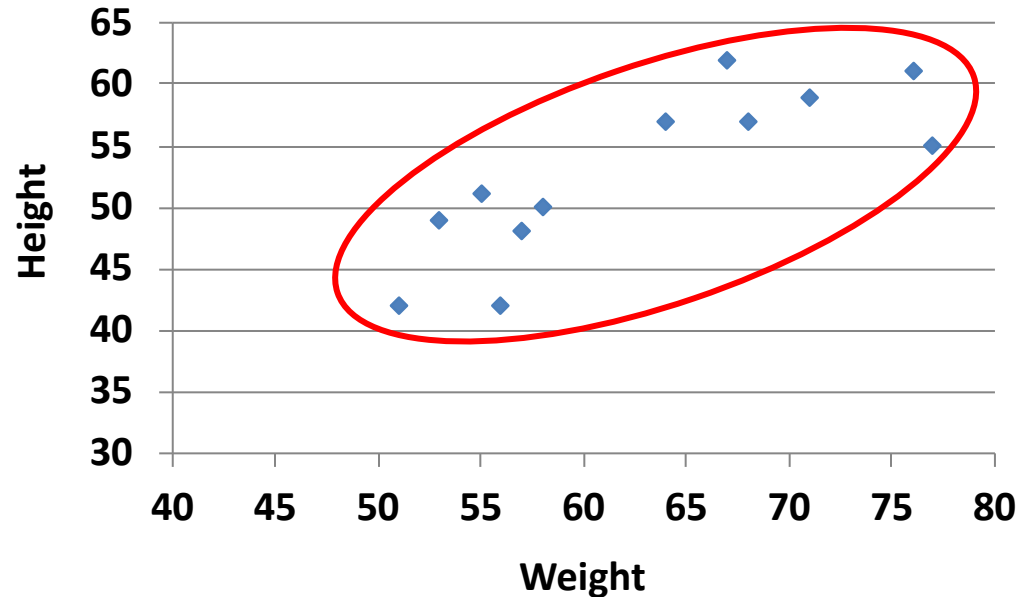
Correlation

What is Wrong With This?

- Consider the following data:

Weight	Height	Weight * Height
64	57	3648
71	59	4189
53	49	2597
67	62	4154
55	51	2805
58	50	2900
77	55	4235
57	48	2736
56	42	2352
51	42	2142
76	61	4636
68	57	3876

$$\begin{array}{lll} E[W] & E[H] & E[W*H] \\ = 62.75 & = 52.75 & = 3355.83 \end{array}$$



$$\begin{aligned} \text{Cov}(W, H) &= E[W*H] - E[W]E[H] \\ &= 3355.83 - (62.75)(52.75) \\ &= 45.77 \end{aligned}$$

The image shows a browser window displaying the Wikipedia article for the Cauchy–Schwarz inequality. The browser's address bar shows the URL https://en.wikipedia.org/wiki/Cauchy–Schwarz_inequali.... The user is logged in as Chris Piech. The article title is "Cauchy–Schwarz inequality". The text explains that in mathematics, the Cauchy–Schwarz inequality, also known as the Cauchy–Bunyakovsky–Schwarz inequality, is a useful inequality encountered in many different settings, such as linear algebra, analysis, probability theory, vector algebra and other areas. It is considered to be one of the most important inequalities in all of mathematics.^[1] It has a number of generalizations, among them Hölder's inequality. The inequality for sums was published by Augustin-Louis Cauchy (1821), while the corresponding inequality for integrals was first proved by Viktor Bunyakovsky (1859). The modern proof of the integral inequality was given by Hermann Amandus Schwarz (1888).^[1]

The article includes a table of contents with the following sections:

- 1 Statement of the inequality
- 2 Proofs
 - 2.1 First proof
 - 2.2 Second proof
 - 2.3 More proofs
- 3 Special cases
 - 3.1 \mathbb{R}^2 (ordinary two-dimensional space)
 - 3.2 \mathbb{R}^n (n -dimensional Euclidean space)
 - 3.3 L^2

$$-\text{Std}(X)\text{Std}(Y) \leq \text{Cov}(X, Y) \leq \text{Std}(X)\text{Std}(Y)$$

Viva La Correlación

- Say X and Y are arbitrary random variables

- Correlation of X and Y , denoted $\rho(X, Y)$:

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

- Note: $-1 \leq \rho(X, Y) \leq 1$

- Correlation measures linearity between X and Y

- $\rho(X, Y) = 1 \quad \Rightarrow \quad Y = aX + b \quad \text{where } a = \sigma_y/\sigma_x$

- $\rho(X, Y) = -1 \quad \Rightarrow \quad Y = aX + b \quad \text{where } a = -\sigma_y/\sigma_x$

- $\rho(X, Y) = 0 \quad \Rightarrow \quad \text{absence of linear relationship}$

- But, X and Y can still be related in some other way!

- If $\rho(X, Y) = 0$, we say X and Y are “uncorrelated”

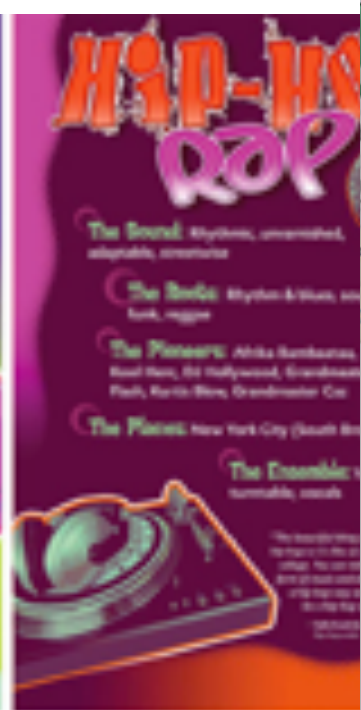
- Note: Independence implies uncorrelated, but **not** vice versa!

Viva La Correlación

- Say X and Y are arbitrary random variables
 - Correlation of X and Y , denoted $\rho(X, Y)$:

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$$

- Say $Y = cX$. Correlation should be 1.



AutoSave OFF Search Sheet Share

Home Insert Page Layout Formulas Data >>

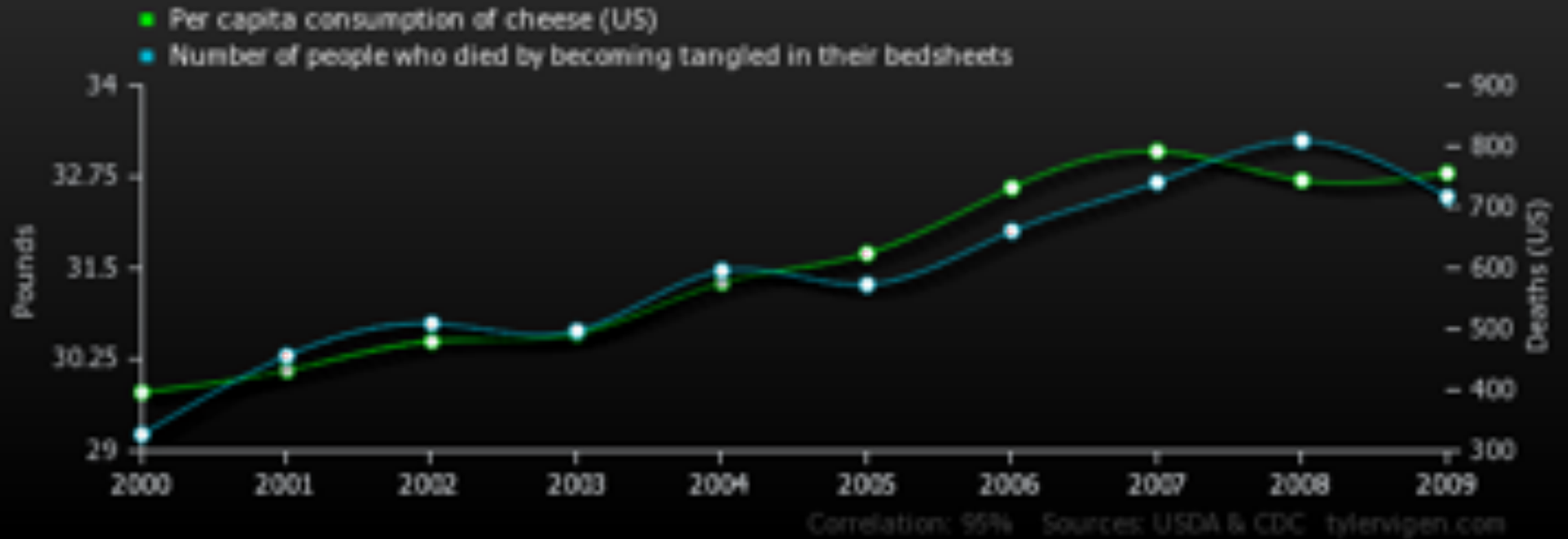
Clipboard Font Alignment Number Conditional Formatting Format as Table Cell Styles

C15 fx 3

	A	B	C	D	E	F	G	H	I
1	Music	Dance	Folk	Country	Classical music	Musical	Pop	Rock	Me
2	5	2	1	2	2	1	5	5	
3	4	2	1	1	1	2	3	5	
4	5	2	2	3	4	5	3	5	
5	5	2	1	1	1	1	2	2	
6	5	4	3	2	4	3	5	3	
7	5	2	3	2	3	3	2	5	
8	5	5	3	1	2	2	5	3	
9	5	3	2	1	2	2	4	5	
10	5	3	1	1	2	4	3	5	
11	5	2	5	2	2	5	3	5	
12	5	3	2	1	2	3	4	3	
13	5	1	1	1	4	1	2	5	
14	5	1	2	1	4	3	3	5	
15	5	5	3	2	1	5	5	2	
16	5	2	1	1	2	3	4	5	
17	1	2	2	3	4	3	3	5	
18	5	3	1	1	1	2	4	4	
19	5	3	3	3	2	2	4	4	
20	5	5	4	3	4	5	5	4	
21	5	3	3	2	4	2	2	4	
22	5	3	2	3	4	3	2	5	
23	5	1	1	3	2	2	2	5	
24	5	3	2	3	3	3	4		
25	5	4	2	2	2	4	4	5	
26	5	3	1	1	4	3	3	5	
27	5	4	2	1	2	3	5	1	
28	5	5	5	4	5	3	4	4	
29	4	3	4	1	3	2	2	4	
30	5	5	1	1	1	1	3	4	
31	5	3	4	2	3	3	3	4	
32	4	4	3	3	3	3	4	4	
33	4	4	1	3	2	3	5	3	
34	5	3	1	3	2	3	3	4	
35	5	2	2	3	4	5	4	3	

Ready music + 100%

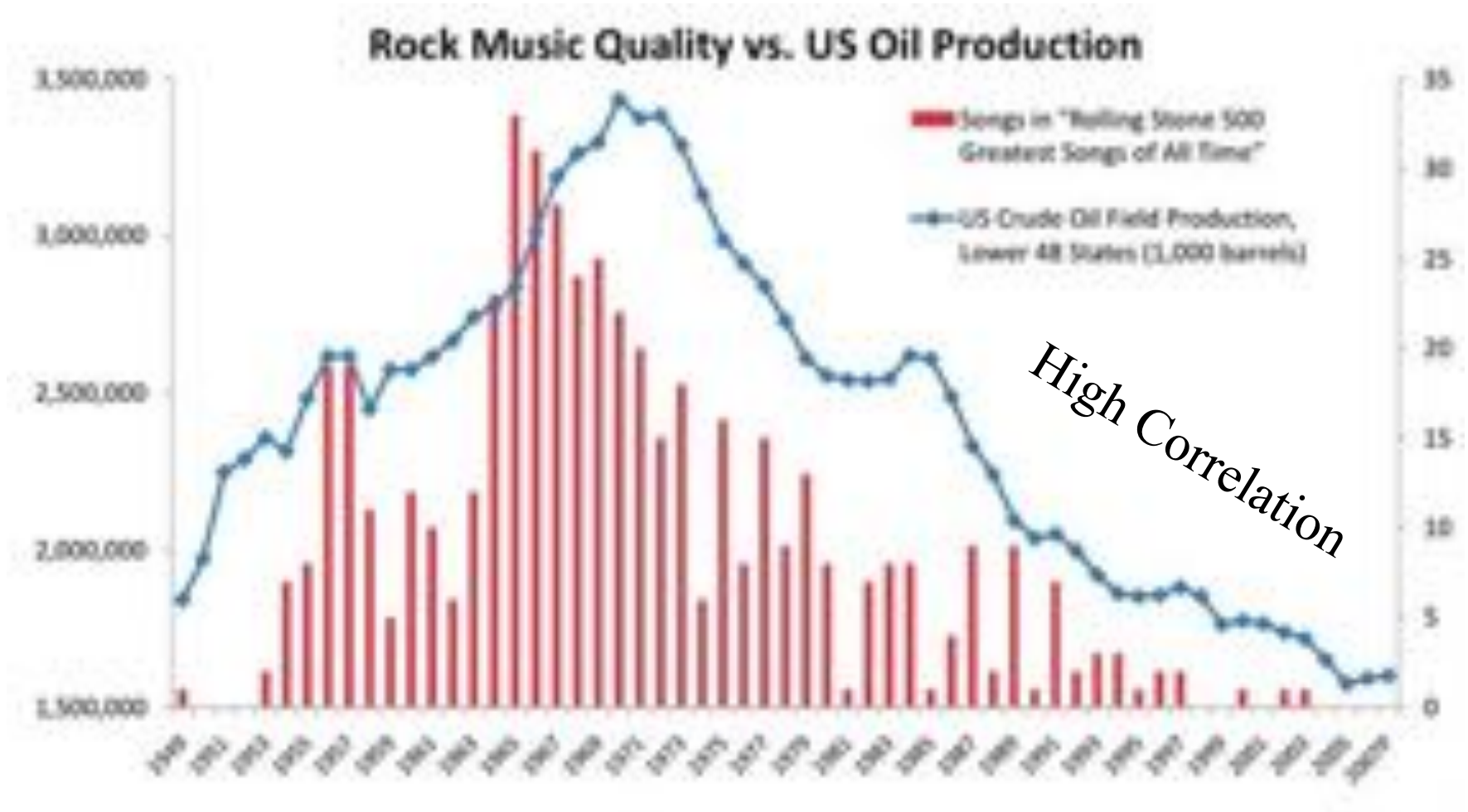
Tell your friends!



	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009
<i>Per capita consumption of cheese (US) Pounds (USDA)</i>	29.8	30.1	30.5	30.6	31.3	31.7	32.6	33.1	32.7	32.8
<i>Number of people who died by becoming tangled in their bedsheets Deaths (US) (CDC)</i>	327	456	509	497	596	573	661	741	809	717

Correlation: 0.947091

Rock Music Vs Oil?



Hubbert Peak Theory

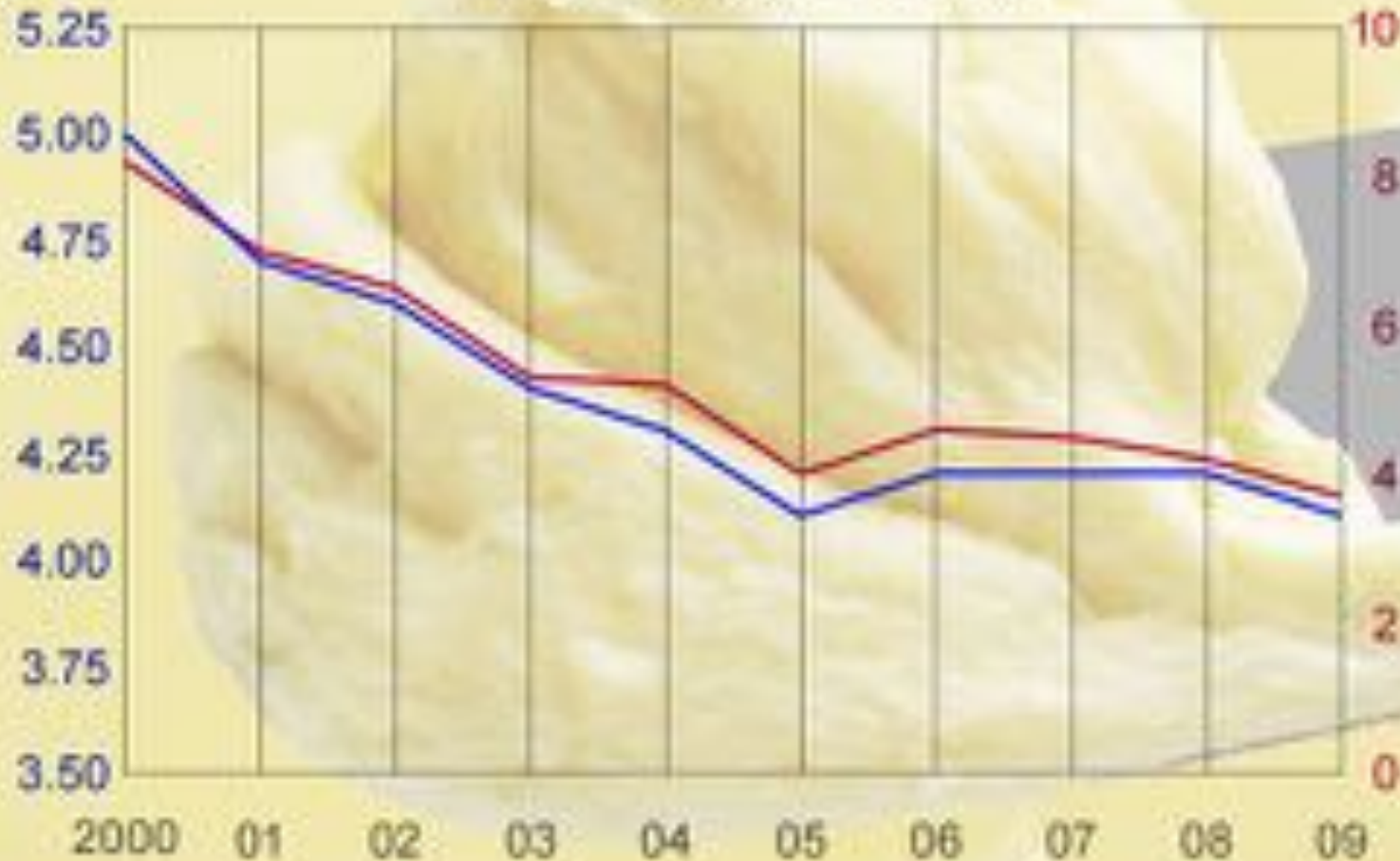
<http://www.aei.org/publication/blog/>

Divorce Vs Butter?

Divorce rate
in Maine per
1,000 people

Per capita
consumption of
margarine (lbs)

Correlation: 99%



Source: US Census, USDA, tylervigen.com

SPL

<http://www.bbc.com/news/magazine-27537142>



Great Expectations

Noah Arthurs

CS109, Stanford University

Conditional Expectation

Conditional Expectation

- X and Y are jointly discrete random variables

- Recall conditional PMF of X given $Y = y$:

$$p_{X|Y}(x | y) = P(X = x | Y = y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}$$

- Define conditional expectation of X given $Y = y$:

$$E[X | Y = y] = \sum_x x P(X = x | Y = y) = \sum_x x p_{X|Y}(x | y)$$

- Analogously, jointly continuous random variables:

$$f_{X|Y}(x | y) = \frac{f_{X,Y}(x, y)}{f_Y(y)} \qquad E[X | Y = y] = \int_{-\infty}^{\infty} x f_{X|Y}(x | y) dx$$

Rolling Dice

- 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 $E[X | Y = 6]$?

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

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

Properties of Conditional Expectation

- X and Y are jointly distributed random variables

$$E[g(X) | Y = y] = \sum_x g(x) p_{X|Y}(x | y) \quad \text{or} \quad \int_{-\infty}^{\infty} g(x) f_{X|Y}(x | y) dx$$

- Expectation of conditional sum:

$$E\left[\sum_{i=1}^n X_i | Y = y\right] = \sum_{i=1}^n E[X_i | Y = y]$$

Conditional Expectation Functions

This is a number:

$$E[X]$$



This is a function of y :

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

$$E[X = 5]$$

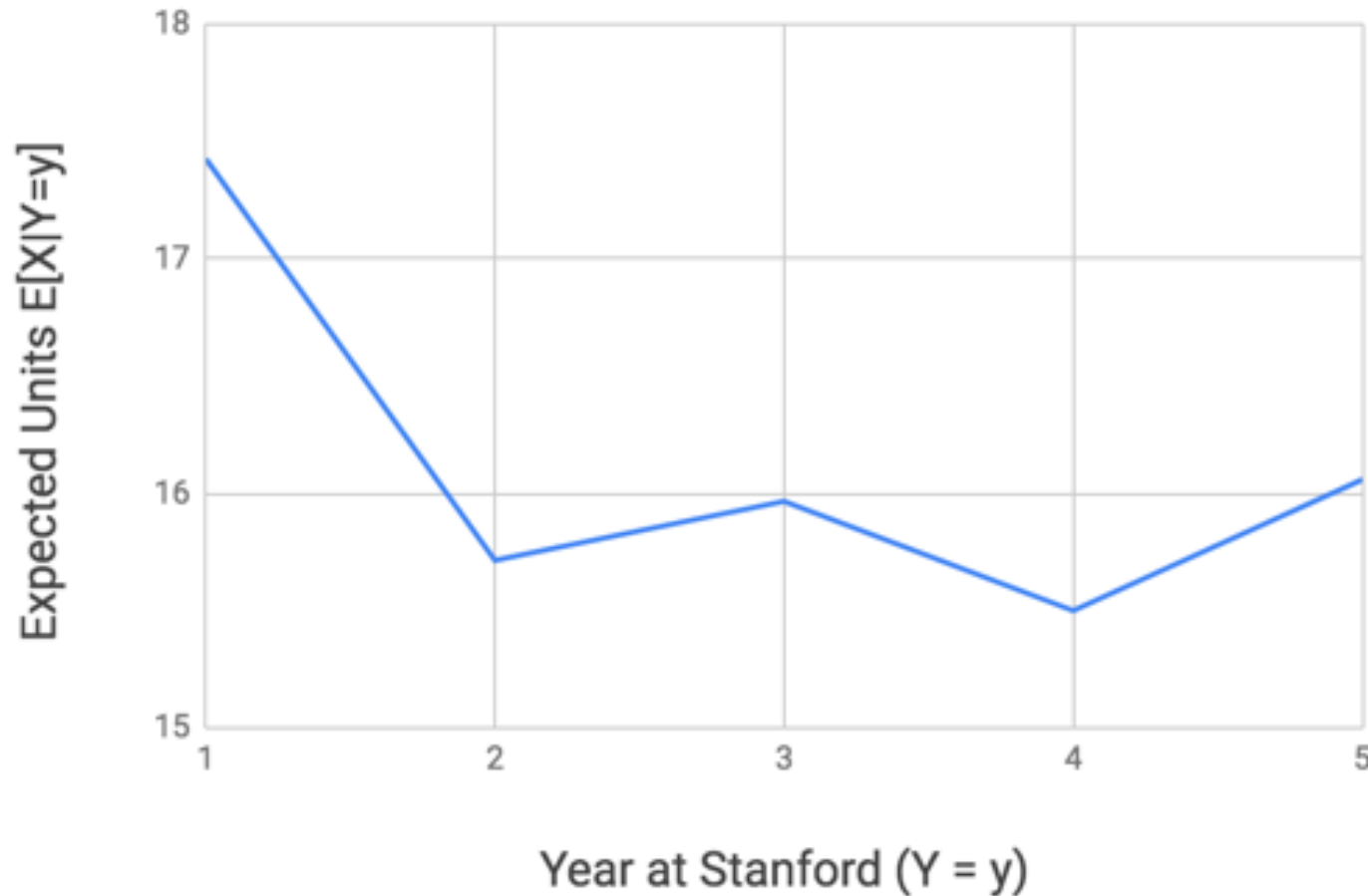
Doesn't make sense. Take expectation of random variables, not events

Conditional Expectation Functions

X = units in fall quarter

Y = year in school

$E[X | Y]$?



Law of Total Expectation

For any random variable X and any discrete random variable Y



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

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 $E[Y]$?

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

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

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

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

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

$$E[Y] = 15$$

Protip: do this in CS161