

A multi-layered cake with white frosting and fresh strawberries on top, set against a light wooden background. The cake is presented on a white stand. The top layer is decorated with a generous amount of white frosting and several fresh strawberries. The middle layers are also visible, showing the same frosting and fruit filling.

CS 109 Review

Noah Arthurs

June 3, 2019

Where we're at

Last week: ML wrap-up, theoretical background for modern ML

This week: course overview, open questions after CS 109

Next week: final exam Tuesday!

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

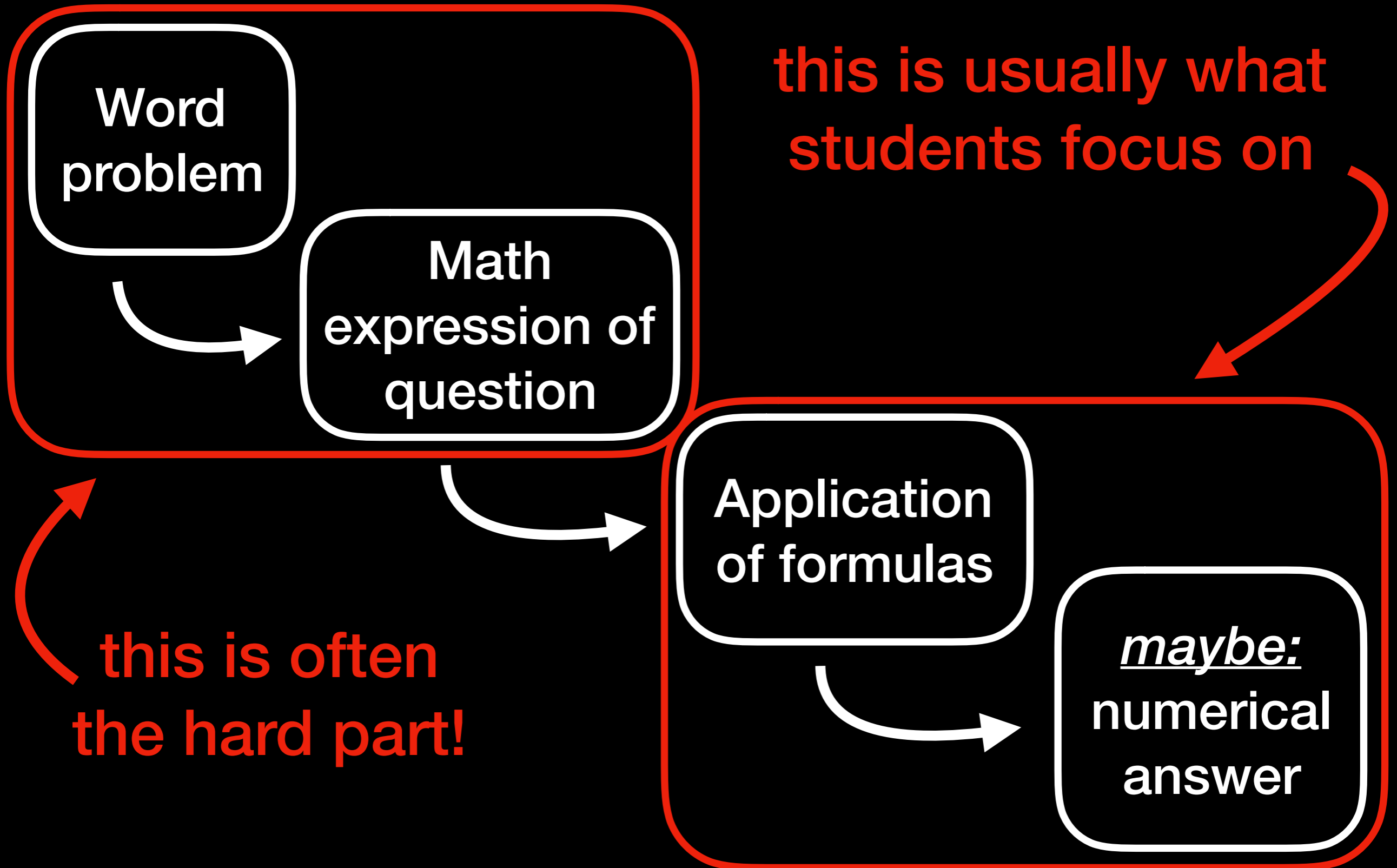
core probability fundamentals

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

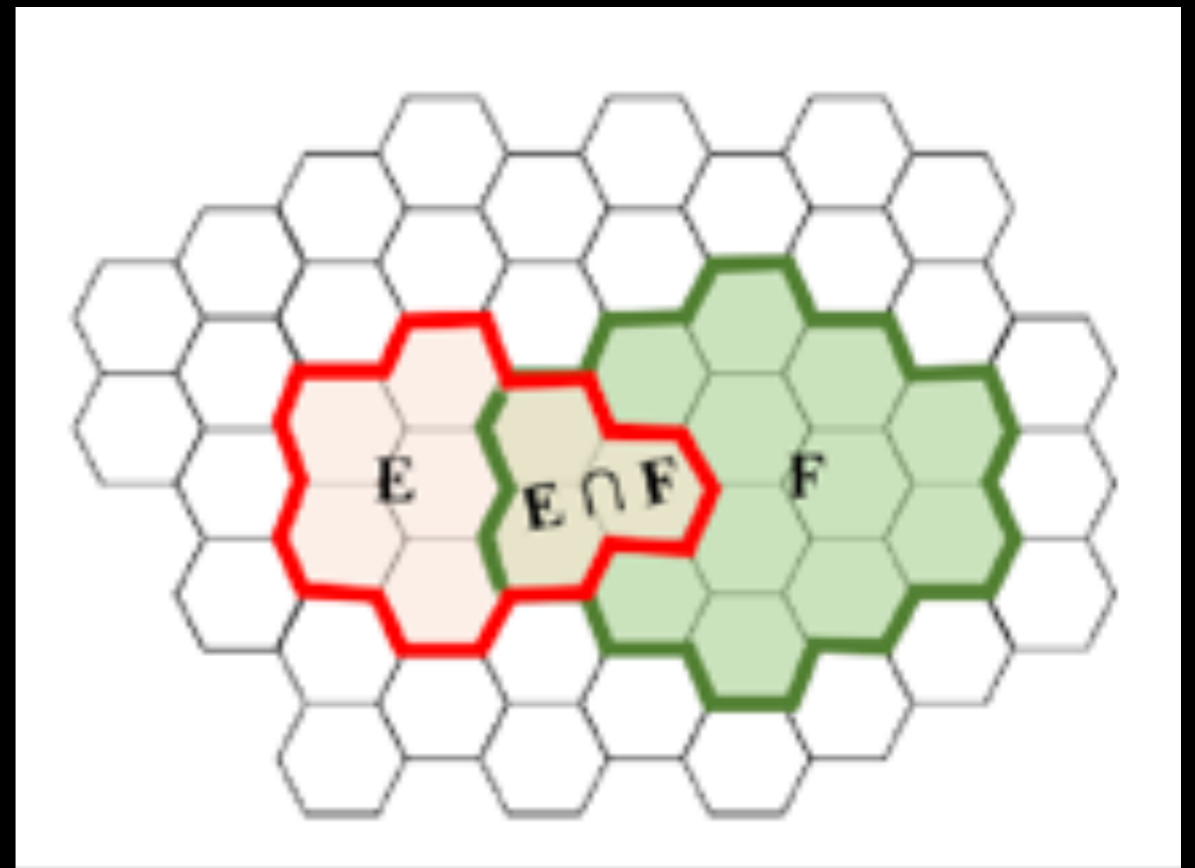
Conditional Probability

definition:

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

Chain Rule:

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



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

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

The diagram illustrates Bayes' Rule with the following components: The word "posterior" is written in pink above the left side of the equation, with a pink arrow pointing to the term $P(E|F)$. The word "likelihood" is written in pink above the numerator, with a pink arrow pointing to the term $P(F|E)$. The word "prior" is written in pink above the numerator, with a pink arrow pointing to the term $P(E)$. The word "normalization constant" is written in pink below the denominator, with a pink arrow pointing to the term $P(F)$.

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

$$P(ABC) = ?$$

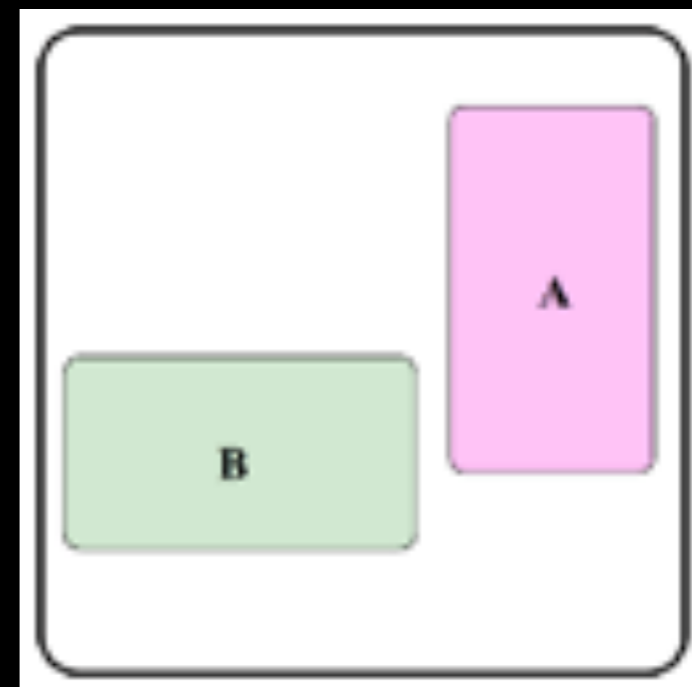
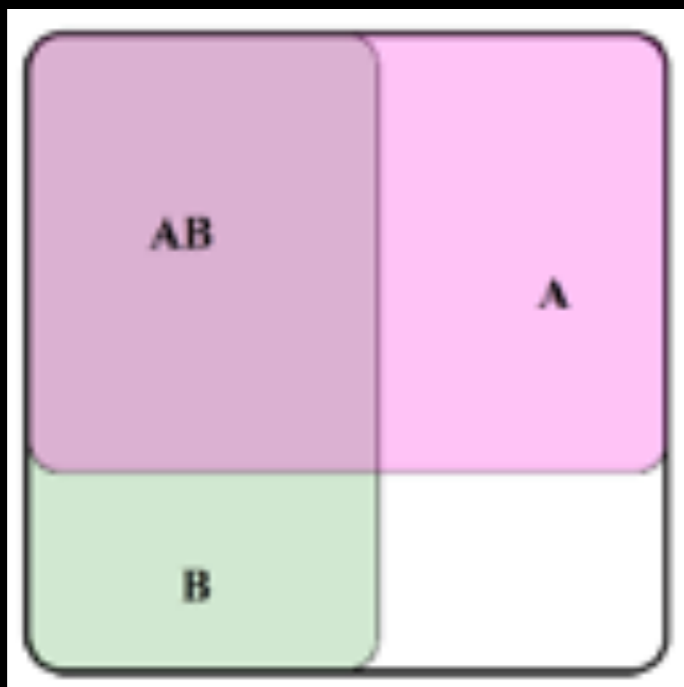
Let $X = AC$:

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

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

The “Gary” Problem

$$\begin{aligned} P(1950|C, \text{Gary}) &= \frac{P(\text{Gary}|1950, C)P(1950|C)}{P(\text{Gary}|C)} \\ &= \frac{P(\text{Gary}|1950, C)P(1950|C)}{\sum_{y=1901}^{2019} P(\text{Gary}|y, C)P(y|C)} \\ &= \frac{P(\text{Gary}|1950)P(1950|C)}{\sum_{y=1901}^{2019} P(\text{Gary}|y)P(y|C)} \\ &= \frac{\left(\frac{\text{count}(1950, \text{Gary})}{\sum_{i=1}^k \text{count}(1950, N_i)} \right) \text{prior}(1950)}{\sum_{y=1901}^{2019} \left(\frac{\text{count}(y, \text{Gary})}{\sum_{i=1}^k \text{count}(y, N_i)} \right) \text{prior}(y)} \end{aligned}$$

The “Gary” Problem

$$\begin{aligned} P(1950|C, \text{Gary}) &= \frac{P(\text{Gary}|1950, C)P(1950|C)}{P(\text{Gary}|C)} && \text{Conditional Bayes} \\ &= \frac{P(\text{Gary}|1950, C)P(1950|C)}{\sum_{y=1901}^{2019} P(\text{Gary}|y, C)P(y|C)} \\ &= \frac{P(\text{Gary}|1950)P(1950|C)}{\sum_{y=1901}^{2019} P(\text{Gary}|y)P(y|C)} \\ &= \frac{\left(\frac{\text{count}(1950, \text{Gary})}{\sum_{i=1}^k \text{count}(1950, N_i)} \right) \text{prior}(1950)}{\sum_{y=1901}^{2019} \left(\frac{\text{count}(y, \text{Gary})}{\sum_{i=1}^k \text{count}(y, N_i)} \right) \text{prior}(y)} \end{aligned}$$

The “Gary” Problem

$$\begin{aligned} P(1950|C, \text{Gary}) &= \frac{P(\text{Gary}|1950, C)P(1950|C)}{P(\text{Gary}|C)} \\ &= \frac{P(\text{Gary}|1950, C)P(1950|C)}{\sum_{y=1901}^{2019} P(\text{Gary}|y, C)P(y|C)} \\ &= \frac{P(\text{Gary}|1950)P(1950|C)}{\sum_{y=1901}^{2019} P(\text{Gary}|y)P(y|C)} \\ &= \frac{\left(\frac{\text{count}(1950, \text{Gary})}{\sum_{i=1}^k \text{count}(1950, N_i)} \right) \text{prior}(1950)}{\sum_{y=1901}^{2019} \left(\frac{\text{count}(y, \text{Gary})}{\sum_{i=1}^k \text{count}(y, N_i)} \right) \text{prior}(y)} \end{aligned}$$

“Full” Law of
Total Probability

The “Gary” Problem

$$\begin{aligned} P(1950|C, \text{Gary}) &= \frac{P(\text{Gary}|1950, C)P(1950|C)}{P(\text{Gary}|C)} \\ &= \frac{P(\text{Gary}|1950, C)P(1950|C)}{\sum_{y=1901}^{2019} P(\text{Gary}|y, C)P(y|C)} \\ &= \frac{P(\text{Gary}|1950)P(1950|C)}{\sum_{y=1901}^{2019} P(\text{Gary}|y)P(y|C)} \\ &= \frac{\left(\frac{\text{count}(1950, \text{Gary})}{\sum_{i=1}^k \text{count}(1950, N_i)} \right) \text{prior}(1950)}{\sum_{y=1901}^{2019} \left(\frac{\text{count}(y, \text{Gary})}{\sum_{i=1}^k \text{count}(y, N_i)} \right) \text{prior}(y)} \end{aligned}$$

Conditional Independence
given in problem

CS 109

topics

machine learning

sampling, making conclusions from data

multivariate distributions

random variables / distributions

discrete RVs

continuous RVs

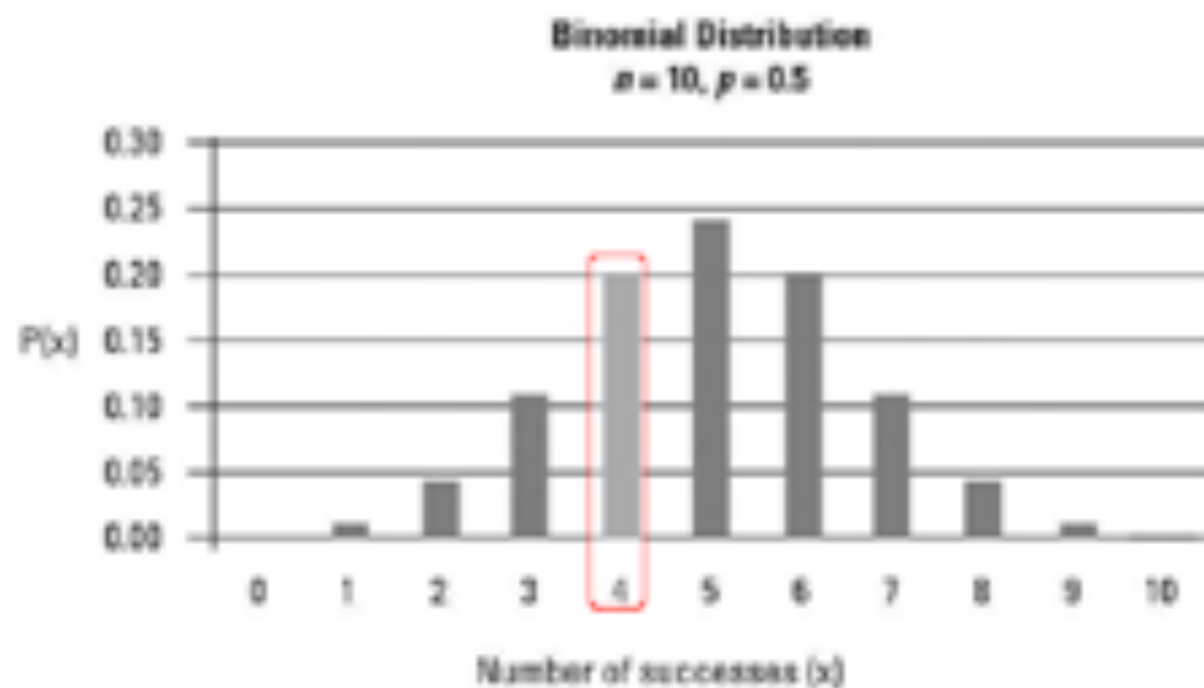
core probability fundamentals

properties of RVs

Probability Distributions

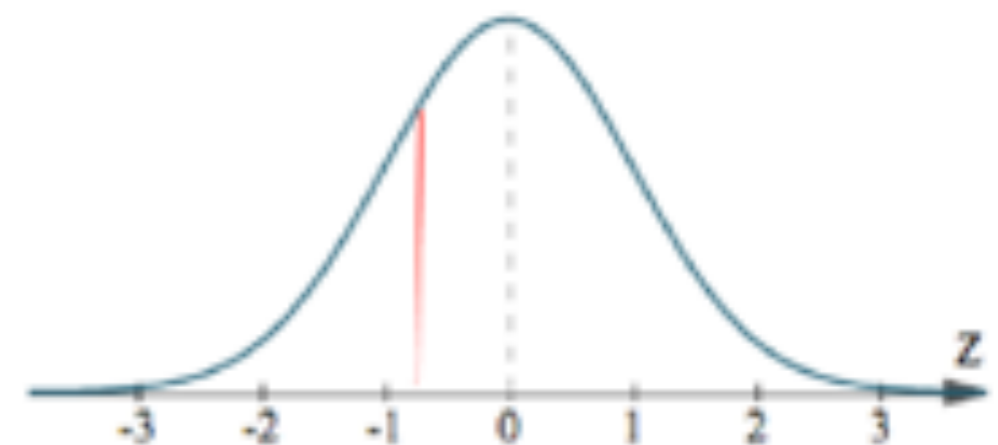
Discrete

PMF:



Continuous

PDF:

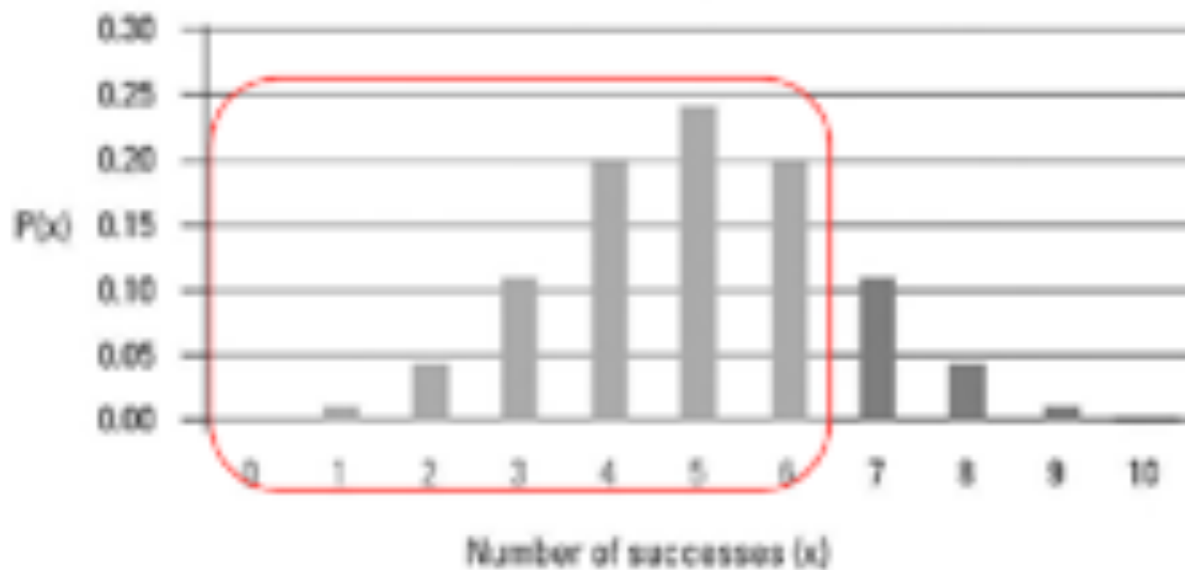


Probability Distributions

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

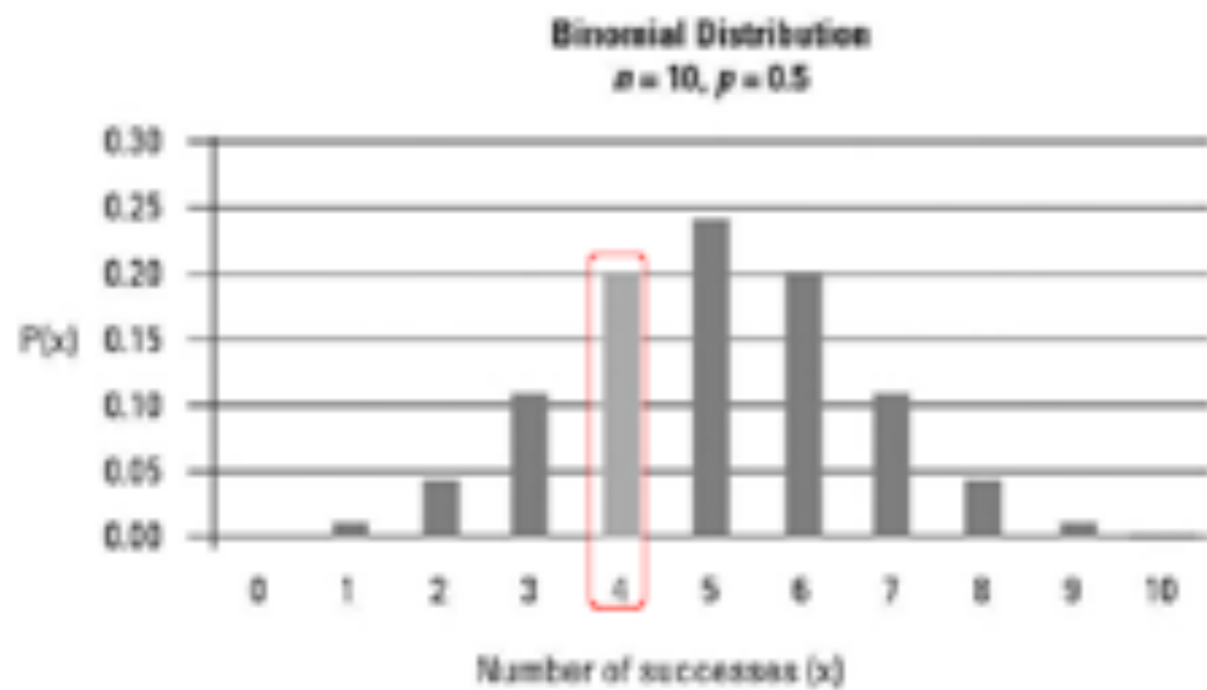
Poisson

p is small

Continuity Correction

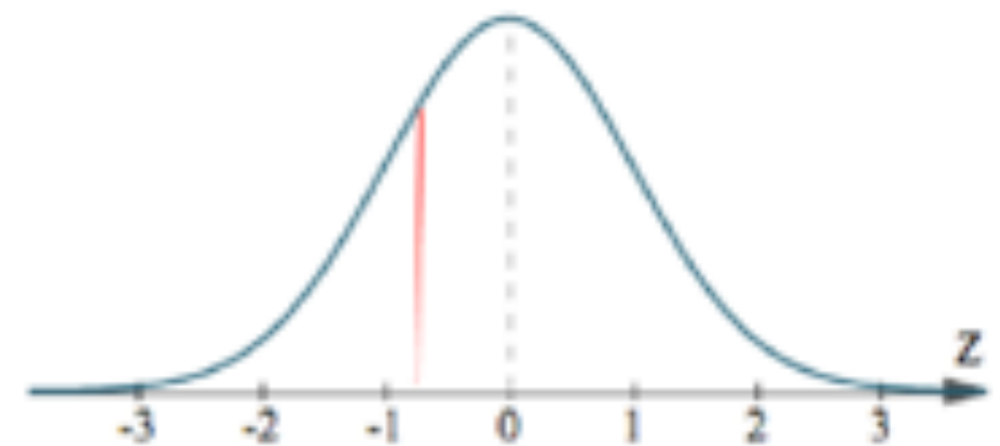
Discrete

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

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

Relationships Between Random Variables

Covariance

the extent to which the deviation of one variable from its mean matches the deviation of the other from its mean

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

Correlation

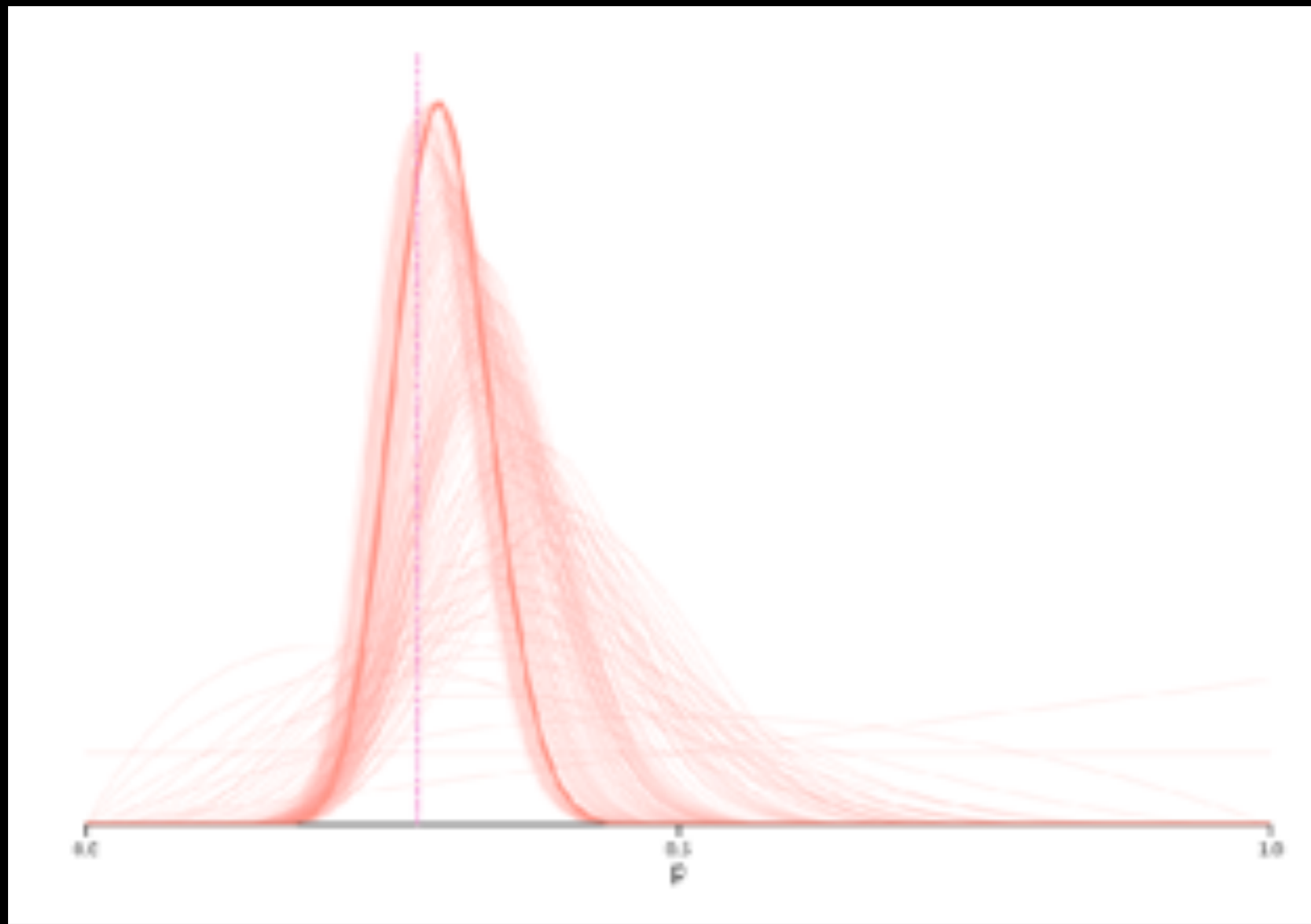
covariance normalized by the variance of each variable
(cancels the units out)

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

if two random variables are independent, they have a covariance of 0
(but not necessarily true the other way around!)

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

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

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.

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

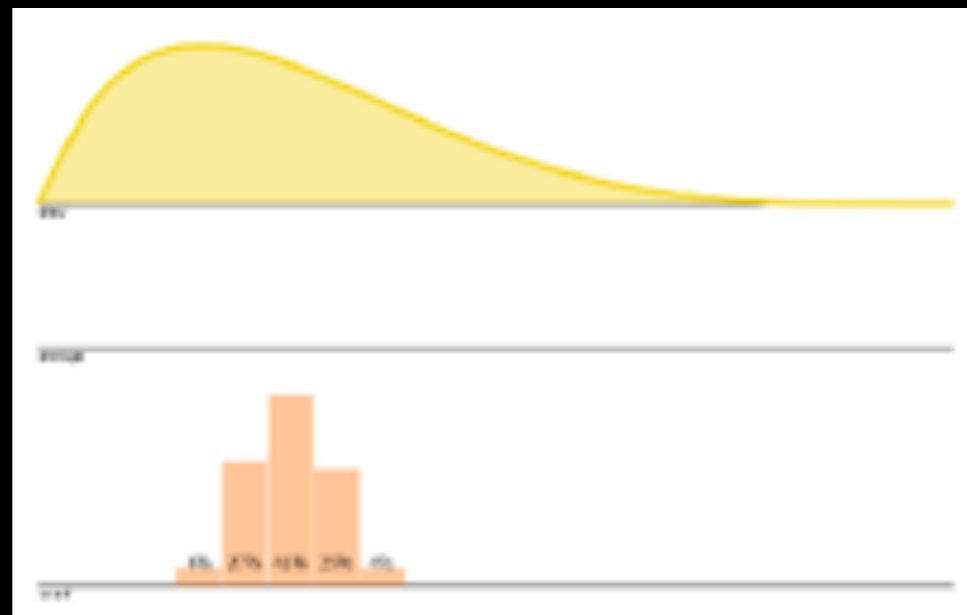
Taking Many Samples

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!



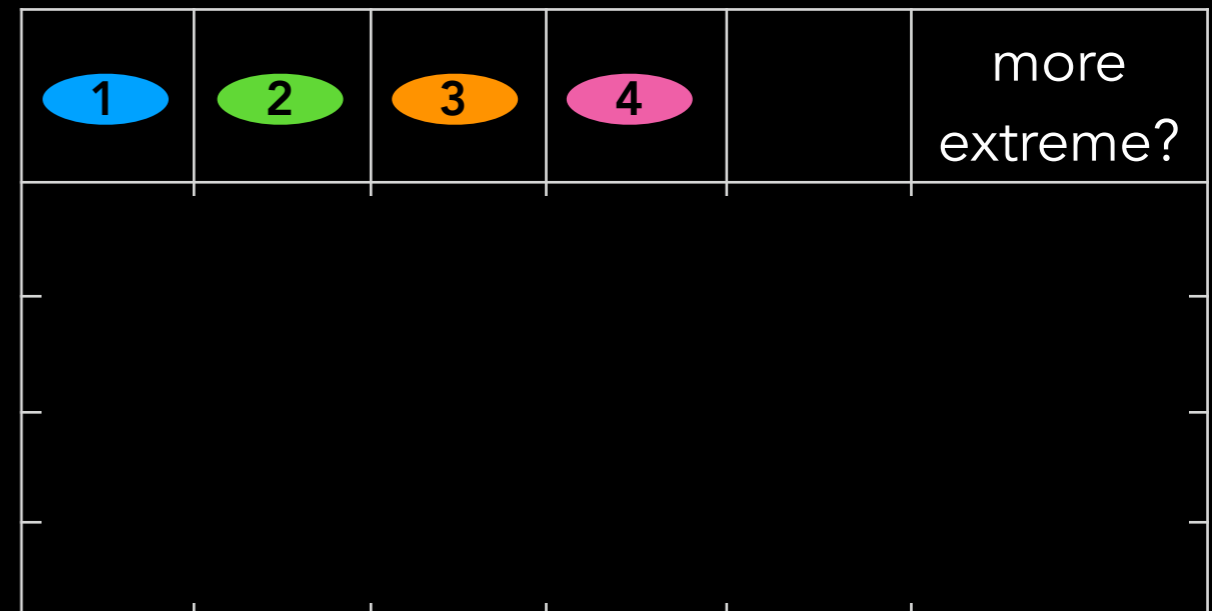
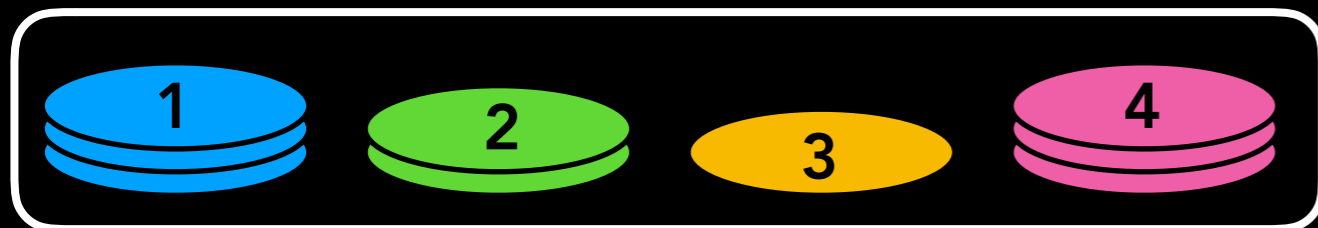
Bootstrapping: Simulating Many Samples From One

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

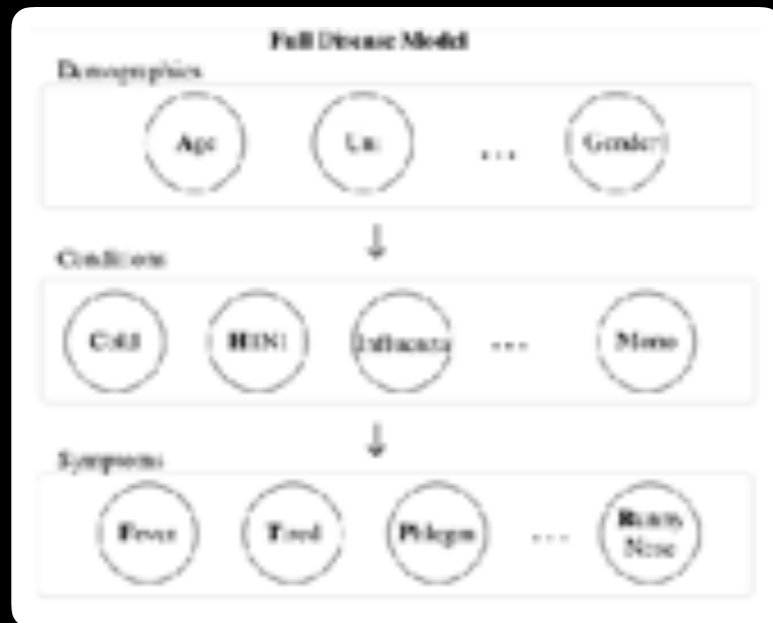
insight

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



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

we can also generate samples where we hold some values fixed (MCMC)

CS 109

topics

machine learning

parameter estimation

classifiers

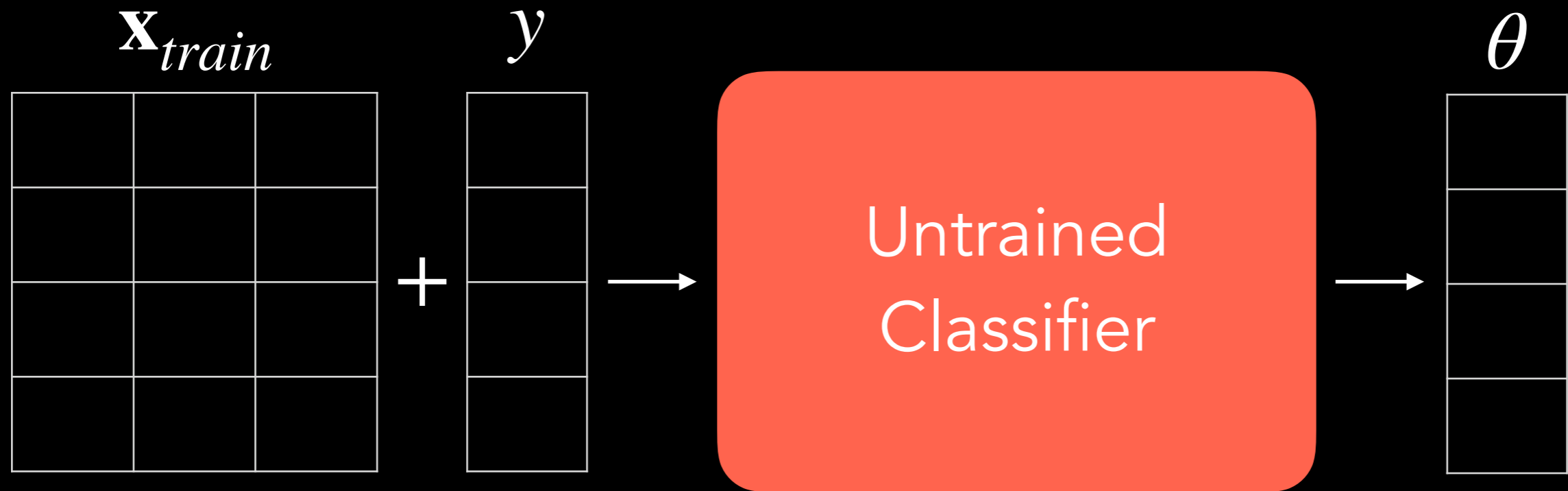
deep learning

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

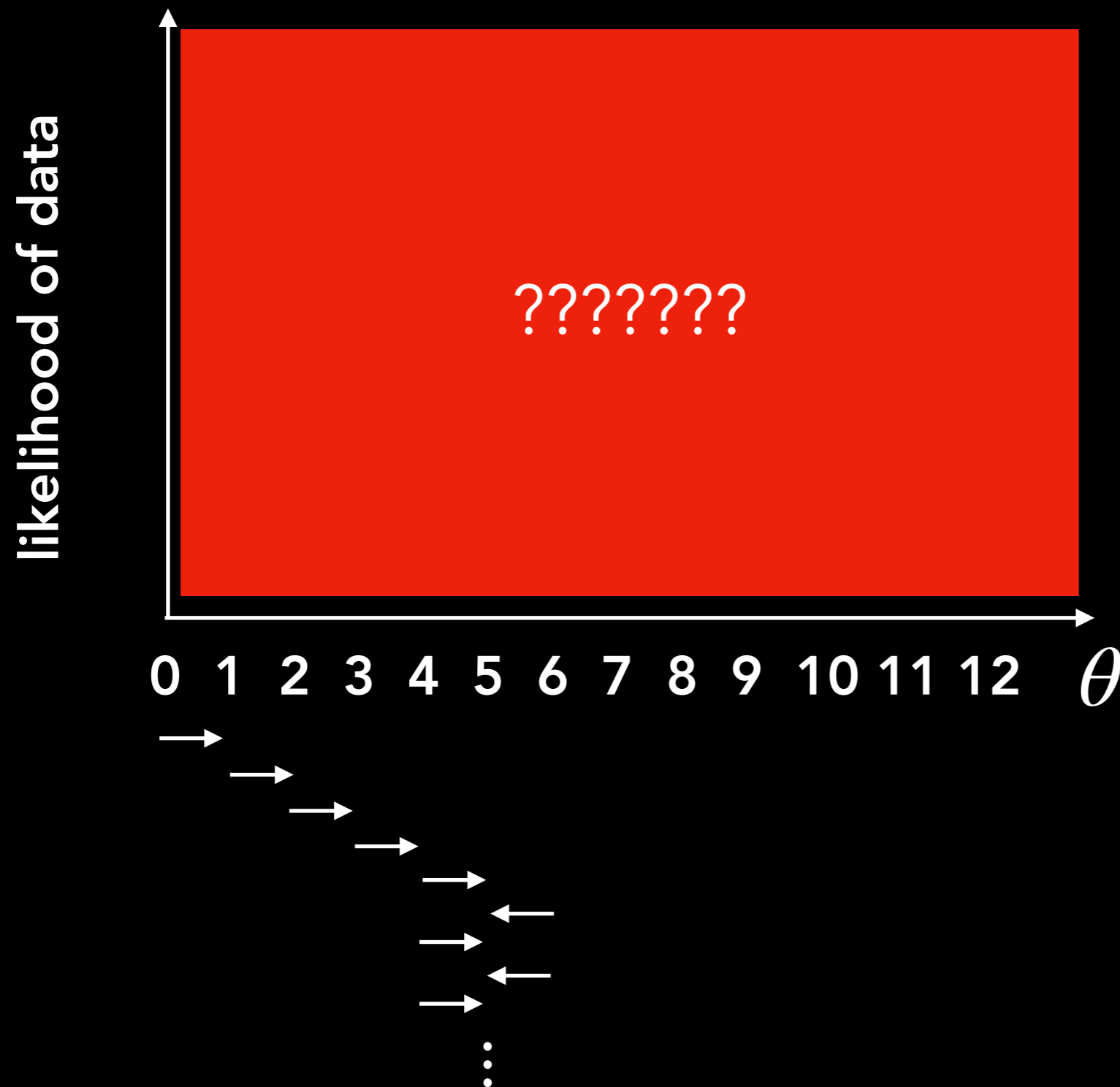
MLE:

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

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

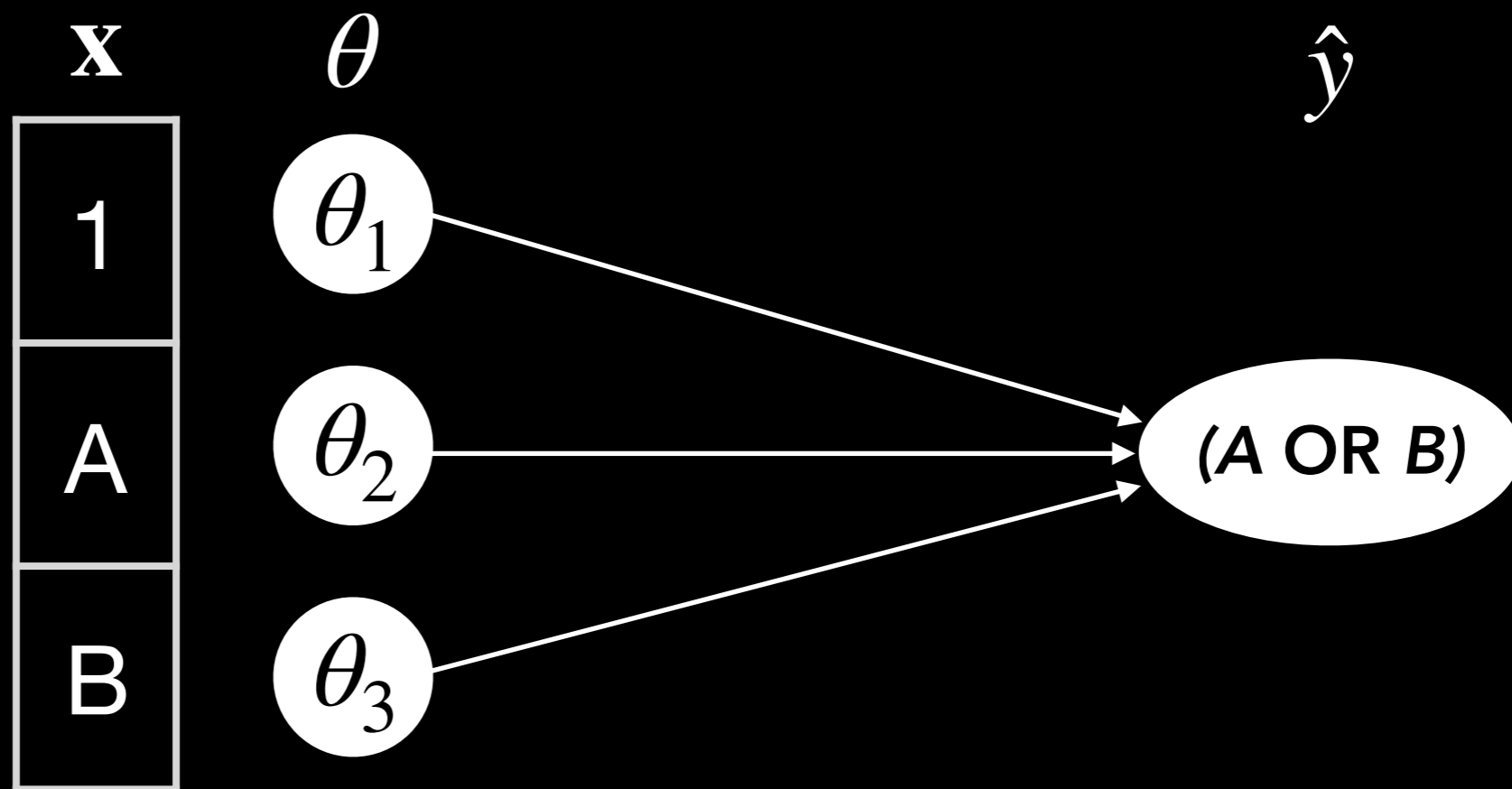
$$= \text{sign} \left[\frac{\partial \text{prob}}{\partial \theta} \right]$$

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

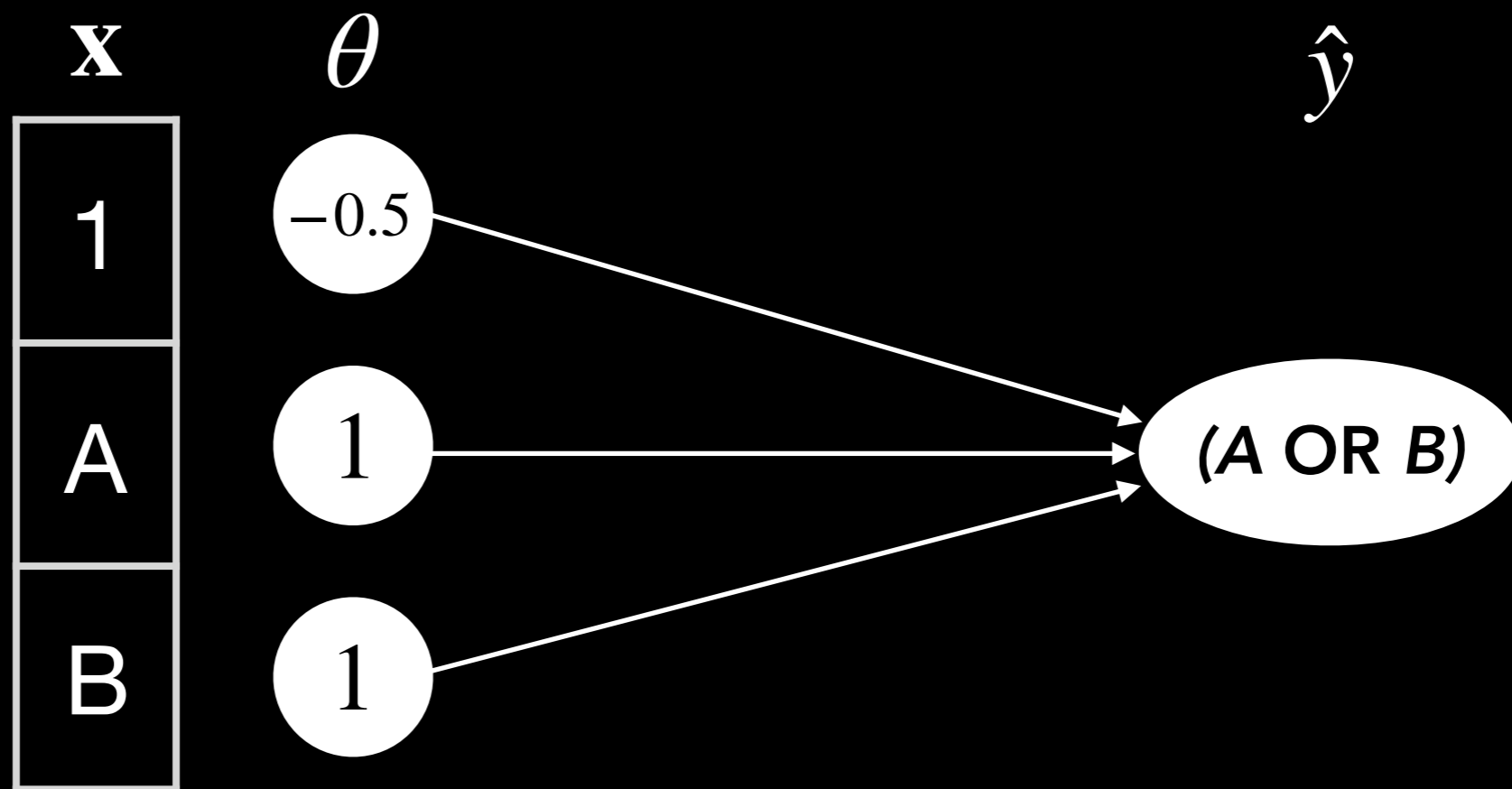
one neuron (logistic regression model)



What weights do we have to learn for θ_1 , θ_2 , θ_3 to perfectly classify data of the form $(A \text{ OR } B)$?

Neural Networks

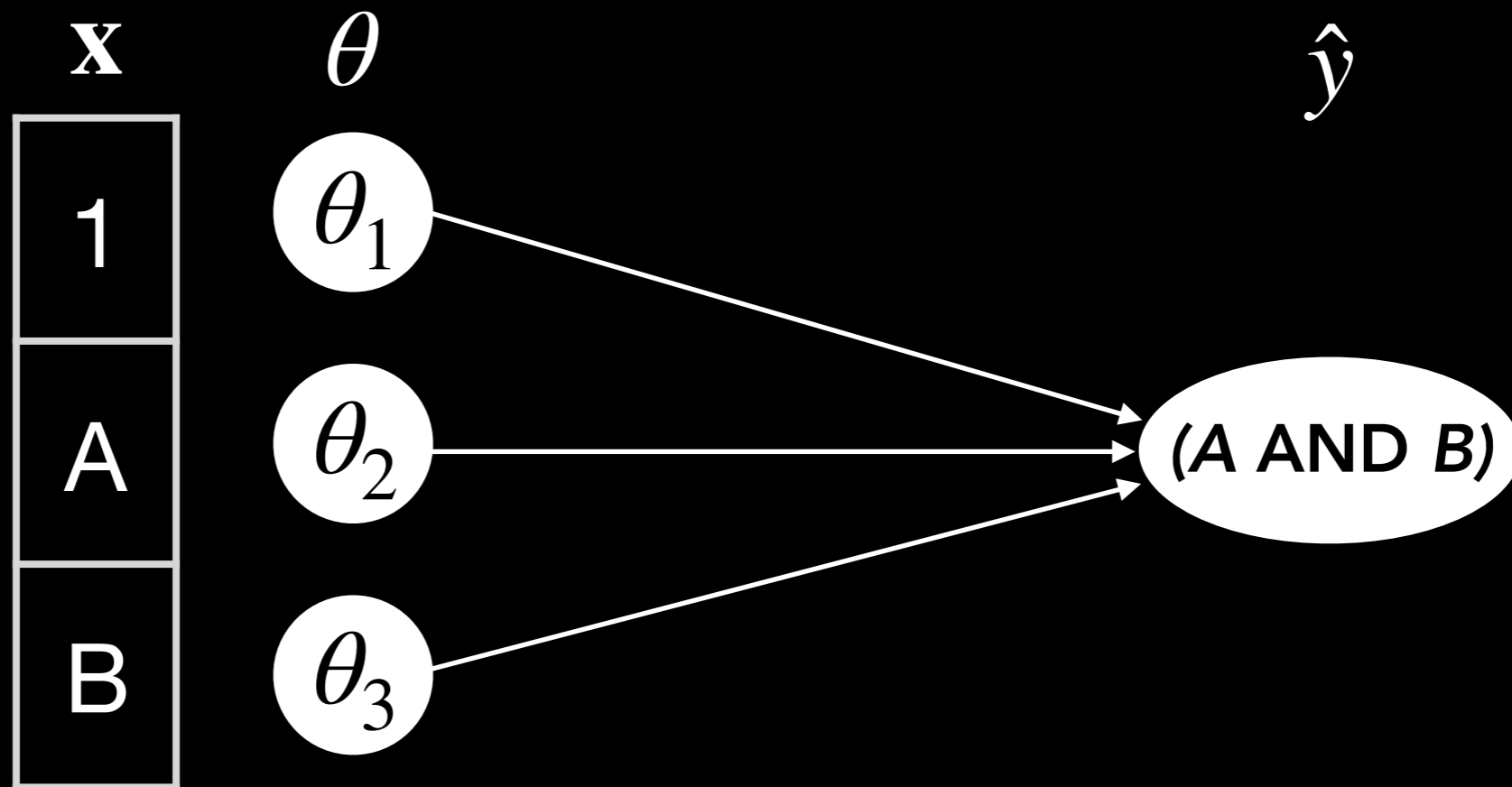
one neuron (logistic regression model)



What weights do we have to learn for $\theta_1, \theta_2, \theta_3$ to perfectly classify data of the form (A OR B)?

Neural Networks

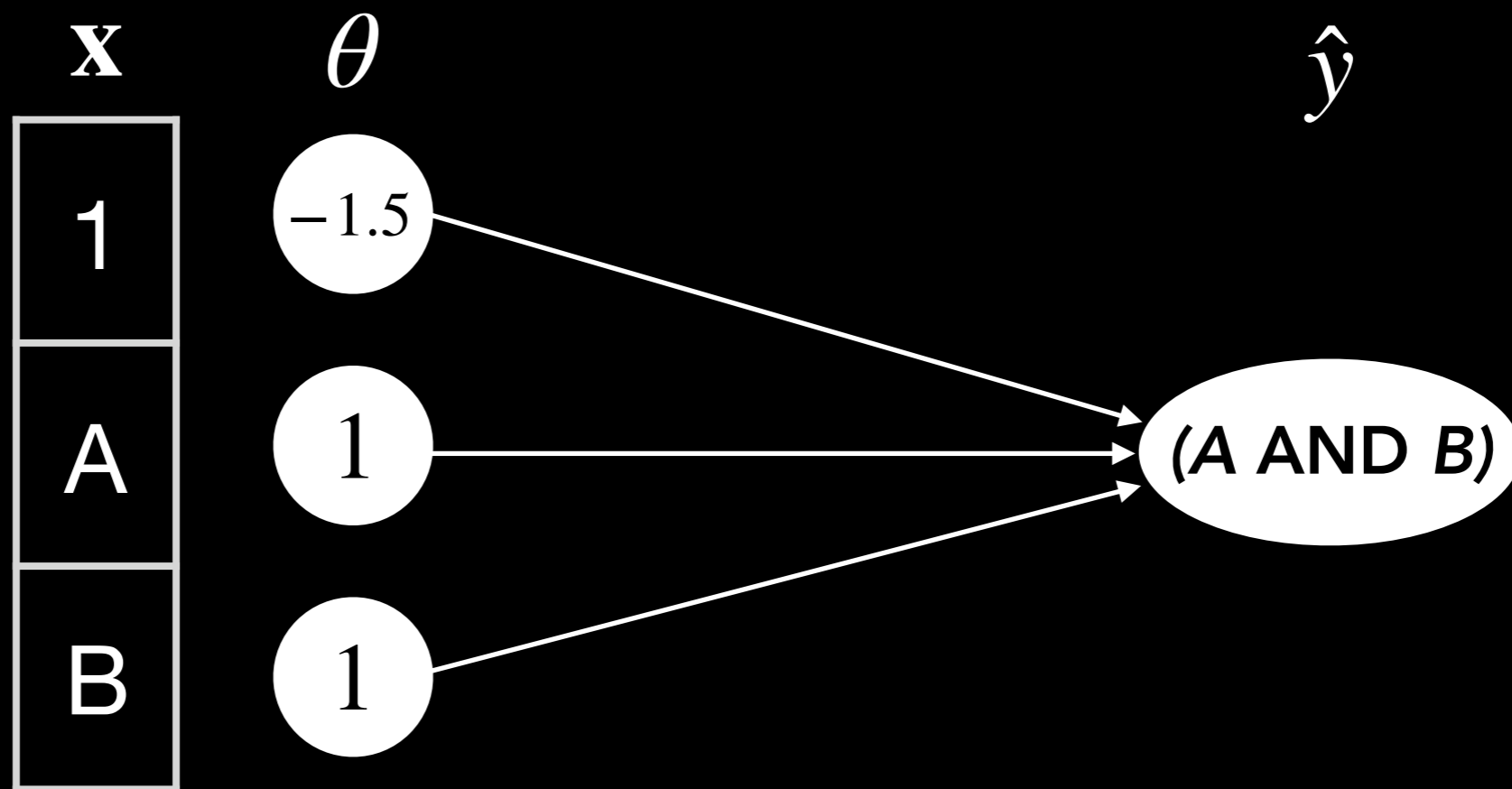
one neuron (logistic regression model)



What weights do we have to learn for θ_1 , θ_2 , θ_3 to perfectly classify data of the form (A AND B)?

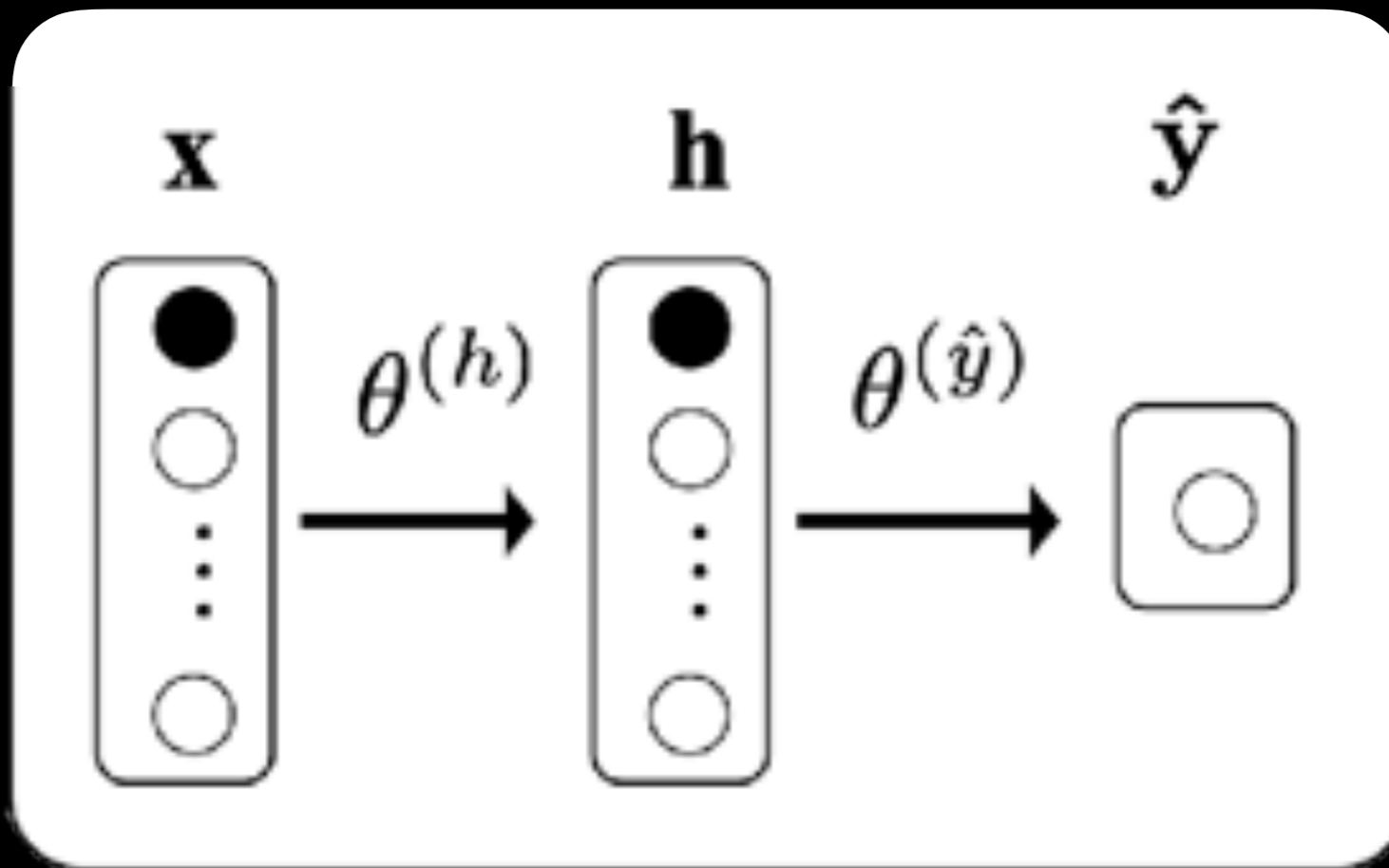
Neural Networks

one neuron (logistic regression model)



What weights do we have to learn for θ_1 , θ_2 , θ_3 to perfectly classify data of the form (A AND B)?

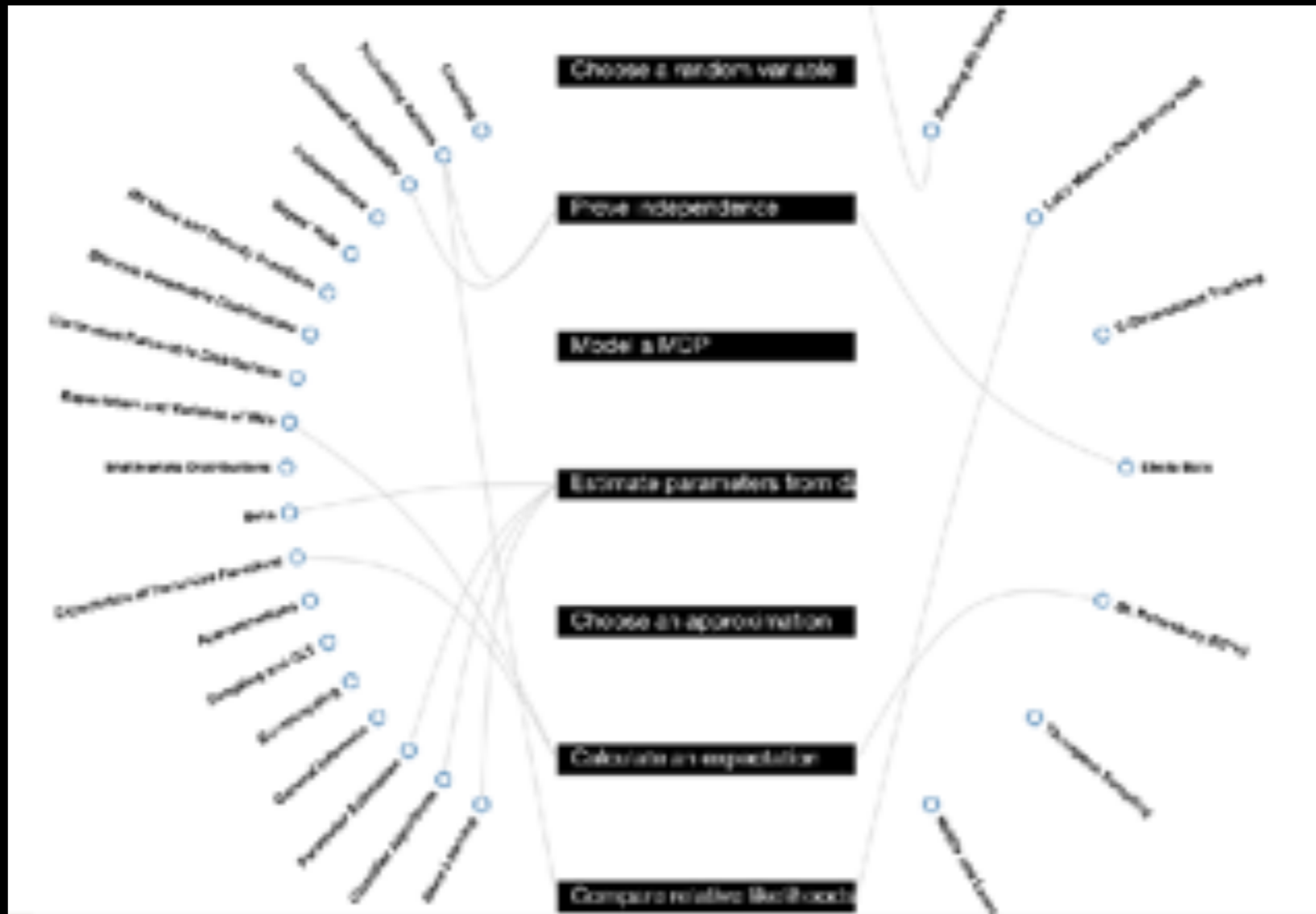
Neural Networks



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!*

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

Concept Organizer



Check out cs109.stanford.edu > Handouts > Big Picture! (live later tonight)

Good luck on the final!

