



CS 109 Review

Oishi Banerjee

March 11, 2020

Adapted from slides by Julia Daniel

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

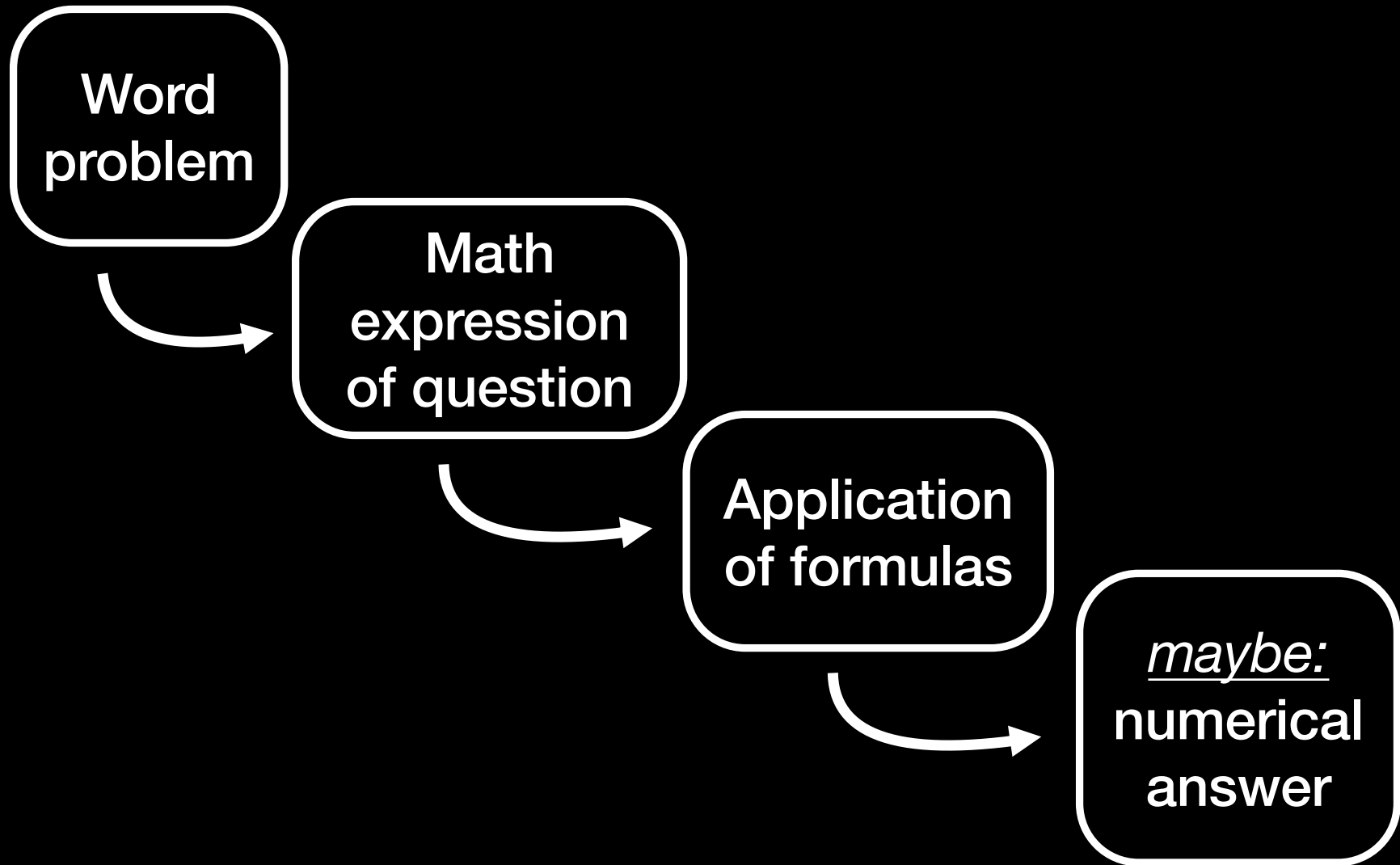
core probability fundamentals

skills

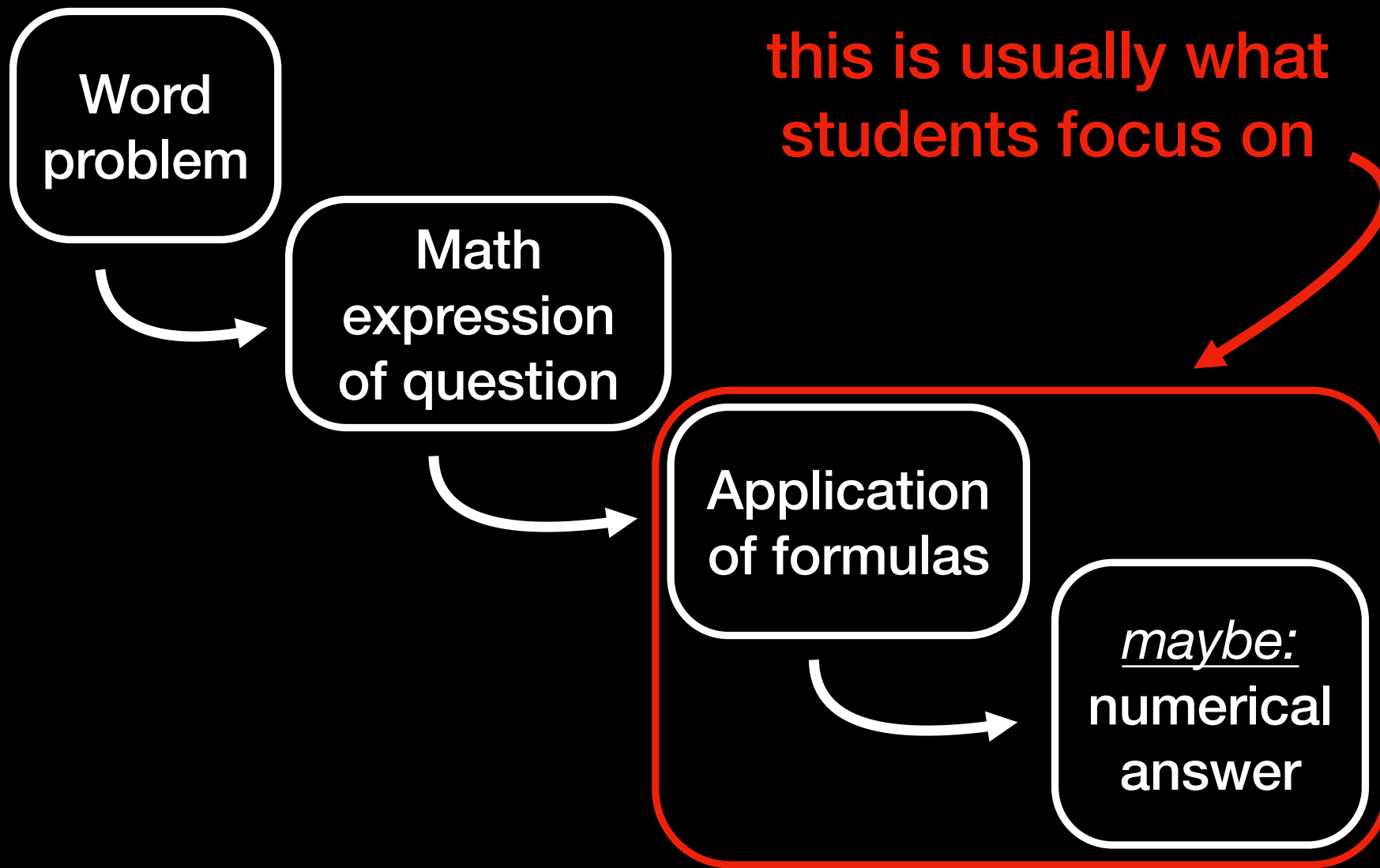
interpreting word problems into math

analyzing and producing code

Solving a CS109 problem



Solving a CS109 problem



Solving a CS109 problem

Word
problem

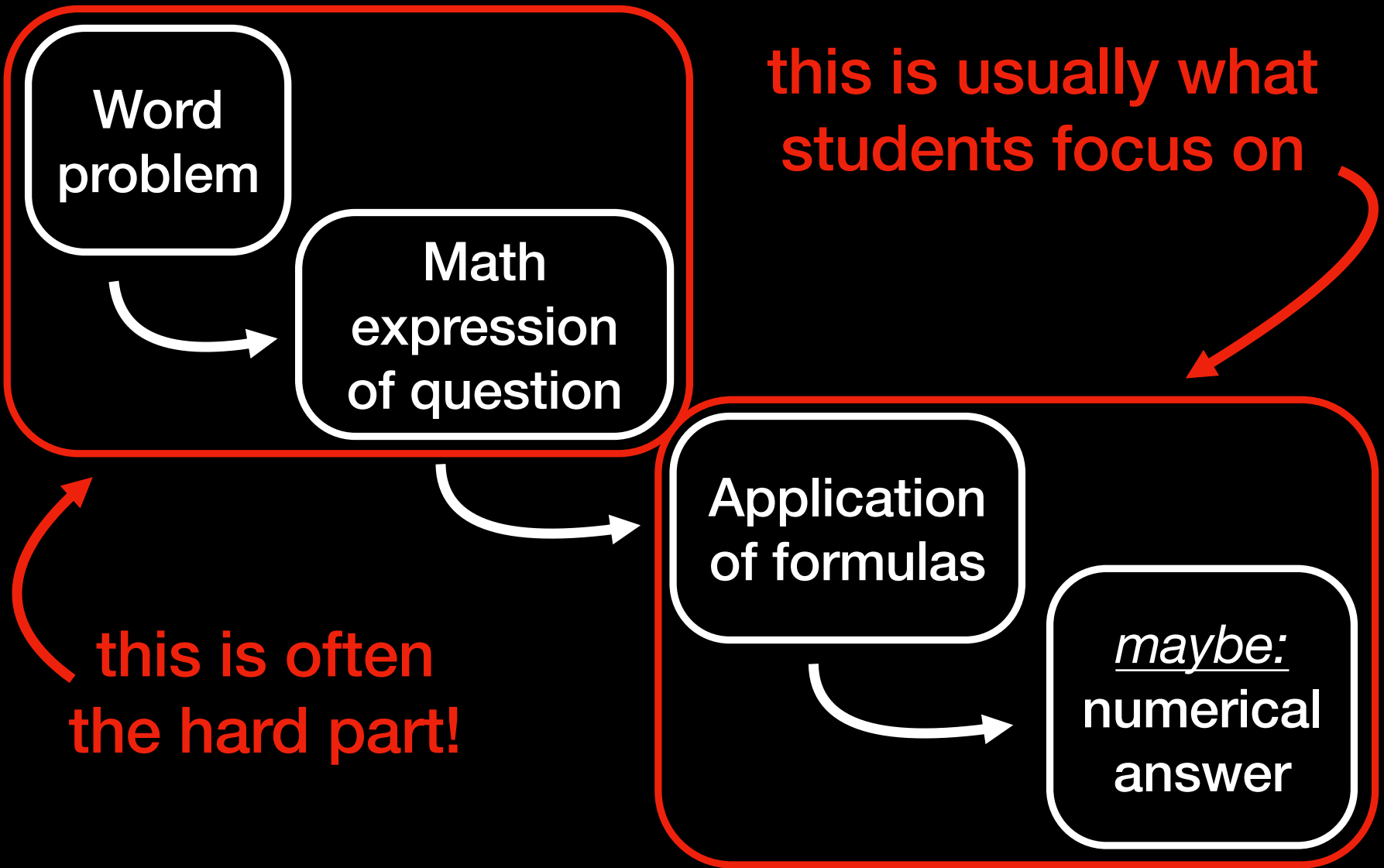
Math
expression
of question

this is usually what
students focus on

this is often
the hard part!

Application
of formulas

maybe:
numerical
answer



Step 1: Defining Your Terms

- What's a 'success'? What's the sample space?
- What does each random variable actually represent, in English? Every definition of an event or a random variable should have a **verb** in it. (' = ' is a verb)
- Make sure units match - particularly important for λ

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 _____”	____ ____
“at least _____”	<i>Flip it: 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!

2 Medical Testing [24 points]

Example

In medicine, there are many circumstances where we would like to detect the presence of a disease in a large population. Suppose that we would like to identify the number of individuals who have measles in a population of 1000 people using a blood test. The test is completely accurate: that is, if there are traces of measles in the blood sample, the test will return true 100% of the time and will otherwise return false. The probability that an individual has measles is 1% for everyone, independently of others.

- a. (8 points) Suppose that we use a blood test on each person, in order, for a total of 1000 blood tests. What is the probability that the tenth test is the first positive test (i.e., the first person we identify with measles)?

2 Medical Testing [24 points]

Example

In medicine, there are many circumstances where we would like to detect the presence of a disease in a large population. Suppose that we would like to identify the number of individuals who have measles in a population of 1000 people using a blood test. The test is completely accurate: that is, if there are traces of measles in the blood sample, the test will return true 100% of the time and will otherwise return false. The probability that an individual has measles is 1% for everyone, independently of others.

- a. (8 points) Suppose that we use a blood test on each person, in order, for a total of 1000 blood tests. What is the probability that the tenth test is the first positive test (i.e., the first person we identify with measles)?

(preamble)

1. 1000 people

2. $P(\text{test_positive} \mid \text{person_has_measles}) = 1$

3. $P(\text{person_has_measles}) = 0.01$

(part a)

1. 1000 independent tests

2. $P(\text{tests 1-9 negative and test 10 is positive})$

2 Medical Testing [24 points]

Example

In medicine, there are many circumstances where we would like to detect the presence of a disease in a large population. Suppose that we would like to identify the number of individuals who have measles in a population of 1000 people using a blood test. The test is completely accurate: that is, if there are traces of measles in the blood sample, the test will return true 100% of the time and will otherwise return false. The probability that an individual has measles is 1% for everyone, independently of others.

- a. (8 points) Suppose that we use a blood test on each person, in order, for a total of 1000 blood tests. What is the probability that the tenth test is the first positive test (i.e., the first person we identify with measles)?

(preamble)

1. 1000 people

2. $P(\text{test_positive} \mid \text{person_has_measles}) = 1$

3. $P(\text{person_has_measles}) = 0.01$

(part a)

1. 1000 independent tests

2. $P(\text{tests 1-9 negative and test 10 is positive})$

What is actually important?

2 Medical Testing [24 points]

Example

In medicine, there are many circumstances where we would like to detect the presence of a disease in a large population. Suppose that we would like to identify the number of individuals who have measles in a population of 1000 people using a blood test. The test is completely accurate: that is, if there are traces of measles in the blood sample, the test will return true 100% of the time and will otherwise return false. The probability that an individual has measles is 1% for everyone, independently of others.

- a. (8 points) Suppose that we use a blood test on each person, in order, for a total of 1000 blood tests. What is the probability that the tenth test is the first positive test (i.e., the first person we identify with measles)?

(part a solution)

1.Independent tests/trials

2.P(test positive) = 0.01

3.P(10 trials until “success”)

4.Implies use Geometric

5.Answer is: $(0.99^9)*(0.01^1)$

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

Code Analysis

Expectation of
binary tree depth
("recursive" expectation)

Bloom Filter Analysis

Expectation of
recursive die roll game

Coding Applications

Dithering

CO2 Levels

Biometric Keystrokes

Titanic

Peer Grading

Thompson Sampling

Code in CS 109

Example

```
int fairRandom() {
    int r1, r2;
    while (true) {
        r1 = unknownRandom();
        r2 = unknownRandom();
        if (r1 != r2) break;
    }
    return r2;
}
```

- a. Show mathematically that `fairRandom` does indeed return a 0 or a 1 with equal probability.

Code in CS 109

Example

```
int fairRandom() {
    int r1, r2;
    while (true) {
        r1 = unknownRandom();
        r2 = unknownRandom();
        if (r1 != r2) break;
    }
    return r2;
}
```

a. Show mathematically that `fairRandom` does indeed return a 0 or a 1 with equal probability.

1. Need to prove $P(\text{fairR}() = 1) = 0.5$
2. We know the function returns, so we break into cases:
3. Case 1: $r1 = 1$ and $r2 = 0 \Rightarrow$ likelihood $p \cdot (1-p)$
4. Case 2: $r1 = 0$ and $r2 = 1 \Rightarrow$ likelihood $p \cdot (1-p)$
5. These are equal \Rightarrow equally likely for $r2$ to be 0 or 1

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 $ceiling(m / n)$ objects.

Combinatorics: Arranging Items

Are your objects all distinct from each other?



All indistinct?



Or is there a mix?



Do you care about the order of your objects? Should you consider



Combinatorics: Arranging Items

ordering n distinct objects

$$n!$$

choosing k out of n distinct objects
to go in an unordered group

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

ordering n objects when some are
indistinct from each other

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

arranging n indistinct objects in r
containers

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

Probability basics

$$P(E) = \lim_{x \rightarrow \infty} \frac{n(E)}{n}$$

in the general case

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)

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

Probability basics

Example

7. If we assume that all possible poker hands (comprised of 5 cards from a standard 52 card deck) are equally likely, what is the probability of being dealt:
- a flush? (A hand is said to be a flush if all 5 cards are of the same suit. Note that this definition means that *straight flushes* (five cards of the same suit in numeric sequence) are also considered flushes.)
 - two pairs? (This occurs when the cards have numeric values a, a, b, b, c , where a, b and c are all distinct.)
 - three of a kind? (This occurs when the cards have numeric values a, a, a, b, c , where a, b and c are all distinct.)

Part a:

Probability basics

Example

7. If we assume that all possible poker hands (comprised of 5 cards from a standard 52 card deck) are equally likely, what is the probability of being dealt:
- a flush? (A hand is said to be a flush if all 5 cards are of the same suit. Note that this definition means that *straight flushes* (five cards of the same suit in numeric sequence) are also considered flushes.)
 - two pairs? (This occurs when the cards have numeric values a, a, b, b, c , where a, b and c are all distinct.)
 - three of a kind? (This occurs when the cards have numeric values a, a, a, b, c , where a, b and c are all distinct.)

Part a:

- 1. Hand rearrangement OK \Rightarrow use unordered sample space**
- 2. Sample space $\Rightarrow 52C5$**
- 3. For event space: choose suit, choose cards $\Rightarrow 4C1 * 13C5$**
- 4. Put it together: $P(\text{a flush}) = 4C1 * 13C5 / 52C5$**

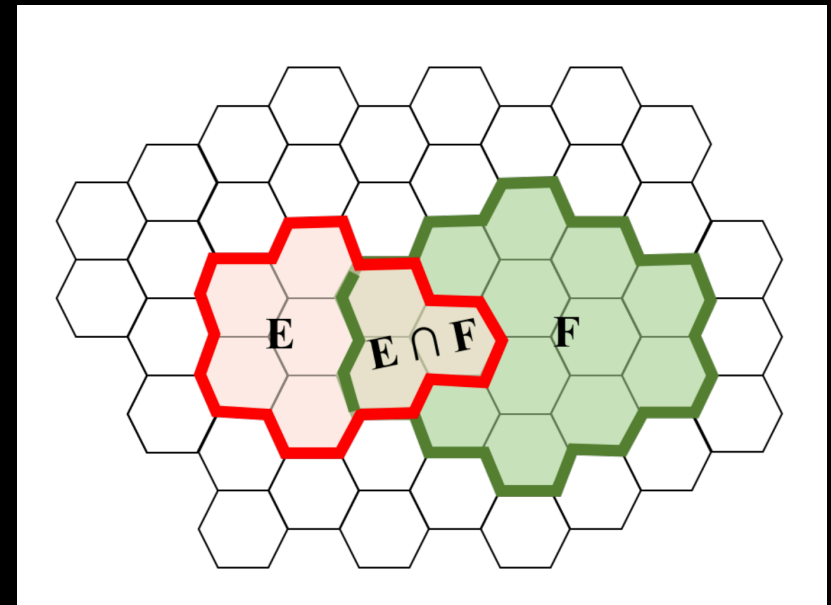
Conditional Probability

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

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

Event W = we walk to class. Event $B = W^C$.

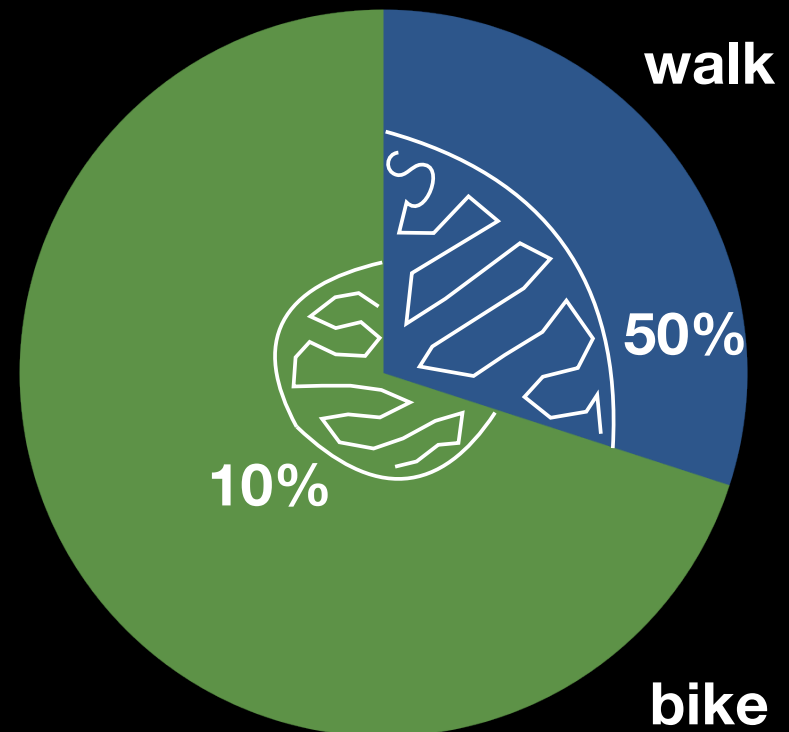
Event L = we are late to class.

$P(L | W) = 0.5$, $P(L | B) = 0.1$.

$P(W) = 0.3$.

$P(L) = ?$

total shaded = ?%
of whole



Law of Total Probability

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

Event W = we walk to class. Event B = we bike = W^C .

Event L = we are late to class.

$P(L | W) = 0.5$, $P(L | B) = 0.1$.

$P(W) = 0.3$.

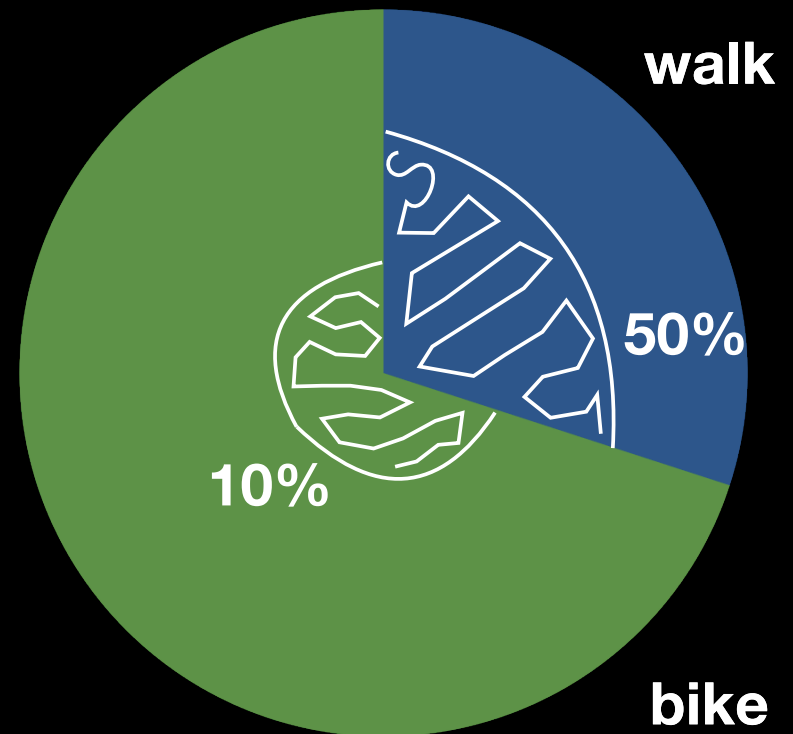
$P(L) = ?$

**total shaded = ?%
of whole**

$$P(L) = P(L | W)P(W) + P(L | W^C)P(W^C)$$

$$= (0.5)(0.3) + (0.1)(0.7)$$

$$= 0.22$$



Law of Total Probability

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

Event W = we walk to class. Event B = we bike = W^C .

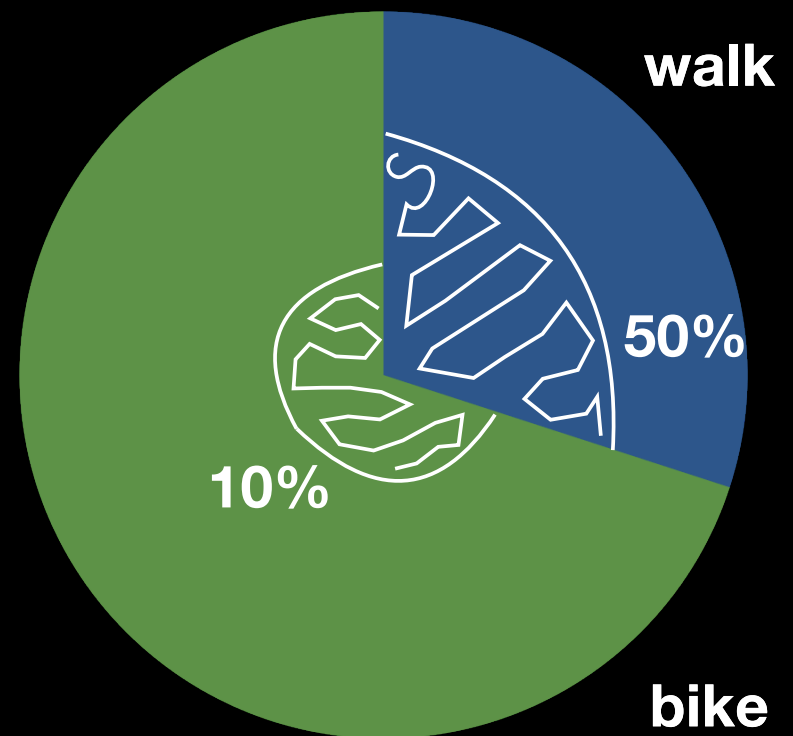
Event L = we are late to class.

$P(L | W) = 0.5$, $P(L | B) = 0.1$.

$P(W) = 0.3$.

$P(L) = ?$

what if we can bike, walk, or
take the Marguerite (> 2 options)?



Law of Total Probability

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

Event W = we walk to class. Event B = we bike = W^C .

Event L = we are late to class.

$P(L | W) = 0.5$, $P(L | B) = 0.1$.

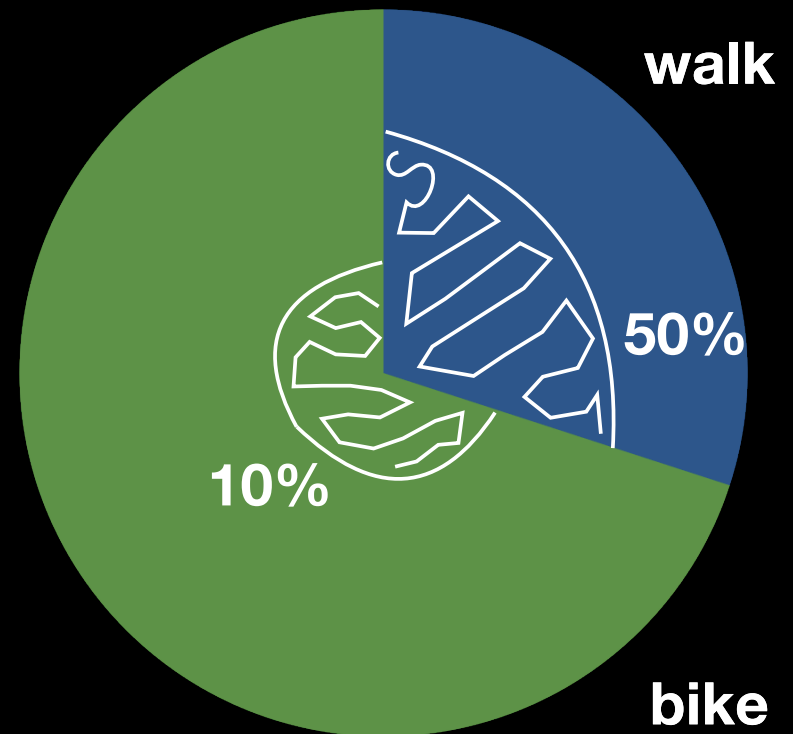
$P(W) = 0.3$.

$P(L) = ?$

what if we can bike, walk, or
take the Marguerite (> 2 options)?

events must be:

- **mutually exclusive**, and
- **exhaustive**



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 term $P(E|F)$, with a pink arrow pointing to it. The word "likelihood" is written in pink above the term $P(F|E)$, with a pink arrow pointing to it. The word "prior" is written in pink above the term $P(E)$, with a pink arrow pointing to it. The word "normalization constant" is written in pink below the term $P(F)$, with a pink arrow pointing to it. The entire equation is rendered in white on a black background.

Bayes' Rule

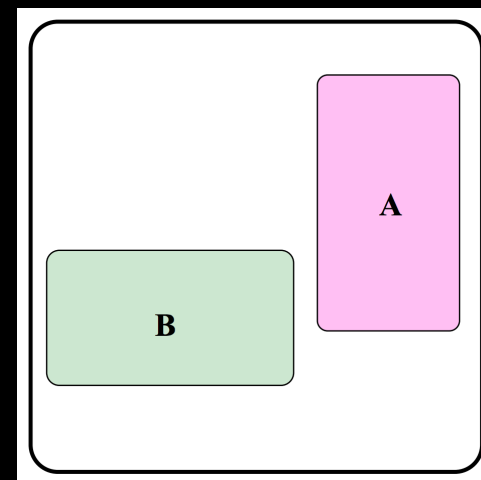
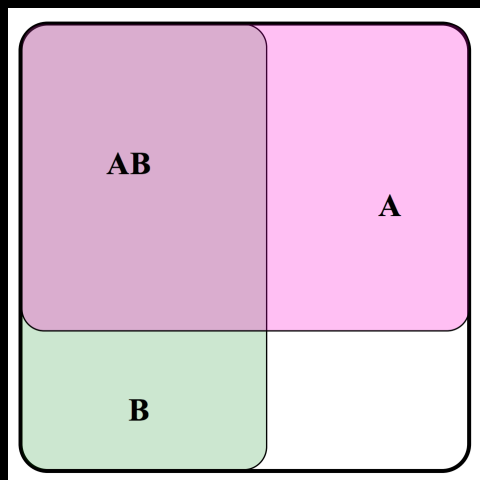
$$P(E|F) = \frac{P(F|E)P(E)}{P(F)}$$
$$P(F|E)P(E) + P(F|E^C)P(E^C)$$

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

Old Principles, New Tricks

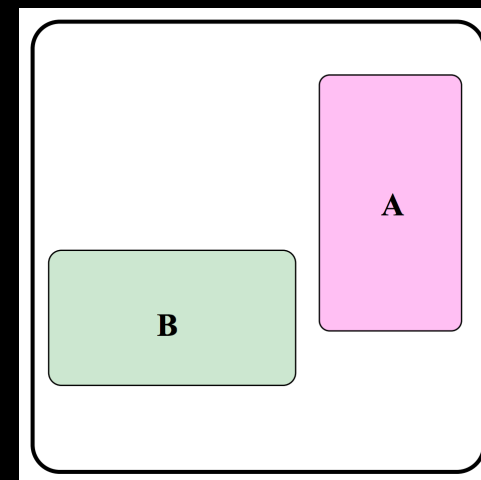
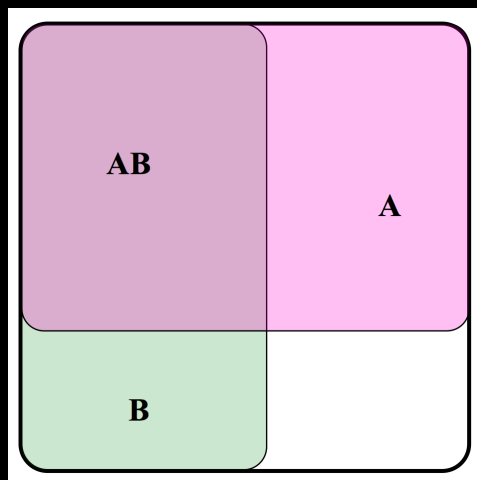
Name of Rule	Original Rule	Conditional Rule
First axiom of probability	$0 \leq P(E) \leq 1$	$0 \leq P(E G) \leq 1$
Complement Rule	$P(E) = 1 - P(E^C)$	$P(E G) = 1 - P(E^C G)$
Chain Rule	$P(EF) = P(E F)P(F)$	$P(EF G) = P(E FG)P(F G)$
Bayes Theorem	$P(E F) = \frac{P(F E)P(E)}{P(F)}$	$P(E FG) = \frac{P(F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

Independence



Independence

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



Independence

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

If E and F are independent.....

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

& vice versa

Example

Beyond the basics

7. Consider a hash table with 15 buckets, of which 9 are empty (have no strings hashed to them) and the other 6 buckets are non-empty (have at least one string hashed to each of them already). Now, 2 new strings are independently hashed into the table, where each string is equally likely to be hashed into any bucket. Later, another 2 strings are hashed into the table (again, independently and equally likely to get hashed to any bucket). What is the probability that both of the final 2 strings are each hashed to empty buckets in the table?

Example

Beyond the basics

7. Consider a hash table with 15 buckets, of which 9 are empty (have no strings hashed to them) and the other 6 buckets are non-empty (have at least one string hashed to each of them already). Now, 2 new strings are independently hashed into the table, where each string is equally likely to be hashed into any bucket. Later, another 2 strings are hashed into the table (again, independently and equally likely to get hashed to any bucket). What is the probability that both of the final 2 strings are each hashed to empty buckets in the table?

How do you begin to break down this problem?

Example

Beyond the basics

7. Consider a hash table with 15 buckets, of which 9 are empty (have no strings hashed to them) and the other 6 buckets are non-empty (have at least one string hashed to each of them already). Now, 2 new strings are independently hashed into the table, where each string is equally likely to be hashed into any bucket. Later, another 2 strings are hashed into the table (again, independently and equally likely to get hashed to any bucket). What is the probability that both of the final 2 strings are each hashed to empty buckets in the table?

How do you begin to break down this problem?

Let event A = first of initial two strings hashed to empty bucket.
Let event B = second of initial two strings hashed to empty bucket.
Let event C = first of final two strings hashed to empty bucket.
Let event D = second of final two strings hashed to empty bucket.

Define events

Beyond the basics

Example

7. Consider a hash table with 15 buckets, of which 9 are empty (have no strings hashed to them) and the other 6 buckets are non-empty (have at least one string hashed to each of them already). Now, 2 new strings are independently hashed into the table, where each string is equally likely to be hashed into any bucket. Later, another 2 strings are hashed into the table (again, independently and equally likely to get hashed to any bucket). What is the probability that both of the final 2 strings are each hashed to empty buckets in the table?

How do you begin to break down this problem?

Let event A = first of initial two strings hashed to empty bucket.
Let event B = second of initial two strings hashed to empty bucket.
Let event C = first of final two strings hashed to empty bucket.
Let event D = second of final two strings hashed to empty bucket.

We compute $P(CD)$ as follows:

$$\begin{aligned} P(CD) &= P(CD \mid AB)P(AB) + P(CD \mid A^C B)P(A^C B) + \\ &\quad P(CD \mid AB^C)P(AB^C) + P(CD \mid A^C B^C)P(A^C B^C) \\ &= (7/15)(6/15)(9/15)(8/15) + (8/15)(7/15)(6/15)(9/15) + \\ &\quad (8/15)(7/15)(9/15)(7/15) + (9/15)(8/15)(6/15)(6/15) \\ &= 12168/154 \approx 0.2404 \end{aligned}$$

Define events

What is the question asking?

Use LOTP

CS 109

topics

machine learning

sampling, making conclusions from data

random variables / distributions

core probability fundamentals

multivariate distributions

discrete RVs

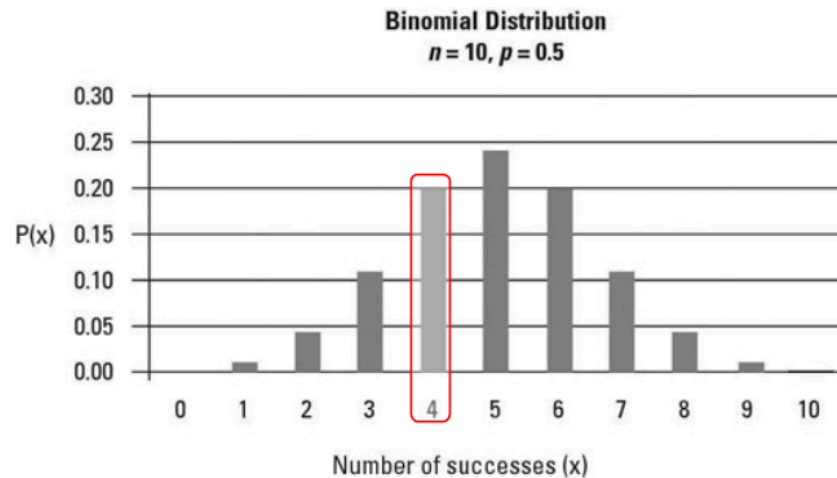
continuous RVs

properties of RVs

Probability Distributions

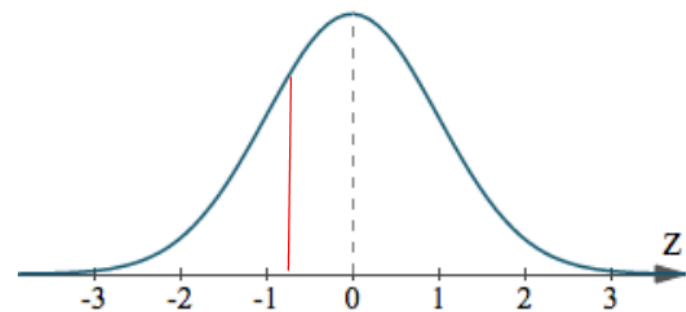
Discrete

PMF:



Continuous

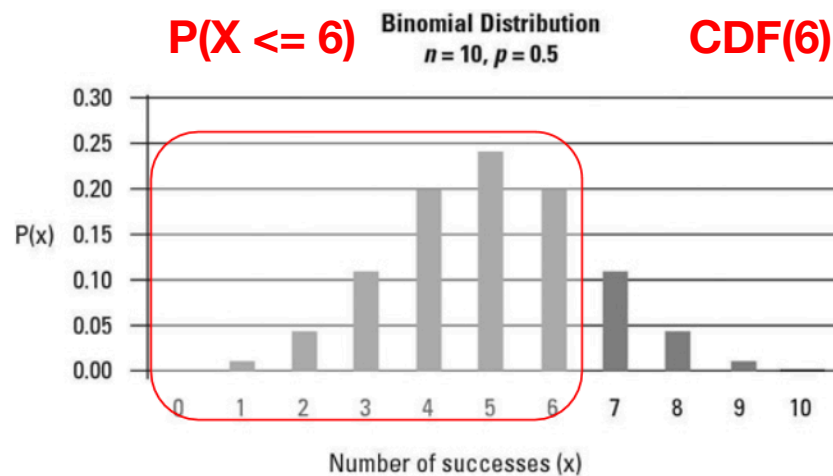
PDF:



Probability Distributions

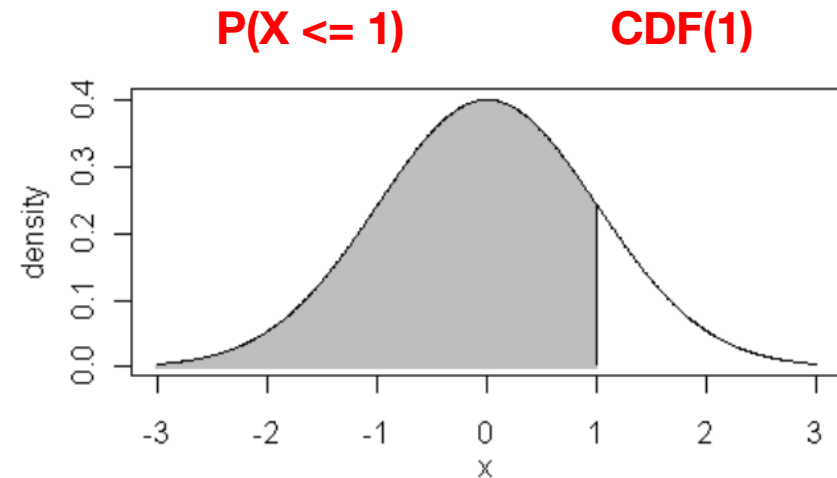
Discrete

CDF: $P(X \leq x)$



Continuous

CDF: $P(X \leq x)$



Expectation & Variance

Discrete definition

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

Continuous definition

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

Expectation & Variance

Discrete definition

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

(mu)

Continuous definition

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

(mu)

Properties of Expectation

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

Expectation & Variance

Discrete definition

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

(mu)

Continuous definition

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

(mu)

Properties of Expectation

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

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

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

Expectation & Variance

Discrete definition

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

(mu)

Continuous definition

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

(mu)

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

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

(mu)

Expectation & Variance

Discrete definition

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

(mu)

Continuous definition

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

(mu)

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

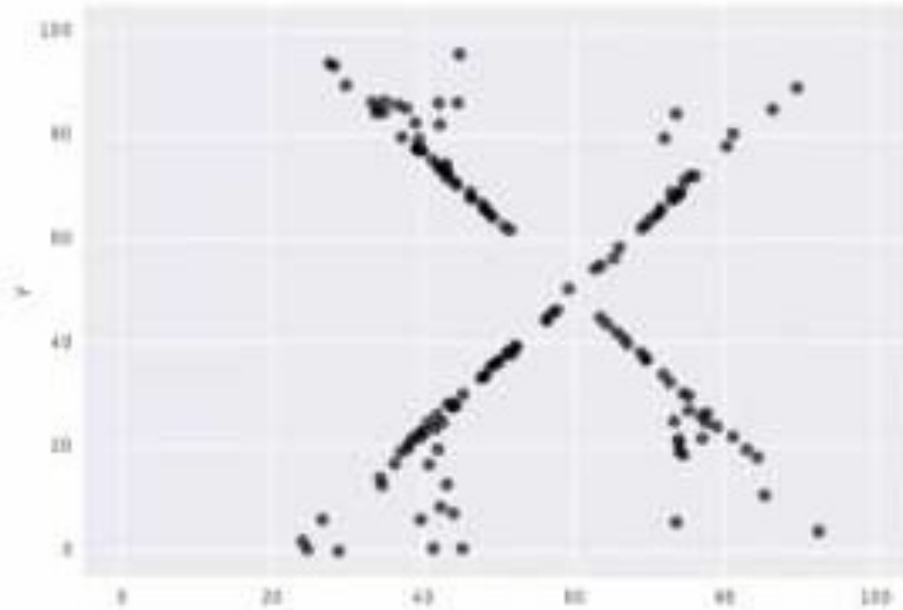
$$Var(X) = E[(X - \mu)^2]$$

(mu)

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

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

Extras



X Mean: 54.2601949
Y Mean: 47.8388784
X SD : 16.7699928
Y SD : 26.9300128
Corr. : -0.0615421

All our (discrete) friends

Ber(p)	Bin(n, p)	Poi(λ)	Geo(p)	NegBin(r, p)
$P(X) = p$	$\binom{n}{k} p^k (1-p)^{n-k}$	$\frac{\lambda^k e^{-\lambda}}{k!}$	$(1-p)^{k-1} p$	$\binom{k-1}{r-1} p^r (1-p)^{k-r}$
Getting candy or not at a random house				

All our (discrete) friends

Ber(p)	Bin(n, p)	Poi(λ)	Geo(p)	NegBin(r, p)
$P(X) = p$	$\binom{n}{k} p^k (1-p)^{n-k}$	$\frac{\lambda^k e^{-\lambda}}{k!}$	$(1-p)^{k-1} p$	$\binom{k-1}{r-1} p^r (1-p)^{k-r}$
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 (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)
$f(x) = \frac{1}{\beta - \alpha}$	$f(x) = \lambda e^{-\lambda x}$	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
$P(a \leq X \leq b) = \frac{b - a}{\beta - \alpha}$	$F(x) = 1 - e^{-\lambda x}$	$F(x) = \Phi\left(\frac{x - \mu}{\sigma}\right)$
thickness of sidewalk pavement between houses		

All our (continuous) friends

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

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

Approximations

When can we approximate a binomial?

n is large

Binomial

Var > 10

Var < 10

Normal

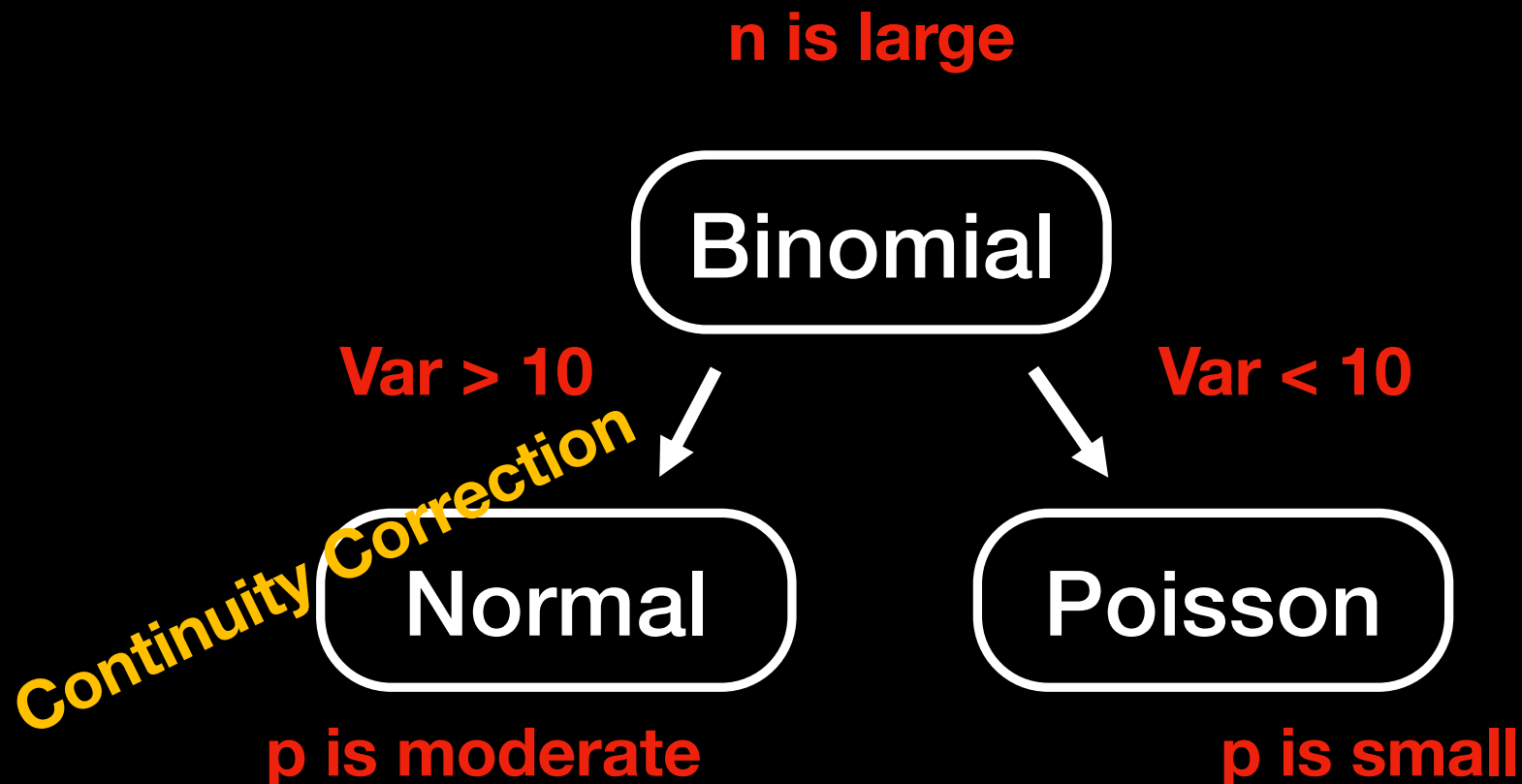
Poisson

p is moderate

p is small

Approximations

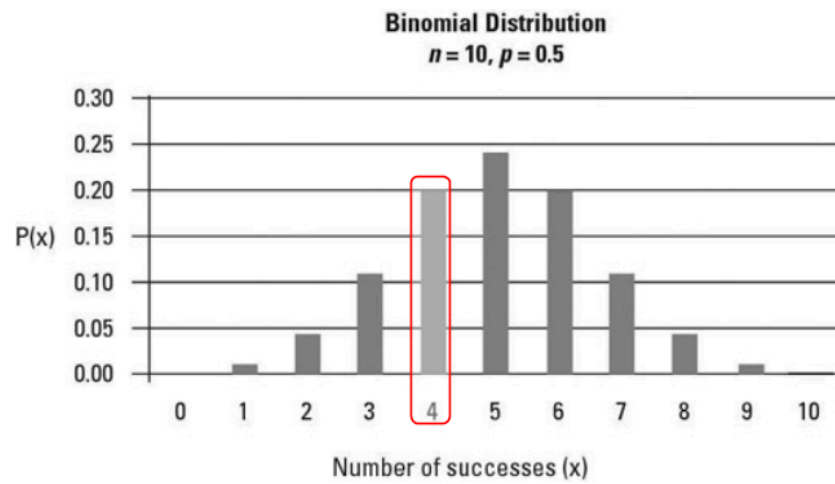
When can we approximate a binomial?



Continuity Correction

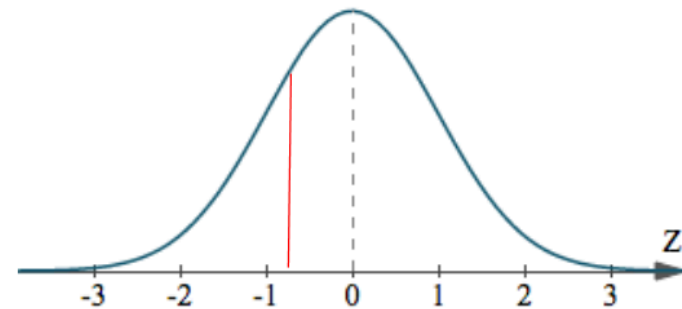
Discrete

PMF:



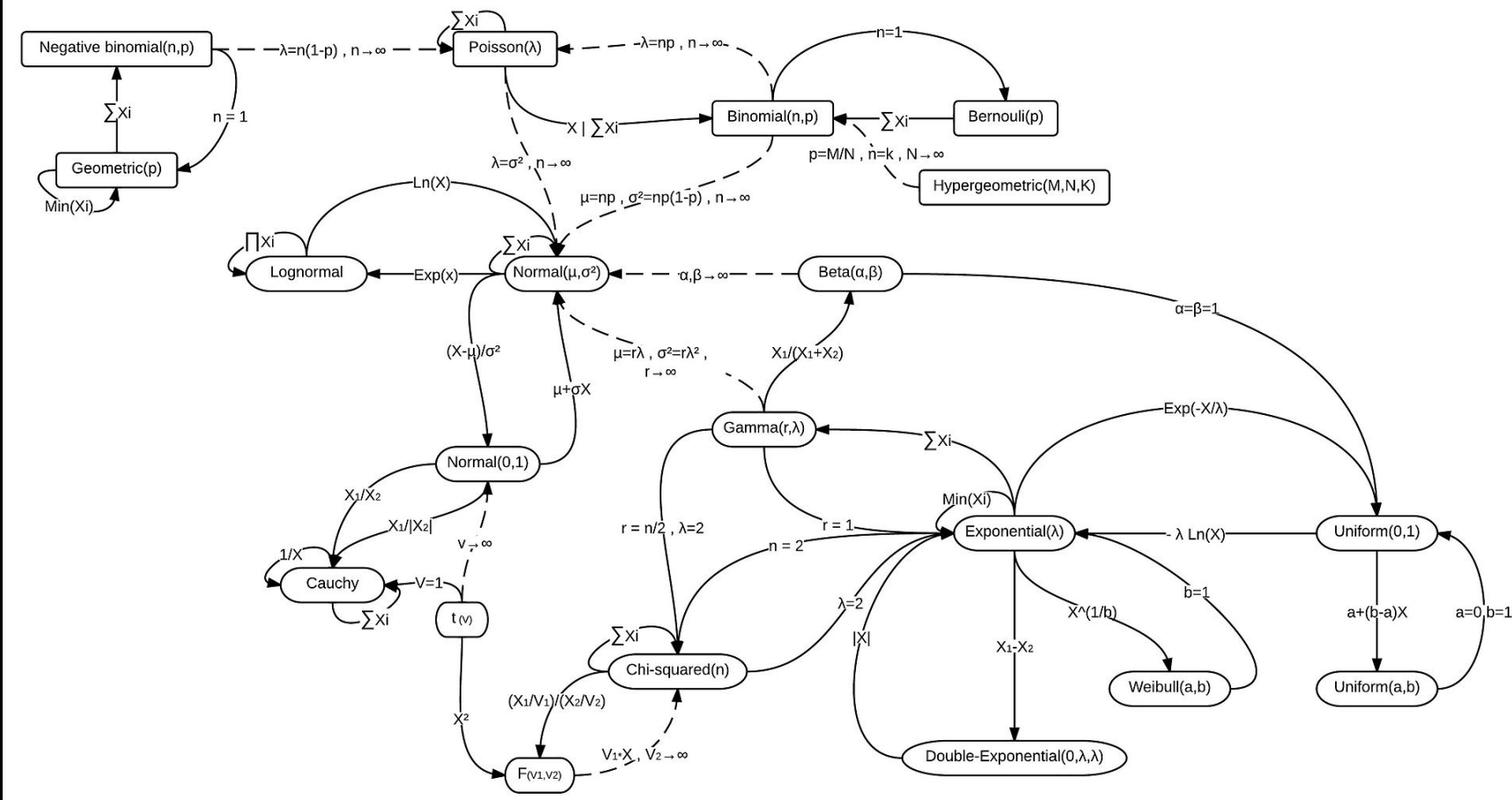
Continuous

PDF:



Only applies to **continuous RVs approximating discrete RVs** - why?

Extras



Distribution onslaught!

Example

Coin flip is heads

Distribution onslaught!

Example

Number of heads in 10 coin flips

Distribution onslaught!

Example

Coin flips until a heads

Distribution onslaught!

Example

Chance of CS109 student sleeping in class is 70%

Number of CS109 students sleeping in class right now?

Distribution onslaught!

Example

Chance of CS109 student sleeping is 70%

Number of CS109 students sleeping right now?
(approximate)

Distribution onslaught!

Example

CS109 students fall asleep on average once a minute.

Time until a CS109 student falls asleep?

Distribution onslaught!

Example

CS109 students fall asleep on average once a minute.

Number of CS109 students who fall asleep in the next 10 minutes?

Joint Distributions

- Discrete case:

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

- Marginalize a variable out:

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

- Continuous case:

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

- Marginalize a variable out:

- For joint distributions to be independent, both their joint probability density functions must be **factorable** and the bounds of the variables must be **separable**.

Joint Distributions

- Discrete case: $p_{x,y}(a, b) = P(X = a, Y = b)$

- Marginalize a variable out: $P_x(a) = \sum_y P_{x,y}(a, y)$

- Continuous case:

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

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

- For joint distributions to be independent, both their joint probability density functions must be **factorable** and the bounds of the variables must be **separable**.

Sums of Indep. RVs

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

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

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

$$f_{X+Y}(a) = \int_{y=-\infty}^{\infty} f_X(a-y)f_Y(y)dy \quad \text{(general case)}$$

Sums of Indep. RVs

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

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

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

$$f_{X+Y}(a) = \int_{y=-\infty}^{\infty} f_X(a-y)f_Y(y)dy \quad \text{(general case)}$$

Caveat: These rules only work for independent X and Y!

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!)

CS 109

topics

machine learning

general inference

sampling, making conclusions from data

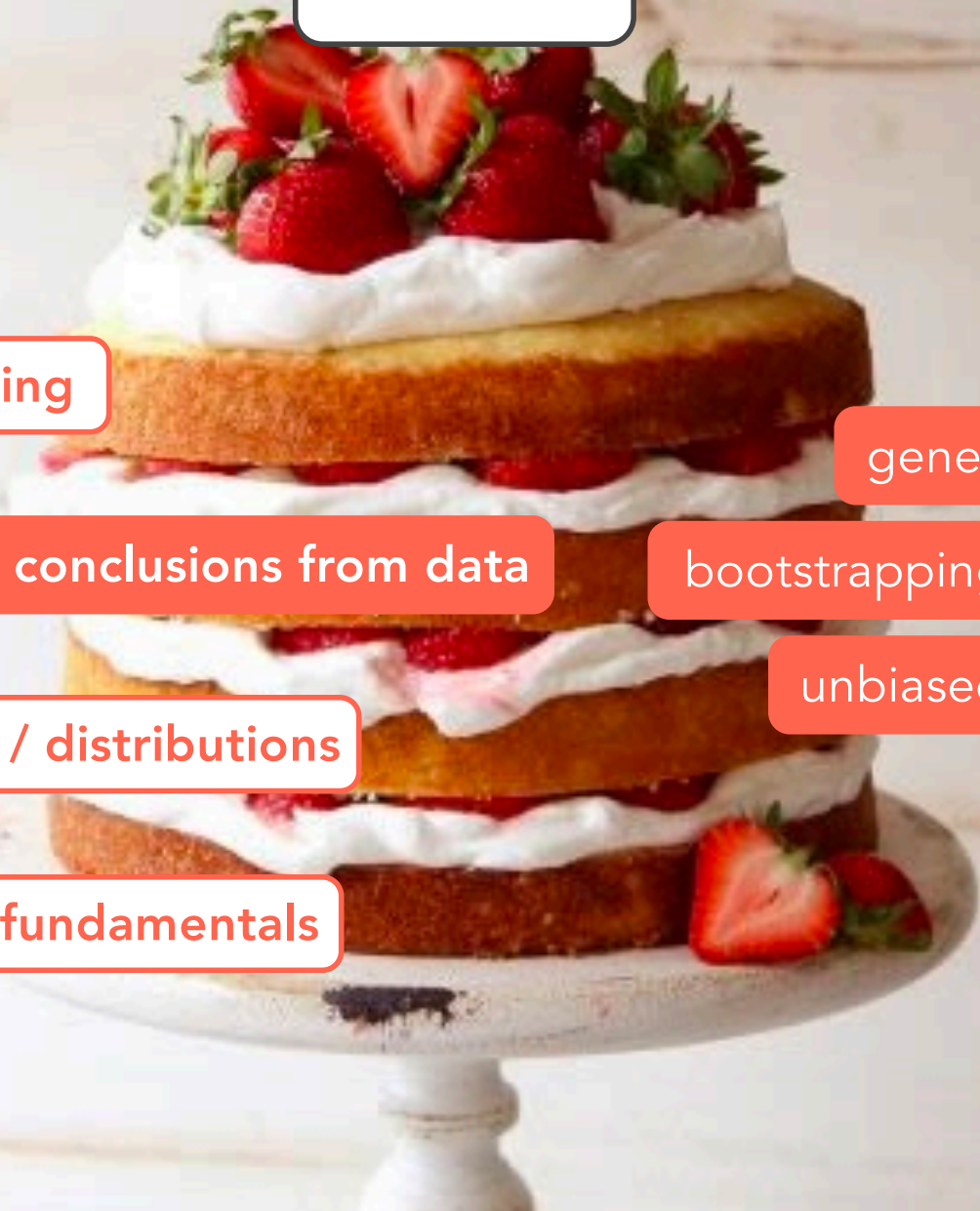
bootstrapping

CLT

unbiased estimators

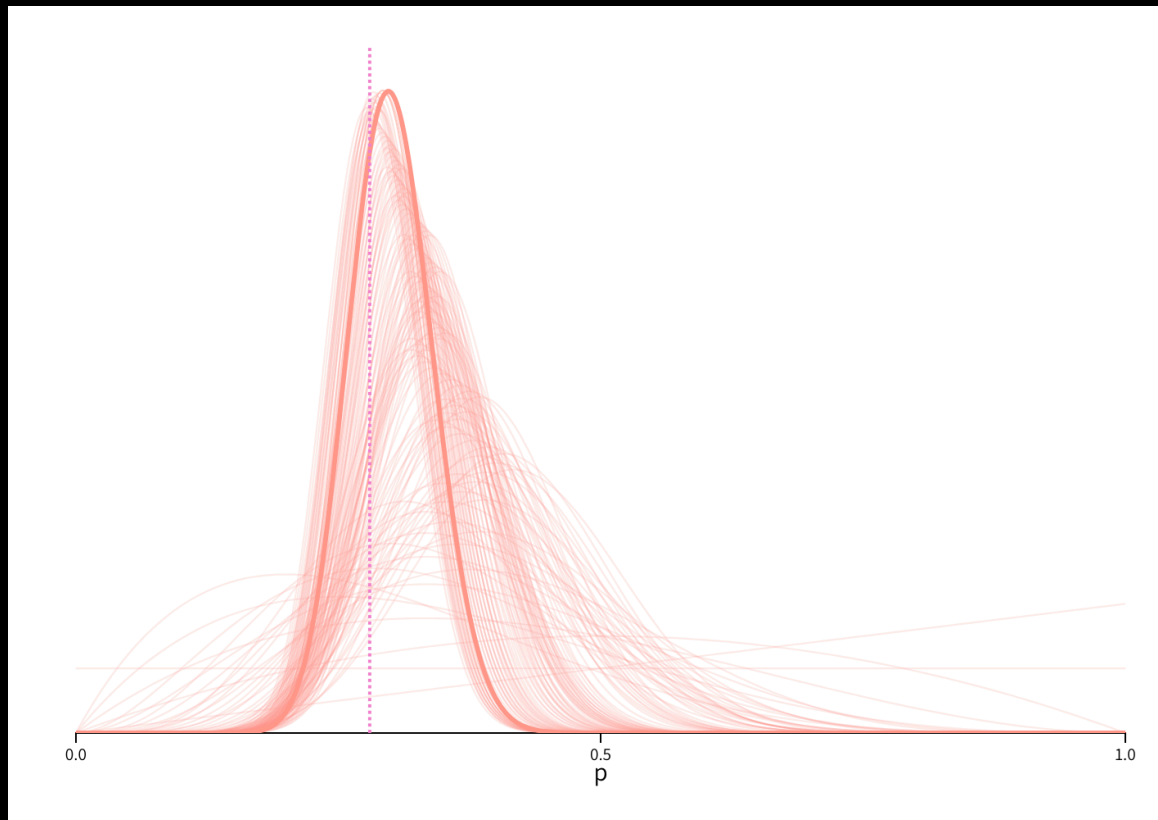
random variables / distributions

core probability fundamentals



Beta

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

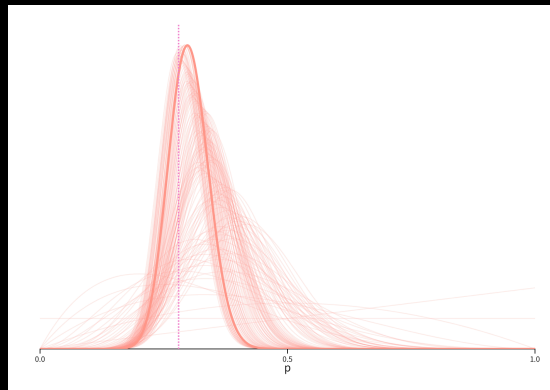


<https://seeing-theory.brown.edu/bayesian-inference/index.html#section3>

Beta

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

Bern(p) — 100101010100010101001010100



Updating belief about Bernoulli parameter p

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....



uh no



Sampling – Conceptual principles

Take a representative sample as large as you can

Sampling – Conceptual principles

Take a representative sample as large as you can

Sample statistics can be helpful in understanding the population

Sampling – Conceptual principles

Take a representative sample as large as you can

Sample statistics can be helpful in understanding the population

Be careful in assuming things about population from sample statistics
(you can bootstrap to better understand your population and statistics)

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!

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

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

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

makes the estimate unbiased

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

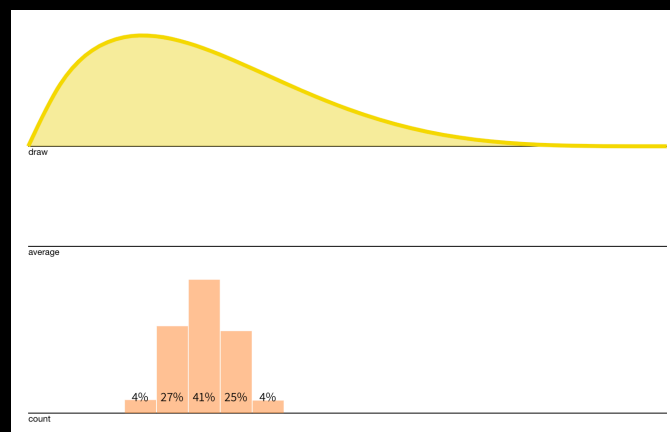
Taking Many Samples

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!



Central Limit Theorem

$$\sum_{i=1}^n X_i \sim \mathcal{N}(n\mu, n\sigma^2) \quad \text{As } n \rightarrow \infty$$

Extras



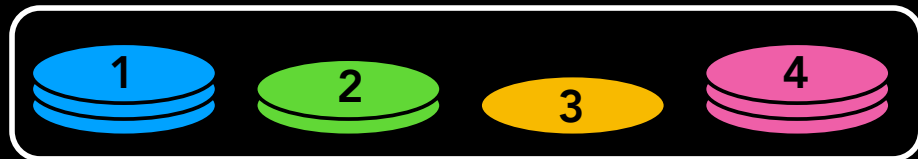
Bootstrapping: Simulating Many Samples From One

challenge

we want to better understand the distribution of our sample statistic(s),
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



1	2	3	4		in range?

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 use these estimates to answer new questions. For example: what was the probability of getting a sample variance between 1.5 and 2?

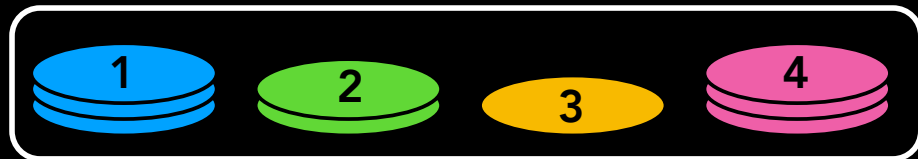
Bootstrapping: Simulating Many Samples From One

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



1	2	3	4		in range?
4	1	2	2		
3	3	3	0		
2	1	1	5		
...		

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 use these estimates to answer new questions. For example: what was the probability of getting a sample variance between 1.5 and 2?

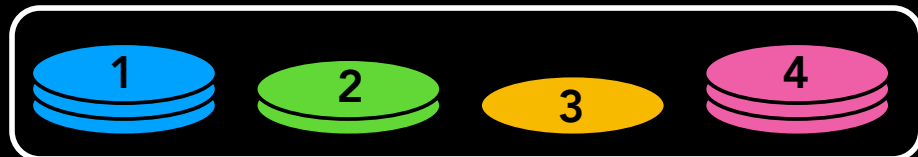
Bootstrapping: Simulating Many Samples From One

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



1	2	3	4	s^2	in range?
4	1	2	2	1.7	
3	3	3	0	0.8	
2	1	1	5	1.8	
...	

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 use these estimates to answer new questions. For example: what was the probability of getting a sample variance between 1.5 and 2?

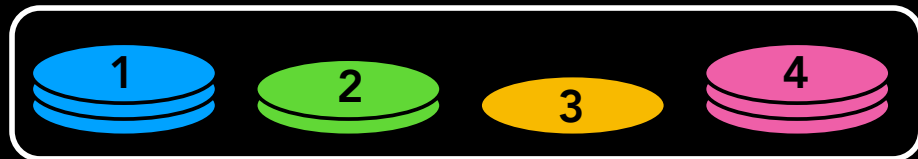
Bootstrapping: Simulating Many Samples From One

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

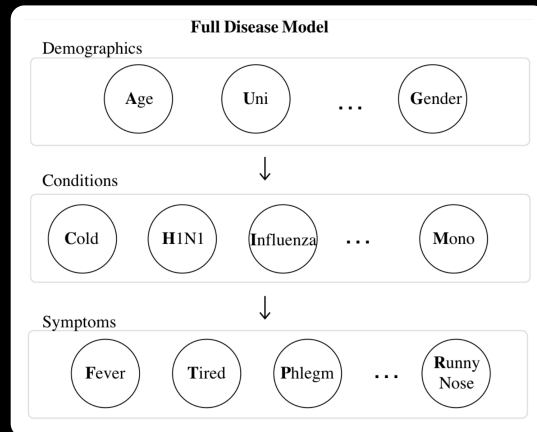
since our sample represents our population, we can sample from the data we have and it's as if we had gone out and collected more



1	2	3	4	s^2	in range?
4	1	2	2	1.7	yes
3	3	3	0	0.8	no
2	1	1	5	1.8	yes
...

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 use these estimates to answer new questions. For example: what was the probability of getting a sample variance between 1.5 and 2?

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



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

Think: What is reasonable to ask on a test about these topics?

Remember the ideas of the algorithms, and practice turning them into high-level pseudocode:

Computing sample statistics
Boot-strapping p values

Joint Sampling
Rejection Sampling
Thompson Sampling

CS 109



topics

machine learning

parameter estimation

classifiers

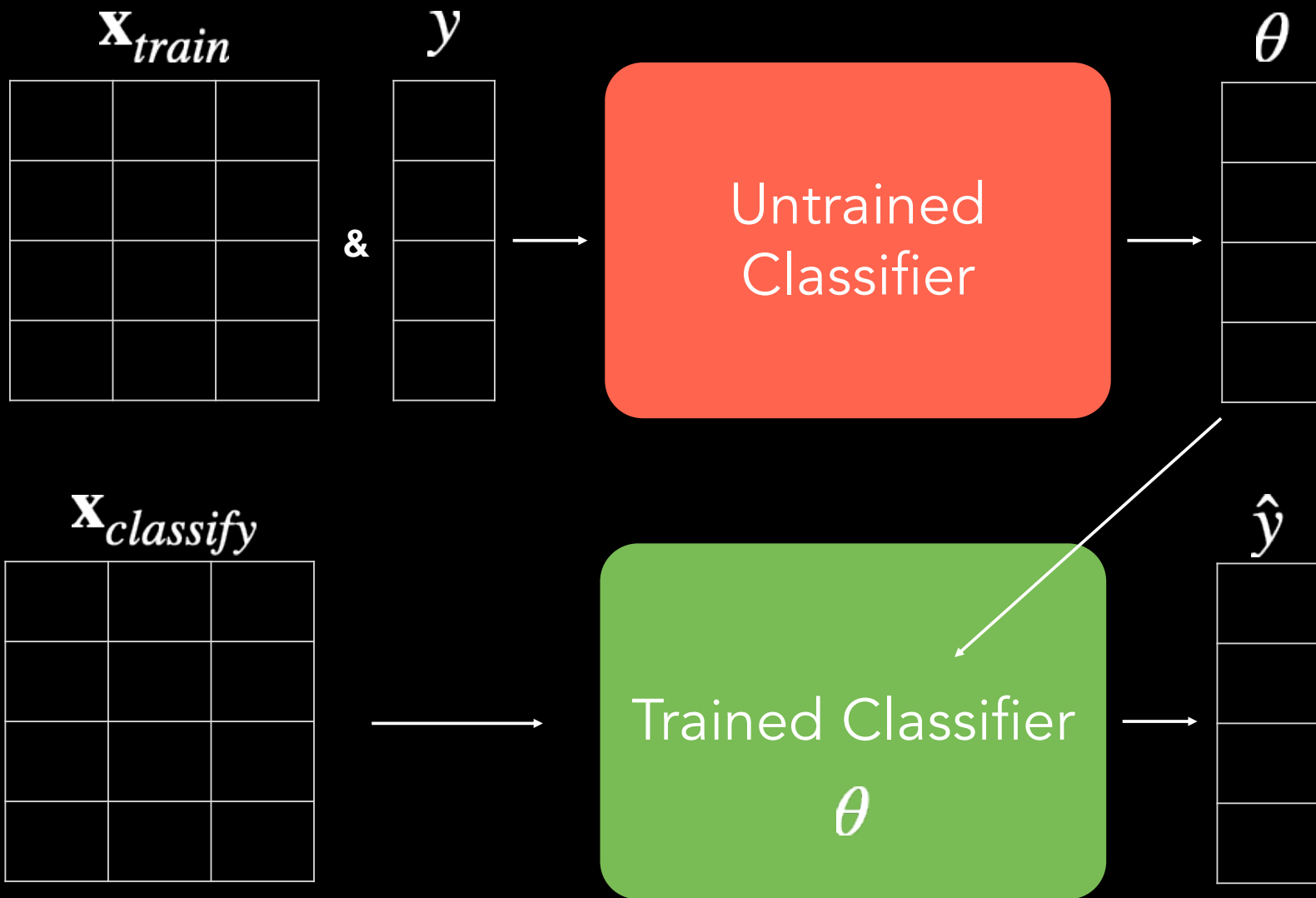
deep learning

sampling, making conclusions from data

random variables / distributions

core probability fundamentals

Classifiers



Parameter Estimation

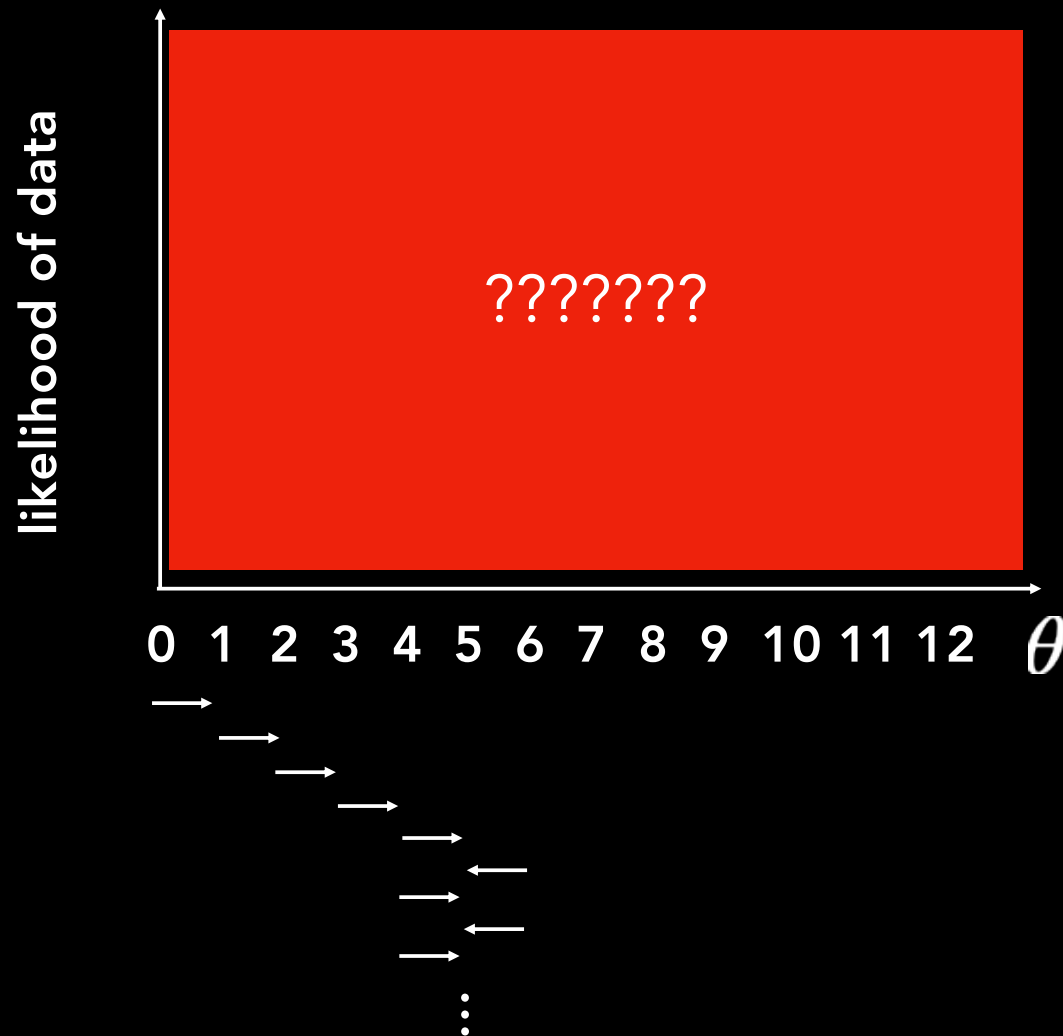
Maximum Likelihood Estimation

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

Gradient Ascent



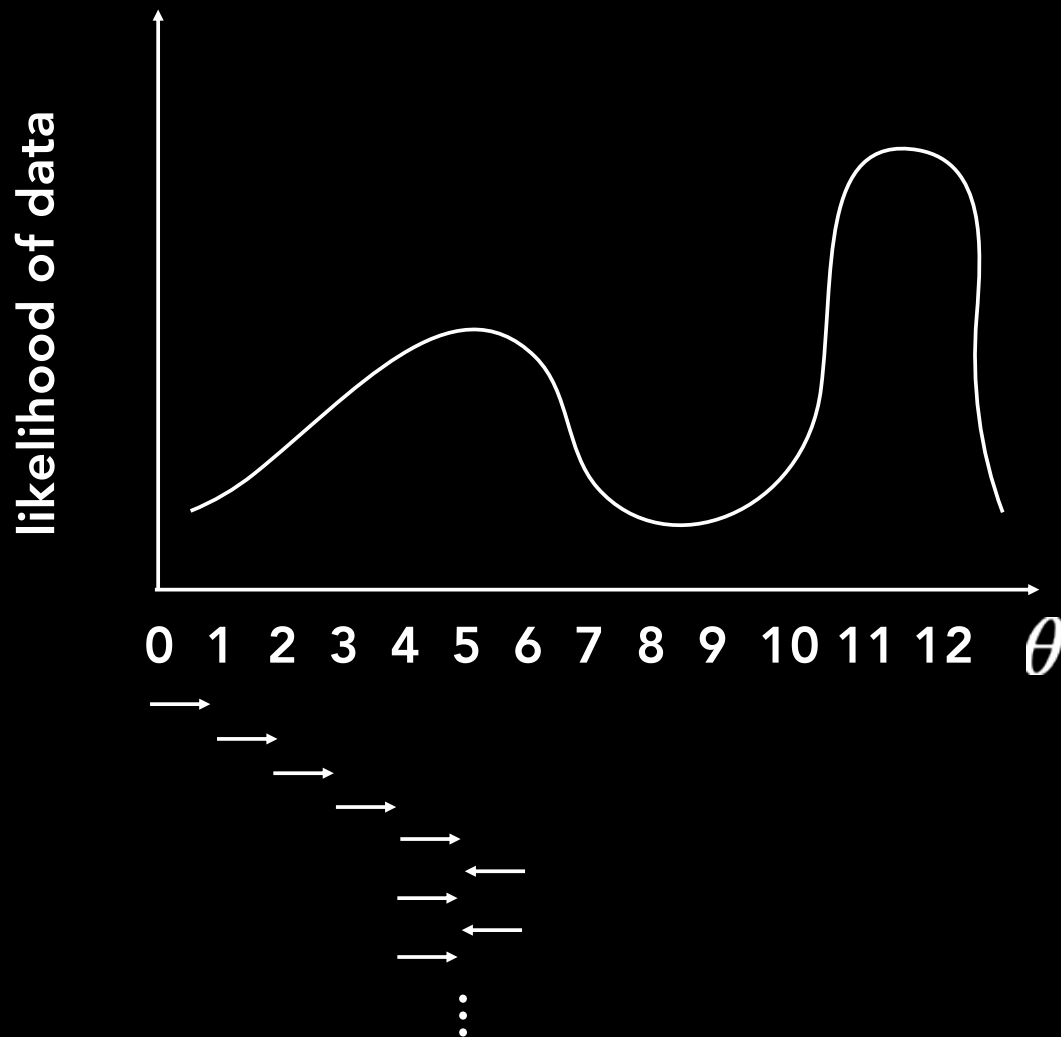
step size

$$\eta = 1$$

step direction

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

Gradient Ascent



step size

$$\eta = 1$$

step direction

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

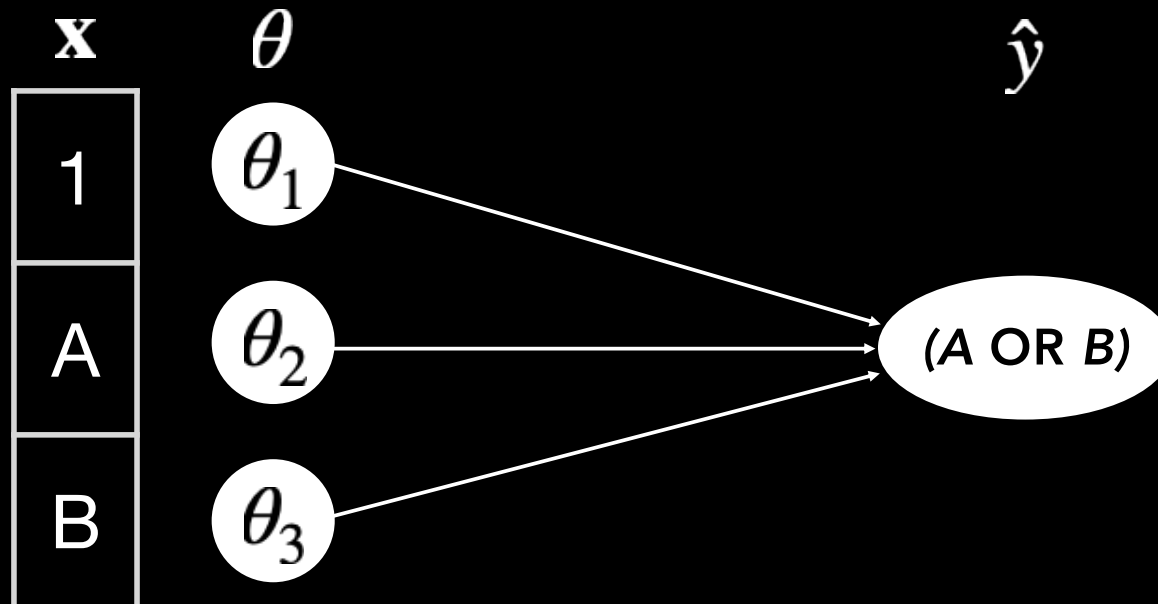
Classifier Algorithms

<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
At train: Best estimates for prior on y and conditional likelihood of data	What are we optimizing/figuring out? At test: Whether $y=0$ or $y=1$ is the best guess	At train: 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 $LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log [1 - \sigma(\theta^T \mathbf{x}^{(i)})]$ $\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)}$

Logistic Regression

Example

one neuron (logistic regression model)

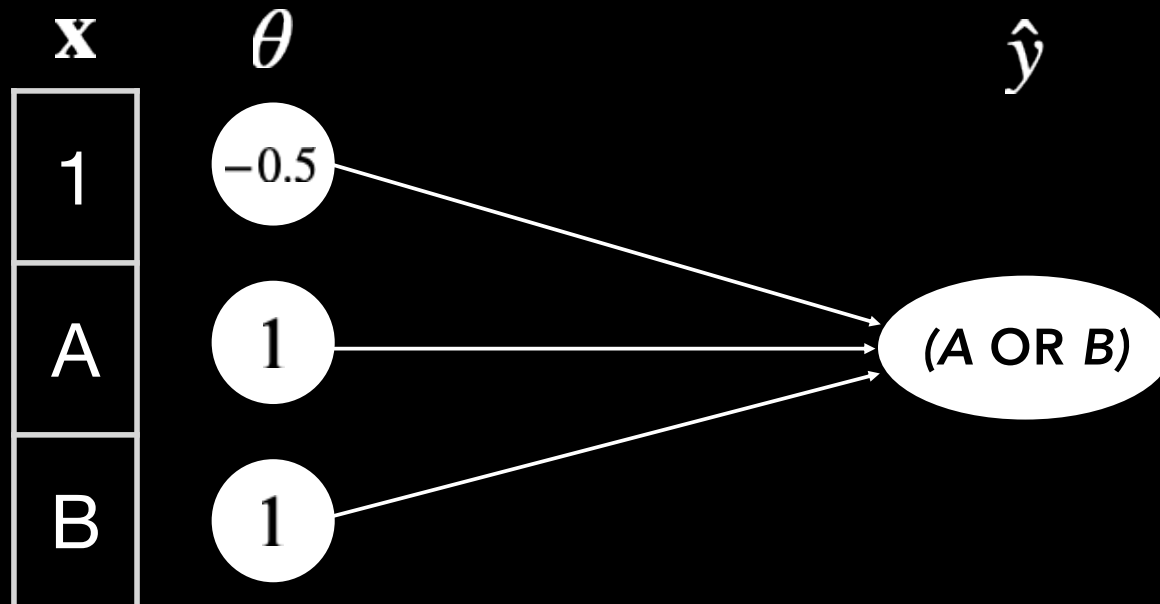


What weights do we have to learn for θ_1 , θ_2 , θ_3 to perfectly classify data of the form (A OR B)?

Logistic Regression

Example

one neuron (logistic regression model)

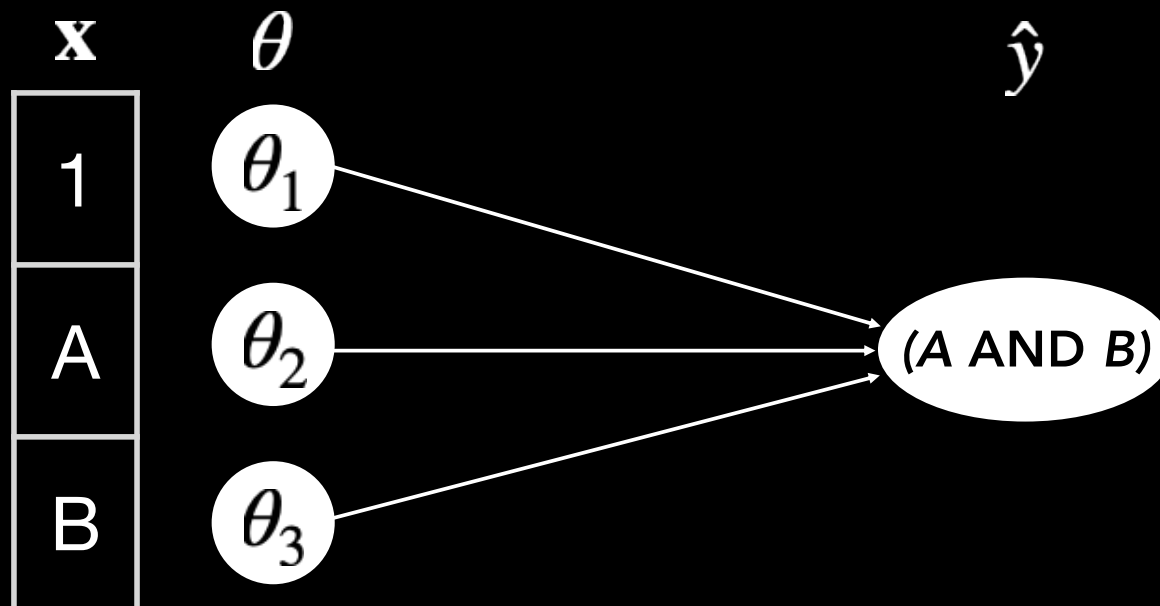


What weights do we have to learn for θ_1 , θ_2 , θ_3 to perfectly classify data of the form (A OR B)?

Logistic Regression

Example

one neuron (logistic regression model)

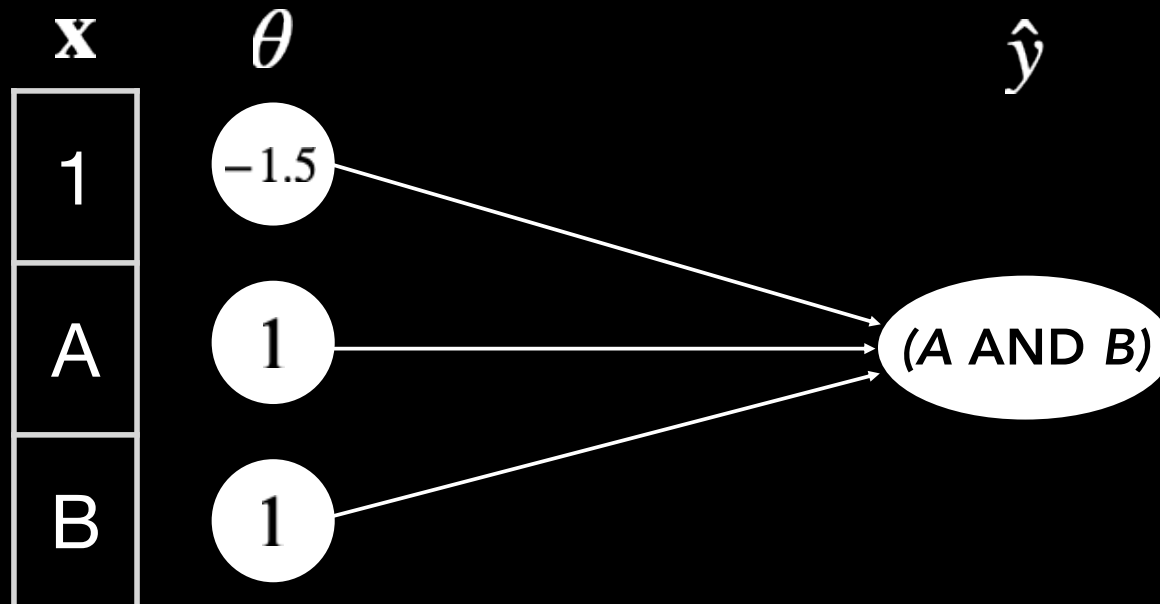


What weights do we have to learn for θ_1 , θ_2 , θ_3 to perfectly classify data of the form (A AND B)?

Logistic Regression

Example

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)?