



Parameter Estimation

Chris Piech
CS109, Stanford University



	titanic.csv — website							
	overview.html		problem12.html		titanic.csv		index.html	
1	Survived	Pclass	Name	Sex	Age	Siblings/Spouses Aboard	Parents/Children Aboard	Fare
2	0	3	"Braund, Mr. Owen Harris"	male	22	1	0	7.25
3	1	1	"Cumings, Mrs. John Bradley (Florence Briggs Thayer)"	female	38	1	0	71.2833
4	1	3	"Heikkinen, Miss. Laina"	female	26	0	0	7.925
5	1	1	"Futrelle, Mrs. Jacques Heath (Lily May Peel)"	female	35	1	0	53.1
6	0	3	"Allen, Mr. William Henry"	male	35	0	0	8.05
7	0	3	"Moran, Mr. James"	male	27	0	0	8.4583
8	0	1	"McCarthy, Mr. Timothy J"	male	54	0	0	51.8625
9	0	3	"Palsson, Master. Gosta Leonard"	male	2	3	1	21.075
10	1	3	"Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)"	female	27	0	2	11.1333
11	1	2	"Nasser, Mrs. Nicholas (Adele Achem)"	female	14	1	0	30.0708
12	1	3	"Sandstrom, Miss. Marguerite Rut"	female	4	1	1	16.7
13	1	1	"Bonnell, Miss. Elizabeth"	female	58	0	0	26.55
14	0	3	"Saunderscock, Mr. William Henry"	male	20	0	0	8.05
15	0	3	"Andersson, Mr. Anders Johan"	male	39	1	5	31.275
16	0	3	"Vestrom, Miss. Hulda Amanda Adolfina"	female	14	0	0	7.8542
17	1	2	"Hewlett, Mrs. (Mary D Kingcome)"	female	55	0	0	16
18	0	3	"Rice, Master. Eugene"	male	2	4	1	29.125
19	1	2	"Williams, Mr. Charles Eugene"	male	23	0	0	13
20	0	3	"Vander Planke, Mrs. Julius (Emilia Maria Vandemoortele)"	female	31	1	0	18
21	1	3	"Masselmanni, Mrs. Fatima"	female	22	0	0	7.225
22	0	2	"Fynney, Mr. Joseph J"	male	35	0	0	26
23	1	2	"Beesley, Mr. Lawrence"	male	34	0	0	13
24	1	3	"McGowan, Miss. Anna ""Annie"""	female	15	0	0	8.0292
25	1	1	"Sloper, Mr. William Thompson"	male	28	0	0	35.5
26	0	3	"Palsson, Miss. Torborg Danira"	female	8	3	1	21.075
27	1	3	"Asplund, Mrs. Carl Oscar (Selma Augusta Emilia Johansson)"	female	38	1	5	31.3875
28	0	3	"Emir, Mr. Farred Chehab"	male	26	0	0	7.225
29	0	1	"Fortune, Mr. Charles Alexander"	male	19	3	2	263
30	1	3	"O'Dwyer, Miss. Ellen ""Nellie"""	female	24	0	0	7.8792
31	0	3	"Todoroff, Mr. Lalio"	male	23	0	0	7.8958
32	0	1	"Uruchurtu, Don. Manuel E"	male	40	0	0	27.7208
33	1	1	"Spencer, Mrs. William Augustus (Marie Eugenie)"	female	48	1	0	146.5208
34	1	3	"Glynn, Miss. Mary Agatha"	female	18	0	0	7.75
35	0	2	"Wheaton, Mr. Edward H"	male	66	0	0	10.5
36								
37	Survived	Pclass	Name	Sex	Age			
38	0	3	Braund, Mr. Owen Harris	male				
39	1	1	Cumings, Mrs. John Bradley (Florence Thayer)	female				
40	1	3	Heikkinen, Miss. Laina	female				
41	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female				
line 27	0	3	Allen, Mr. William Henry	male				
	0	3	Moran, Mr. James	male				
	0	1	McCarthy, Mr. Timothy J	male				
	0	3	Palsson, Master. Gosta Leonard	male				

Microsoft Excel - titanic

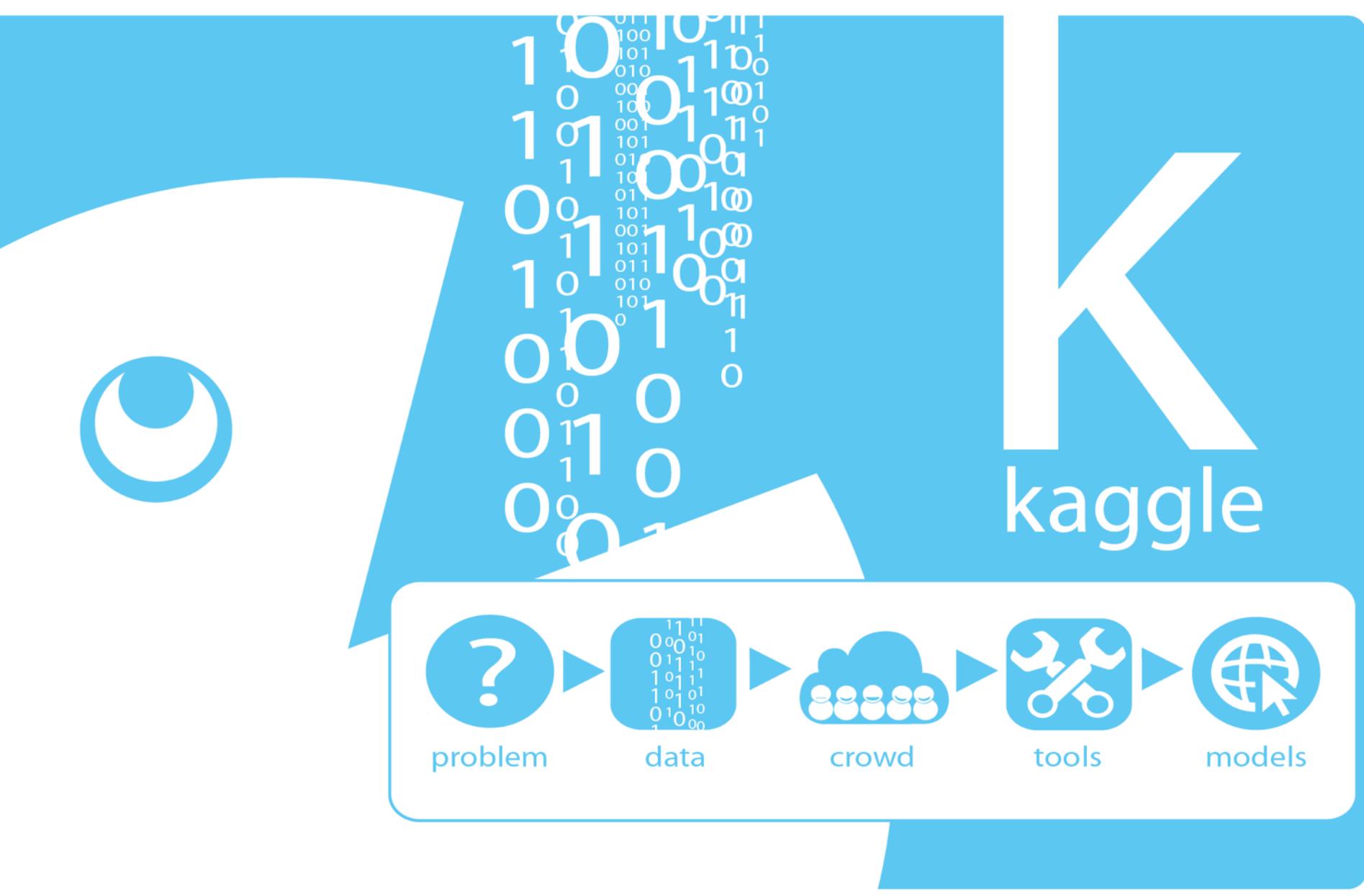
Search Sheet

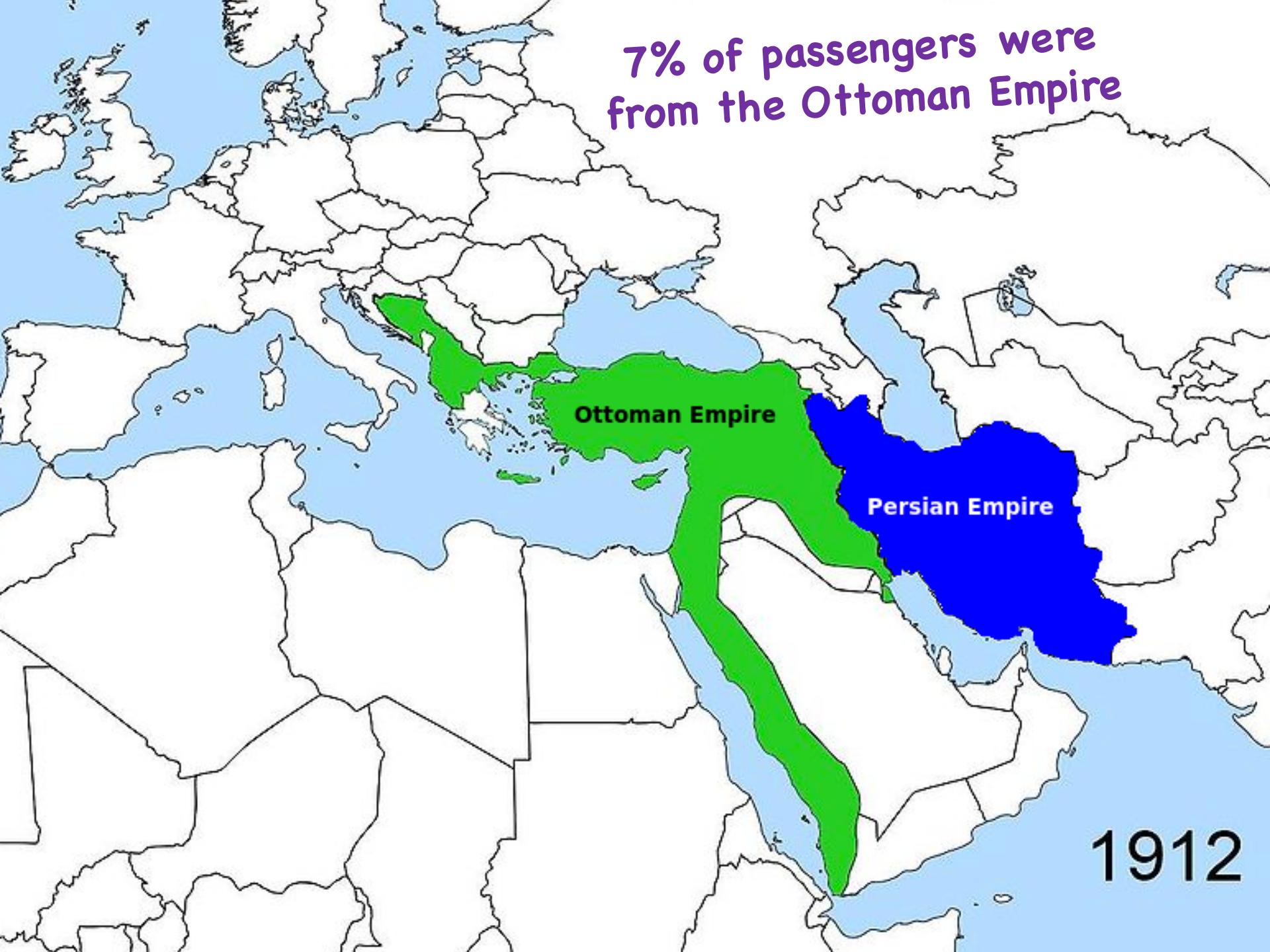
Home Insert Page Layout Formulas Data Review View

K3

Get External Data Refresh All Edit Links Sort Filter Advanced Text to Columns Group Ungroup Subtotal

	A	B	C	D	E	F	G	H	I
1	Survived	Pclass	Name	Sex	Age	Siblings/Spouses Aboard	Parents/Children Aboard	Fare	
2	0	3	Braund, Mr. Owen Harris	male	22	1	0	7.25	
3	1	1	Cumings, Mrs. John Bradley (Florence ^{she})	female	38	1	0	71.2833	
4	1	3	Heikkinen, Miss. Laina	female	26	0	0	7.925	
5	1	1	Futrelle, Mrs. Jacques Heath (Lily May ^{she})	female	35	1	0	53.1	
6	0	3	Allan, Mr. William Henry	male	35	0	0	8.05	
7	0	3	Moran, Mr. James	male	27	0	0	8.4583	
8	0	1	McCarthy, Mr. Timothy J	male	54	0	0	51.8625	
9	0	3	Palsson, Master, Gosta Leonard	male	2	3	1	21.075	
10	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina ^{she})	female	27	0	2	11.1333	
11	1	2	Nasser, Mrs. Nicholas (Adele Acherman ^{she})	female	14	1	0	30.0708	
12	1	3	Sandstrom, Miss. Marguerite Rut	female	4	1	1	16.7	
13	1	1	Bonnell, Miss. Elizabeth	female	58	0	0	26.55	
14	0	3	Saundercock, Mr. William Henry	male	20	0	0	8.05	
15	0	3	Andersson, Mr. Anders Johansen	male	39	1	5	31.275	
16	0	3	Vestrom, Miss. Hulda Amanda Adolfi	female	14	0	0	7.8542	
17	1	2	Hewlett, Mrs. (Mary D'Gama ^{she})	female	55	0	0	16	
18	0	3	Rice, Master, Eugene	male	2	4	1	29.125	
19	1	2	Williams, Mr. Charles Eugene	male	23	0	0	13	
20	0	3	Vander Plank, Mrs. Julius (Emelia M ^{she})	female	31	1	0	18	
21	1	3	Masselmani, Mrs. Fatima	female	22	0	0	7.225	
22	0	2	Fynney, Mr. Joseph J	male	35	0	0	26	
23	1	2	Beesley, Mr. Lawrence	male	34	0	0	13	
24	1	3	McGowan, Miss. Anna "Annie"	female	15	0	0	8.0292	
25	1	1	Sloper, Mr. William Thompson	male	28	0	0	35.5	
26	0	3	Palsson, Miss. Torborg Danira	female	8	3	1	21.075	
27	1	3	Asplund, Mrs. Carl Oscar (Selma Augusta ^{she})	female	38	1	5	31.3875	
28	0	3	Emir, Mr. Farred Chehab	male	26	0	0	7.225	
29	0	1	Fortune, Mr. Charles Alexander	male	19	3	2	263	
30	1	3	O'Dwyer, Miss. Ellen "Nellie"	female	24	0	0	7.8792	
31	0	3	Todoroff, Mr. Lalio	male	23	0	0	7.8958	
32	0	1	Uruchurtu, Don. Manuel E	male	40	0	0	27.7208	
33	1	1	Spencer, Mrs. William Augustus (Marie ^{she})	female	48	1	0	146.5208	
34	1	3	Glynn, Miss. Mary Agatha	female	18	0	0	7.75	
35	0	2	Wheadon, Mr. Edward H	male	66	0	0	10.5	
36	0	1	Meyer, Mr. Edgar Joseph	male	28	1	0	82.1708	
37	0	1	Holmstrom, Mr. Alexander Oskar	male	42	1	0	52	
									292
									1.05
									18
									417
	Siblings/Spouses Aboard			Parents/Children Aboard			Fare		
22				1			0	7.25	
38				1			0	71.2833	
26				0			0	7.925	
35				1			0	53.1	
35				0			0	8.05	
27				0			0	8.4583	
54				0			0	51.8625	
2				3			1	21.075	





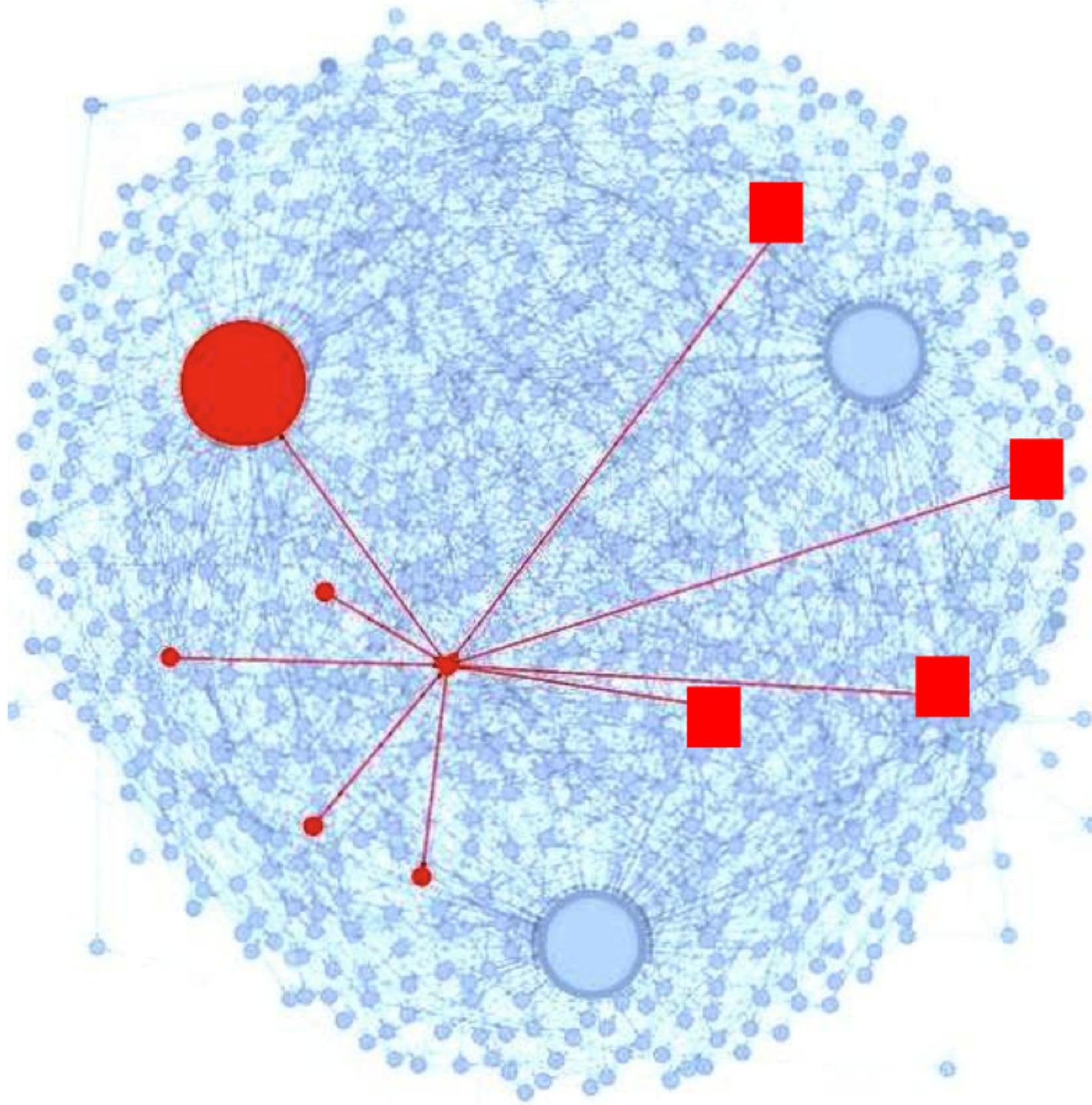
7% of passengers were
from the Ottoman Empire

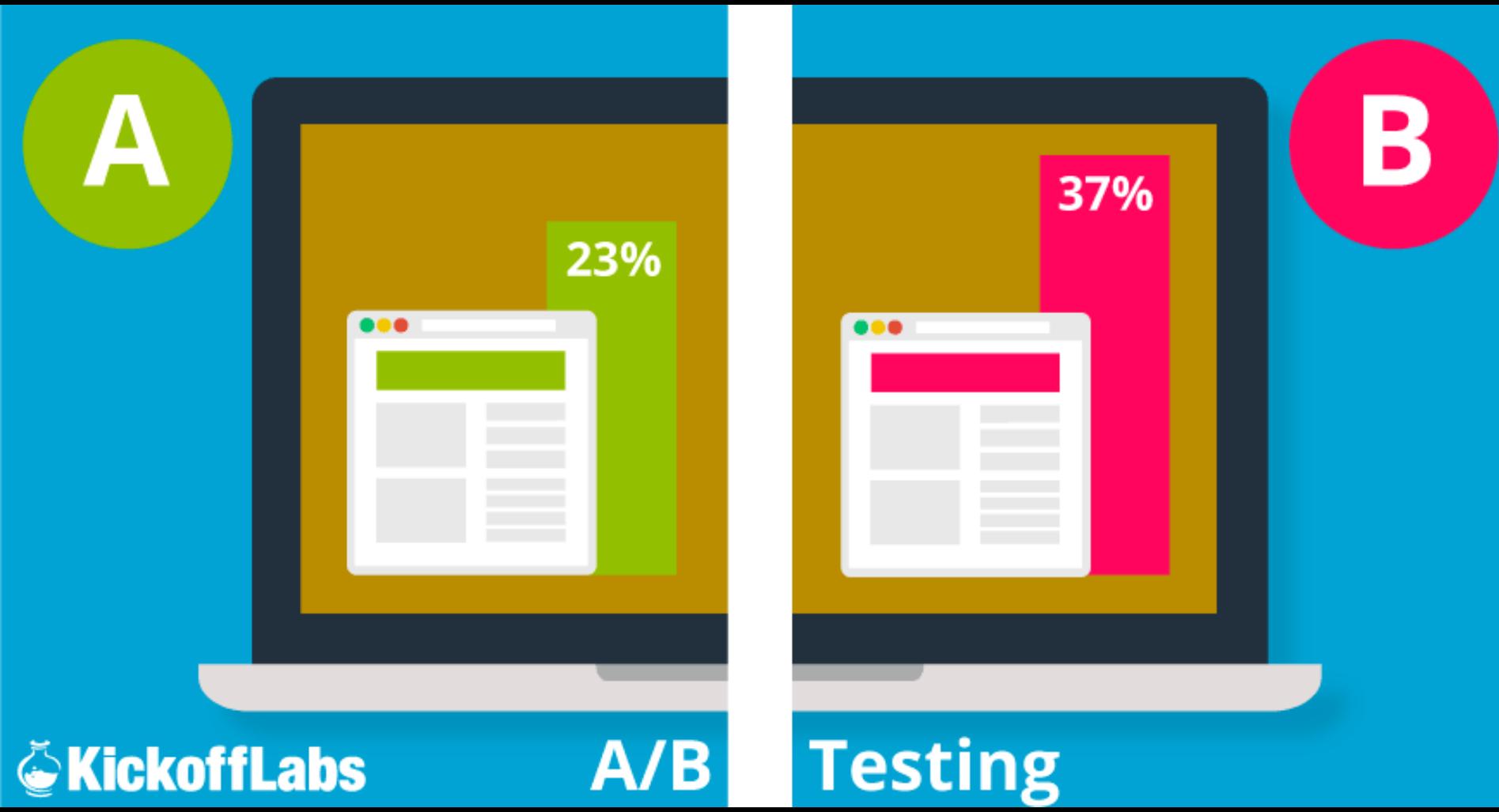
Ottoman Empire

Persian Empire

1912







Conditional Expectation

Conditional Expectation

- X and Y are jointly discrete random variables
 - Recall conditional PMF of X given $Y = y$:

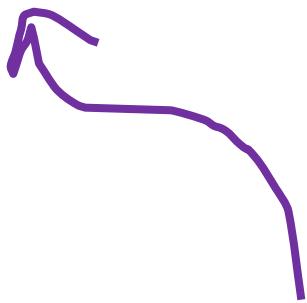
$$p_{X|Y}(x | y) = P(X = x | Y = y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}$$

- 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:

$$f_{X|Y}(x | y) = \frac{f_{X,Y}(x, y)}{f_Y(y)}$$

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

$$X | Y = 5$$



This is just a random variable

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 $E[X | Y = 6]$?

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

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

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:

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

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 $E[Y]$?

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

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

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

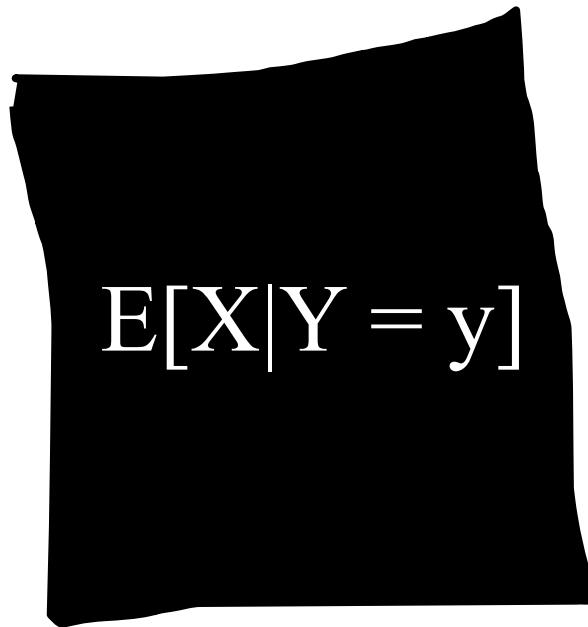
$$E[Y | X = 3] = E[7 + Y] = 7 + E[Y]$$

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

$$E[Y] = 15$$

Expectations of Conditional Expectation

- Define $g(Y) = E[X | Y]$
 - For any $Y = y$, $g(Y) = E[X | Y = y]$
 - This is just function of Y , since we sum over all values of X

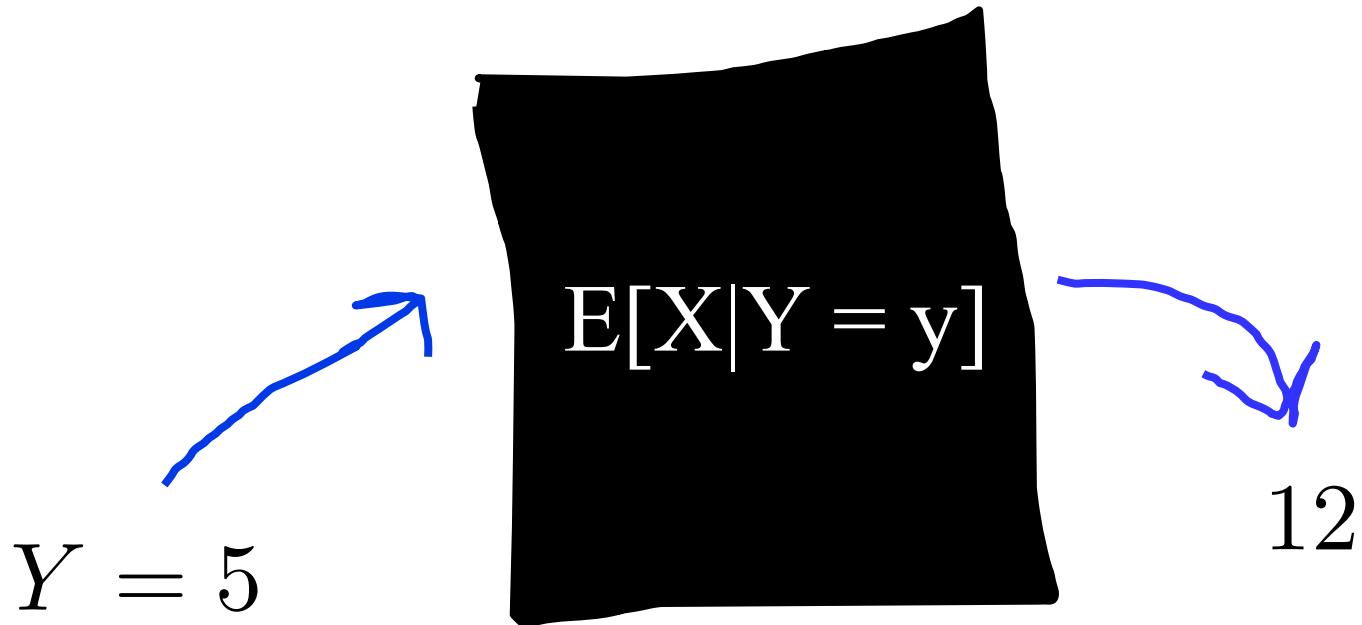


$E[X|Y = y]$

This is a function

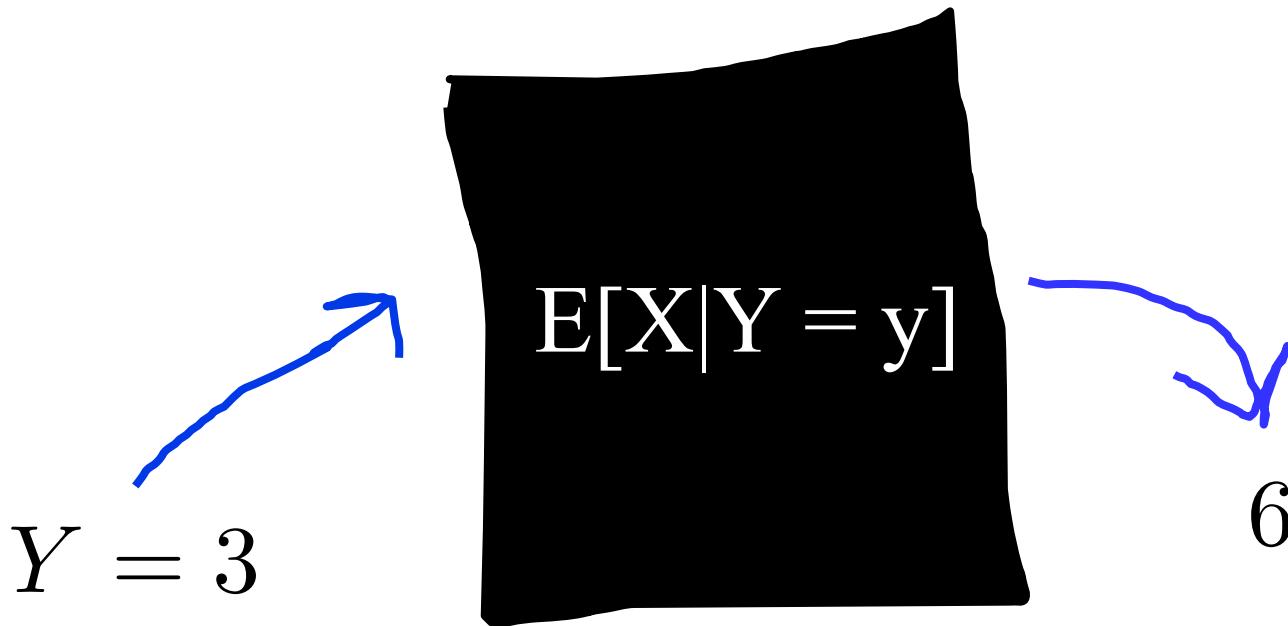
Expectations of Conditional Expectation

- Define $g(Y) = E[X | Y]$
 - For any $Y = y$, $g(Y) = E[X | Y = y]$
 - This is just function of Y , since we sum over all values of X



Expectations of Conditional Expectation

- Define $g(Y) = E[X | Y]$
 - For any $Y = y$, $g(Y) = E[X | Y = y]$
 - This is just function of Y , since we sum over all values of X



Central Limit Theorem

The Central Limit Theorem

- Consider I.I.D. random variables X_1, X_2, \dots
 - X_i have distribution with $E[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2$
 - Let: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$
 - Central Limit Theorem:

Version 1

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \text{ as } n \rightarrow \infty$$

Version 2

$$\frac{X_1 + X_2 + \dots + X_n - n\mu}{\sigma\sqrt{n}} \rightarrow N(0, 1) \text{ as } n \rightarrow \infty$$

Get Good at Manipulating Normals

$$X \sim N(\mu, \sigma^2)$$

$$aX + b \sim N(a\mu, a^2\sigma^2)$$

$$X \sim N(\mu_1, \sigma_1) \qquad Y \sim N(\mu_2, \sigma_2)$$

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

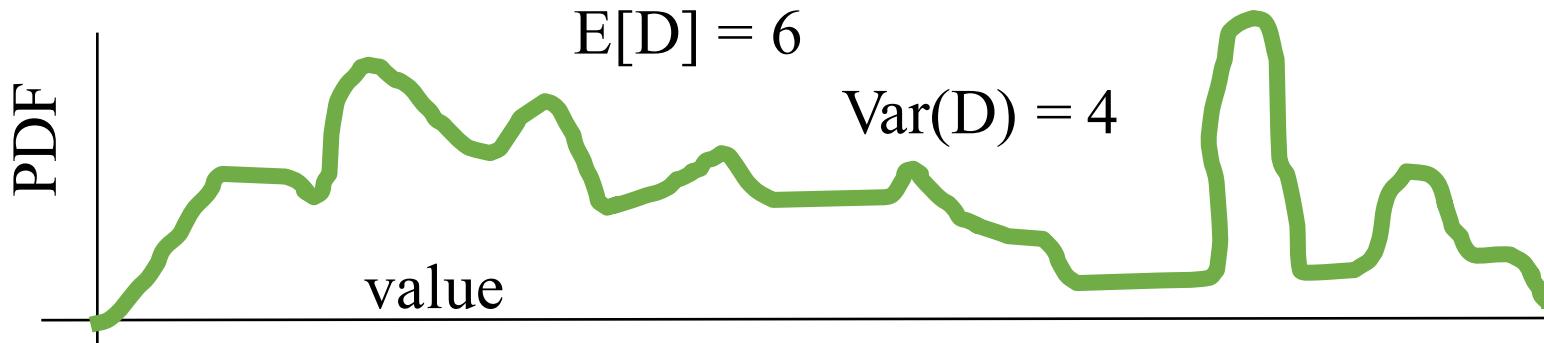
(If X and Y are independent)

How To Find Central Limit Theorem?

Are you averaging many (>10) I.I.D. random variables?

Are you adding many (>10) I.I.D. random variables?

You are tracking an object on a 1D line and know its location X . Your radar goes down and you don't get to observe it for 20 time steps. Each time step you assume that its change in position is IID with this pdf:



What is the distribution of your belief about the location of the object after 20 time steps?

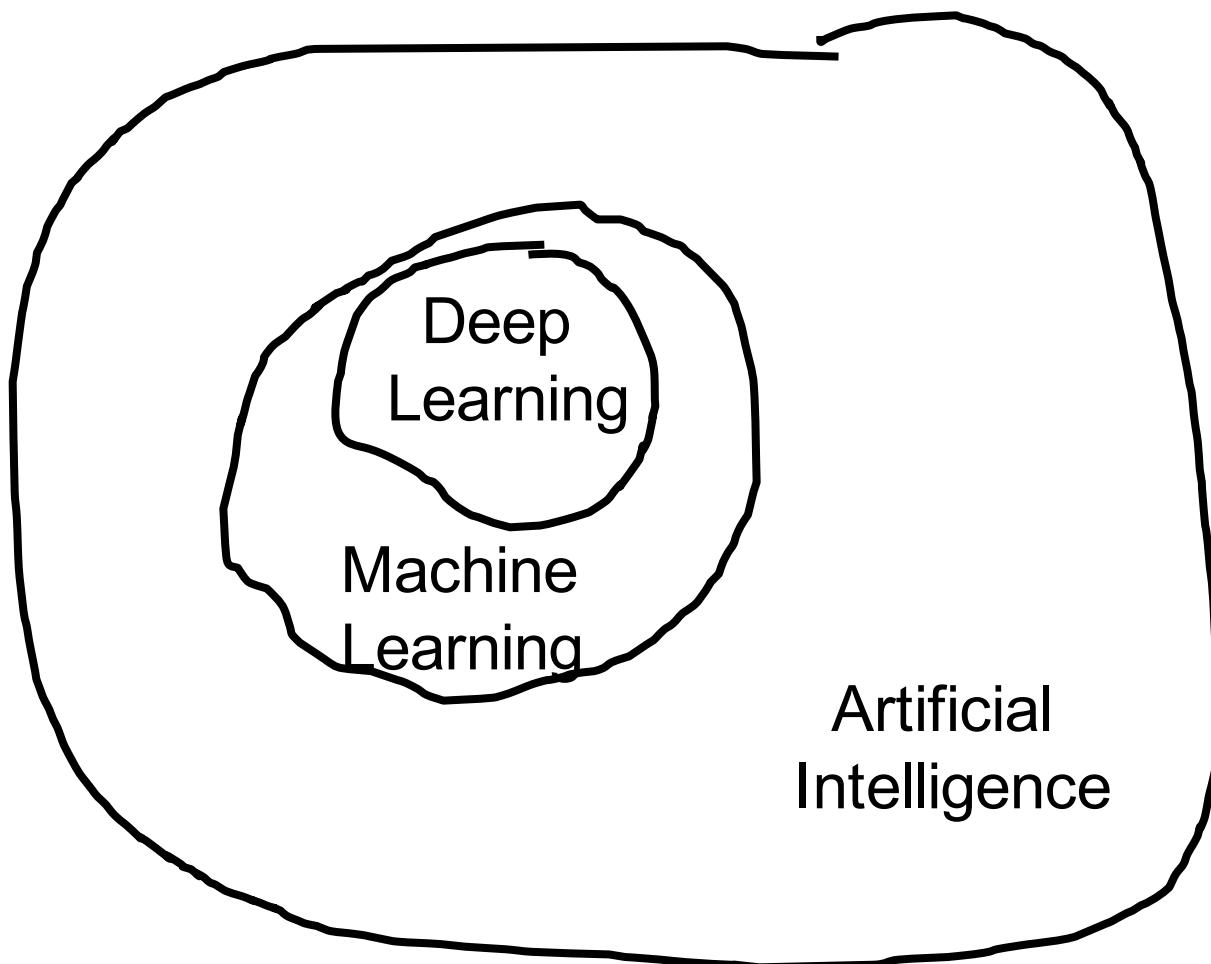
End Review

What is AI?

[suspense]

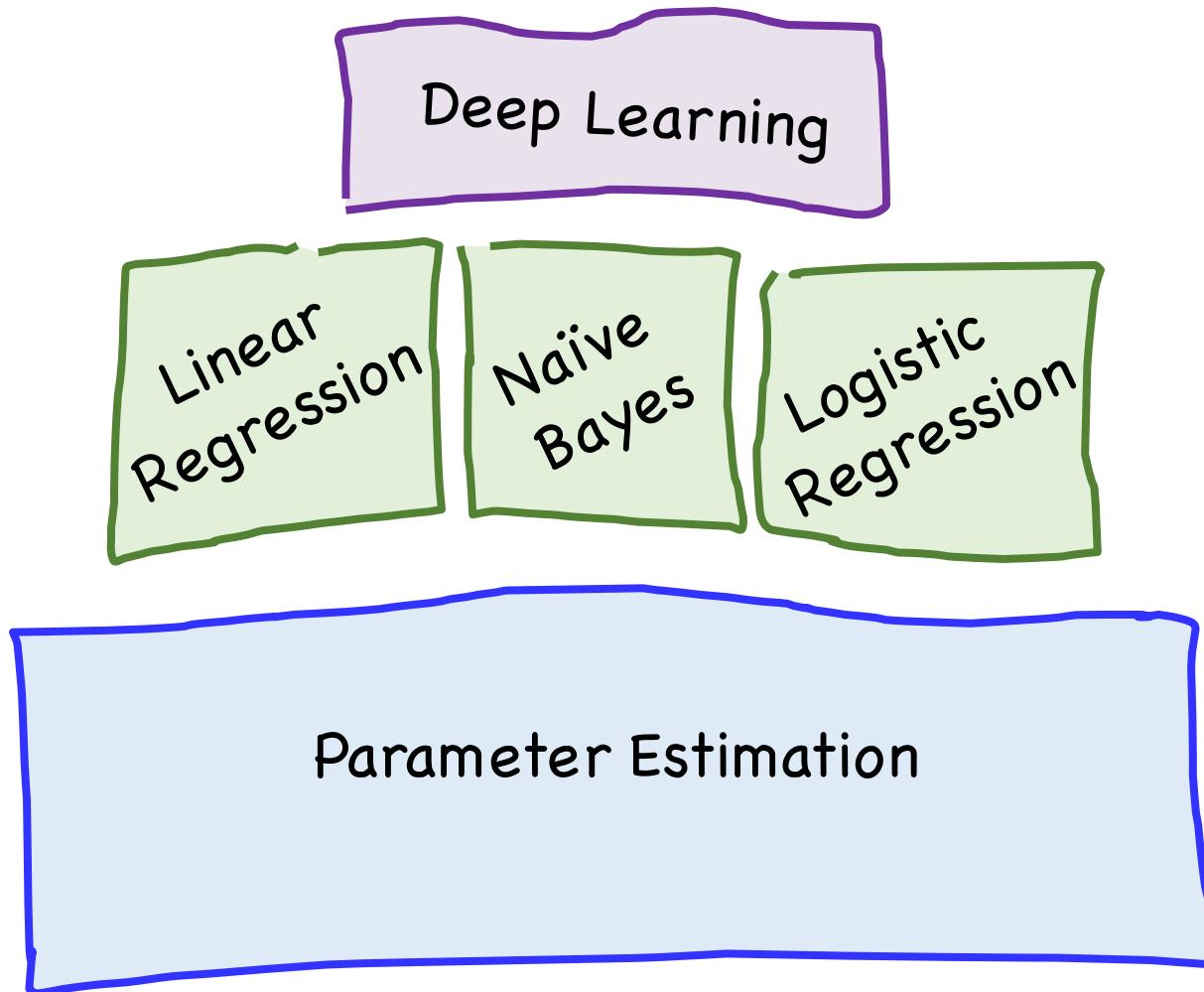
AI: The study and design
of intelligent **agents**

AI and Machine Learning

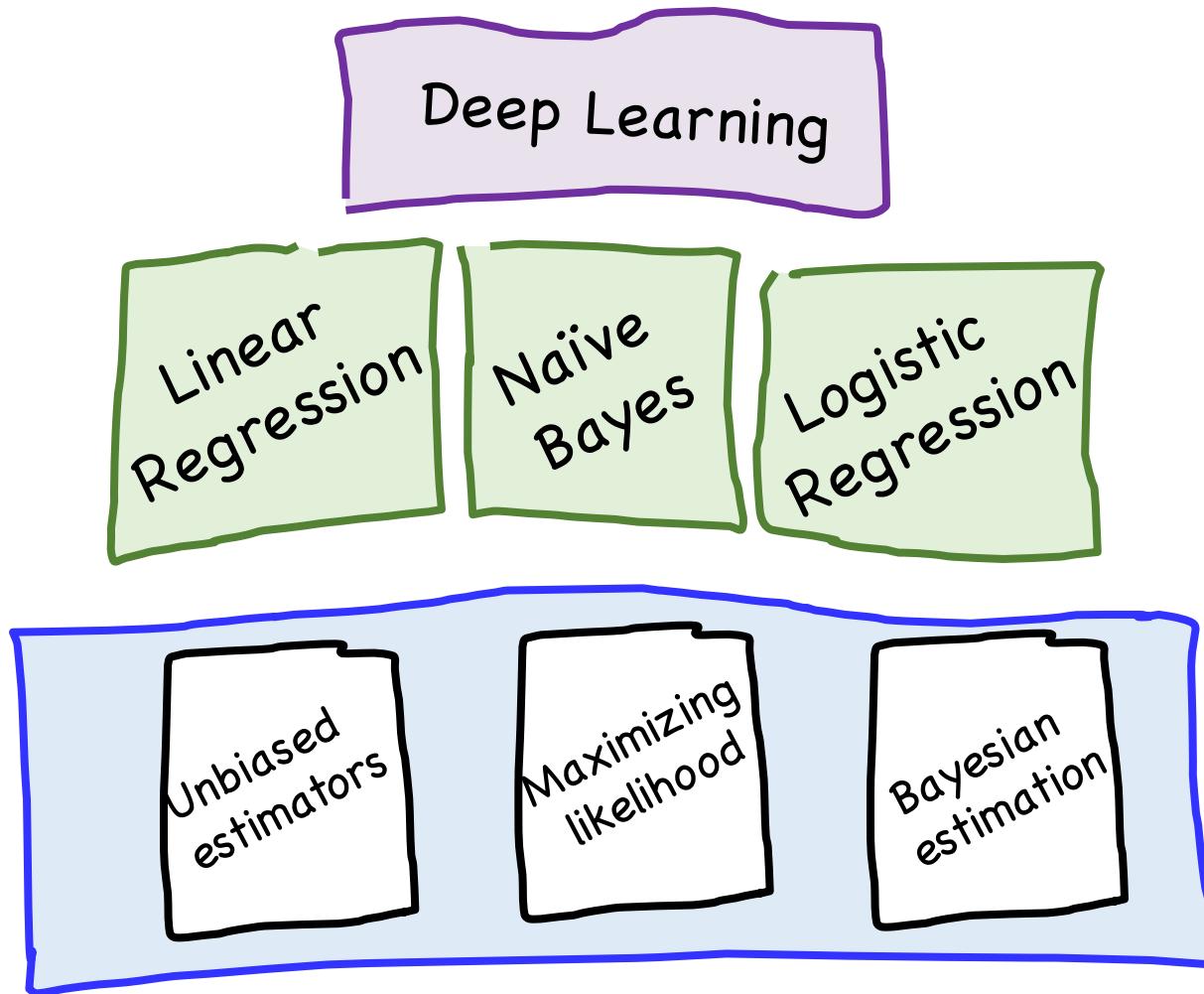


ML: Rooted in probability theory

Our Path



Our Path



Jump Straight to Deep Learning?

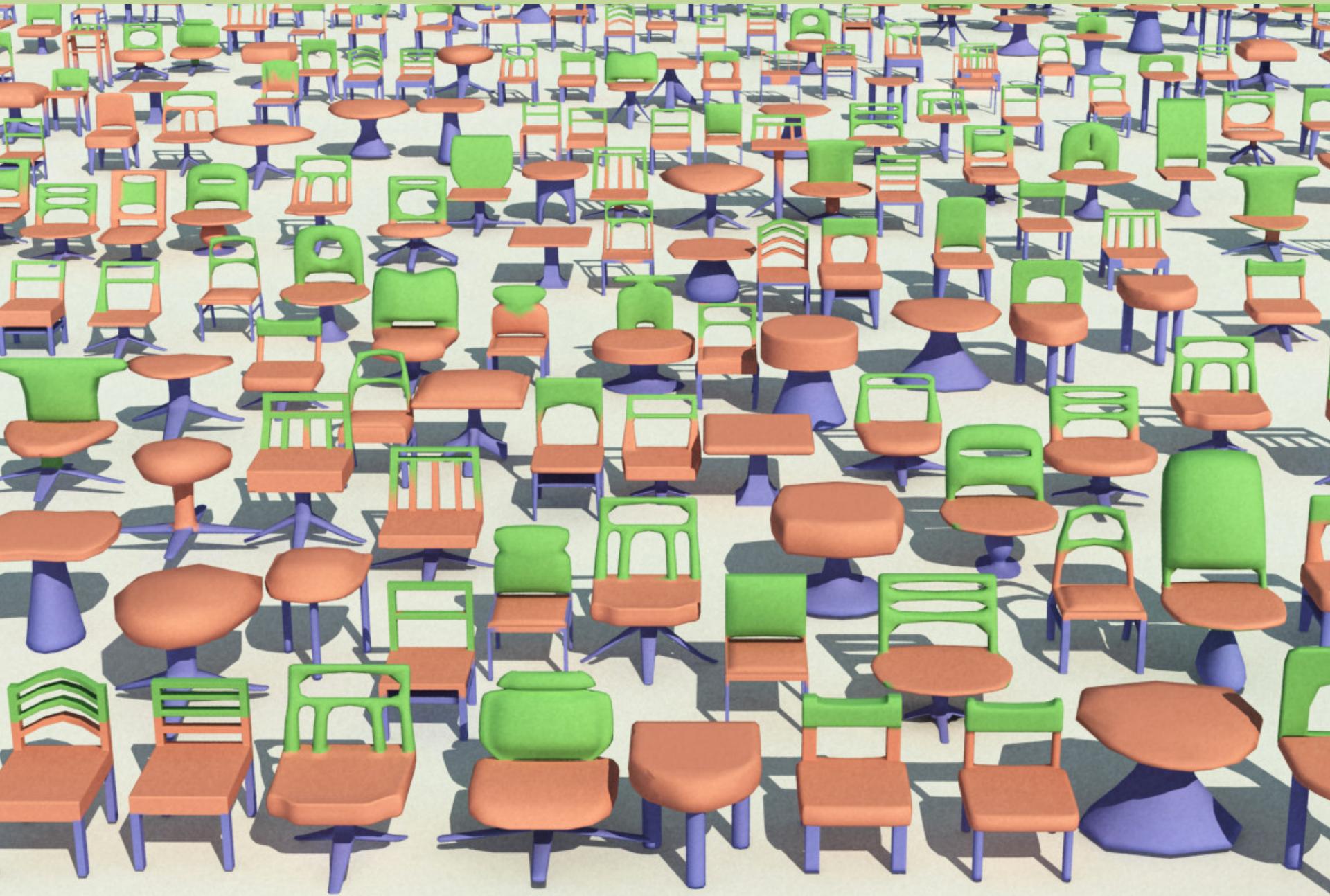
Tensor Flow



Understand the theory to help you debug

But another reason...

Machine Learning Uses a Lot of Data



One Shot Learning

Single training example:

କୁ

Test set:

a	ଶ	ଅ	ଶ
କୁ	ଅ	ପ୍ଲ	କୁ
ମ	କୁ	ହେ	କୁ
ମ	ଅ	କୁ	ନ୍ତର

One Shot Learning

Single
training
example:



Computers can't do that.

Understand the theory to push on the grand challenges

A silhouette of the iconic Disney castle is positioned in the center of the background, partially obscured by a dark, star-filled foreground.

WALT DISNEY
PICTURES

Once upon a time...

...there was parameter estimation

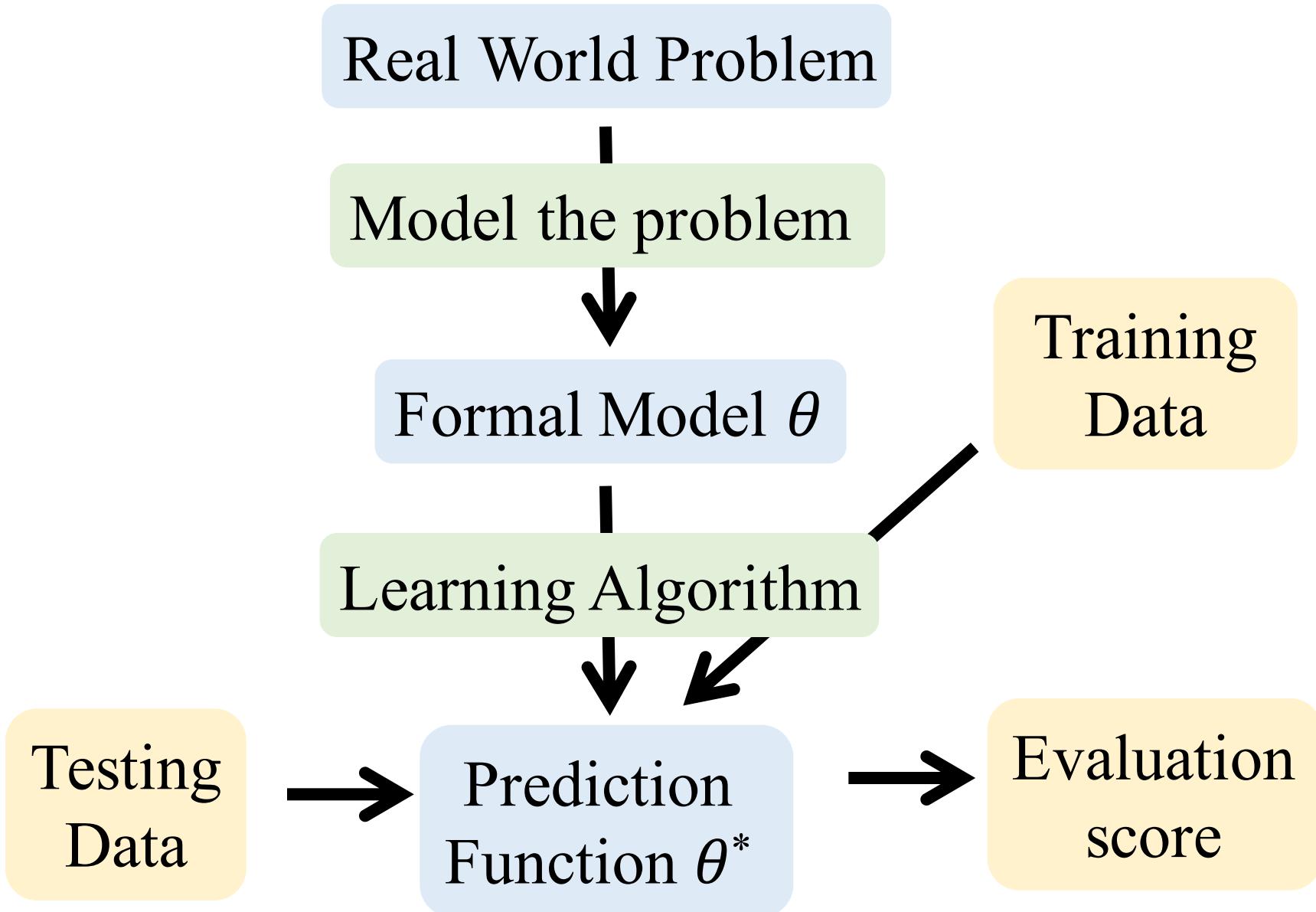
What are Parameters?

- Consider some probability distributions:
 - $\text{Ber}(p)$ $\theta = p$
 - $\text{Poi}(\lambda)$ $\theta = \lambda$
 - $\text{Uni}(\alpha, \beta)$ $\theta = (\alpha, \beta)$
 - $\text{Normal}(\mu, \sigma^2)$ $\theta = (\mu, \sigma^2)$
 - $Y = mX + b$ $\theta = (m, b)$
 - etc...
- Call these “parametric models”
- Given model, parameters yield actual distribution
 - Usually refer to parameters of distribution as θ
 - Note that θ that can be a vector of parameters

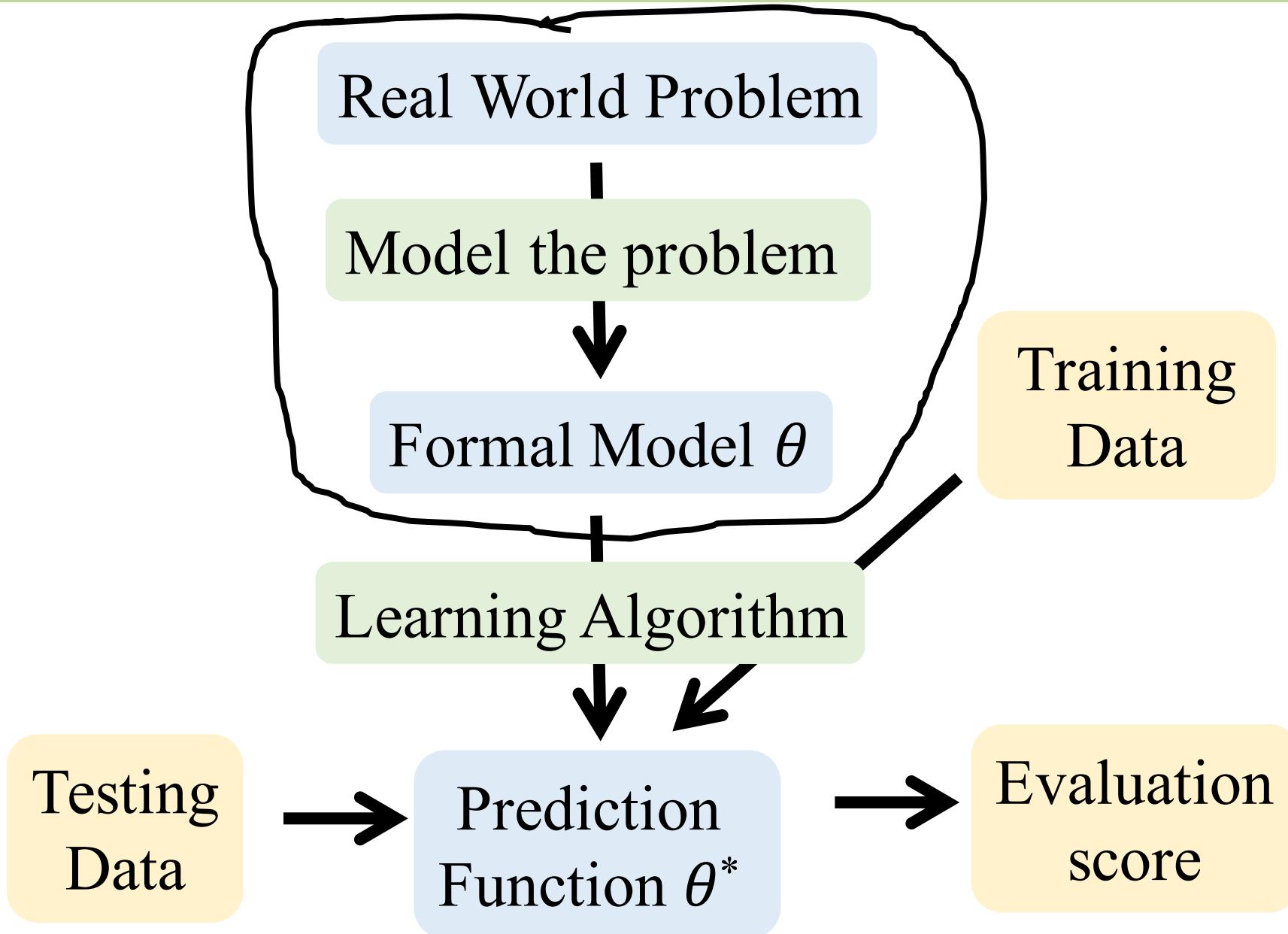
Why Do We Care?

- In real world, don't know "true" parameters
 - But, we do get to observe data
 - E.g., number of times coin comes up heads, lifetimes of disk drives produced, number of visitors to web site per day, etc.
 - Need to estimate model parameters from data
 - "Estimator" is random variable estimating parameter
- Estimate of parameters allows:
 - Better understanding of process producing data
 - Future predictions based on model
 - Simulation of processes

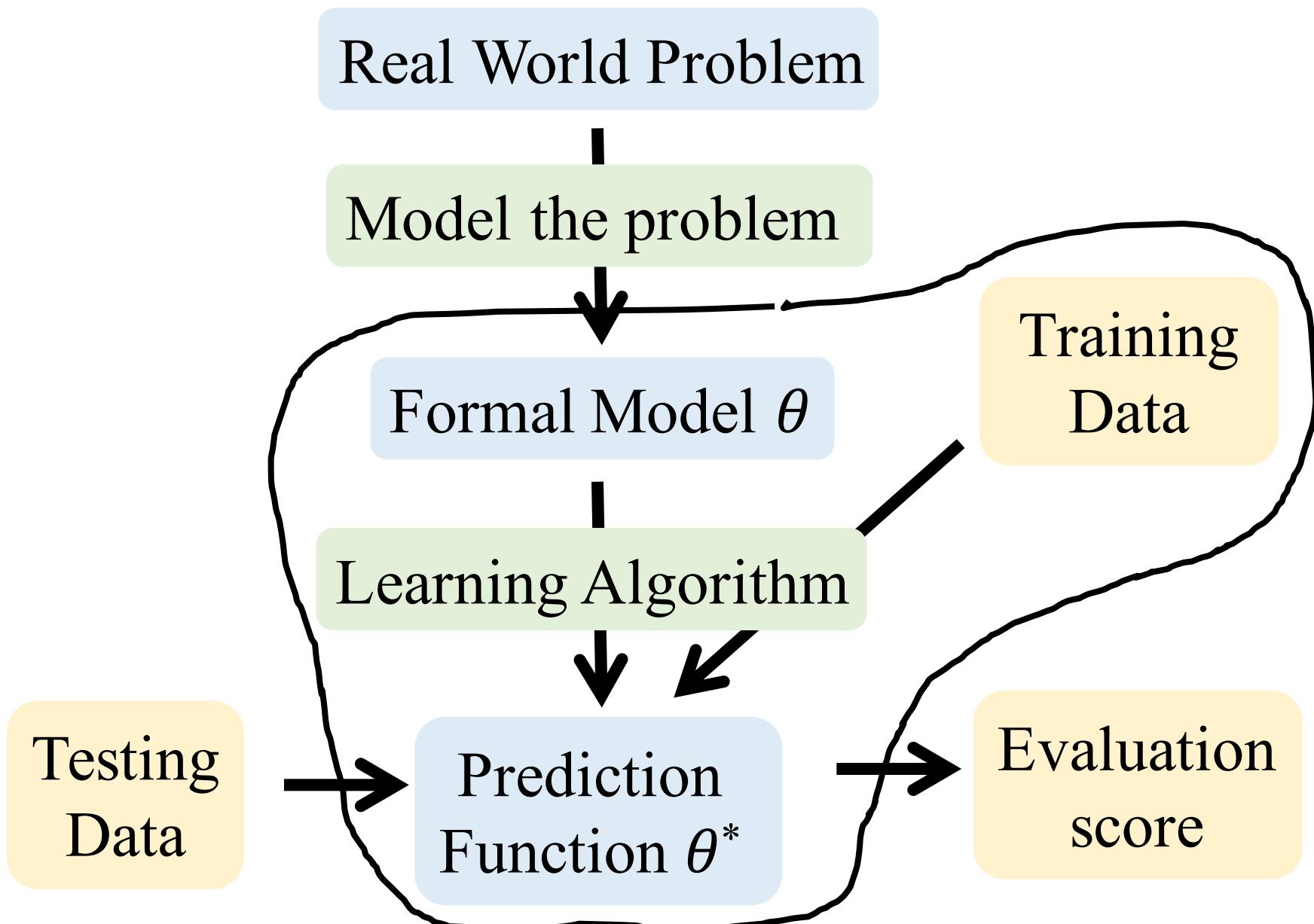
Supervised Learning



Modelling



Training



Testing

Real World Problem

Model the problem

Formal Model θ

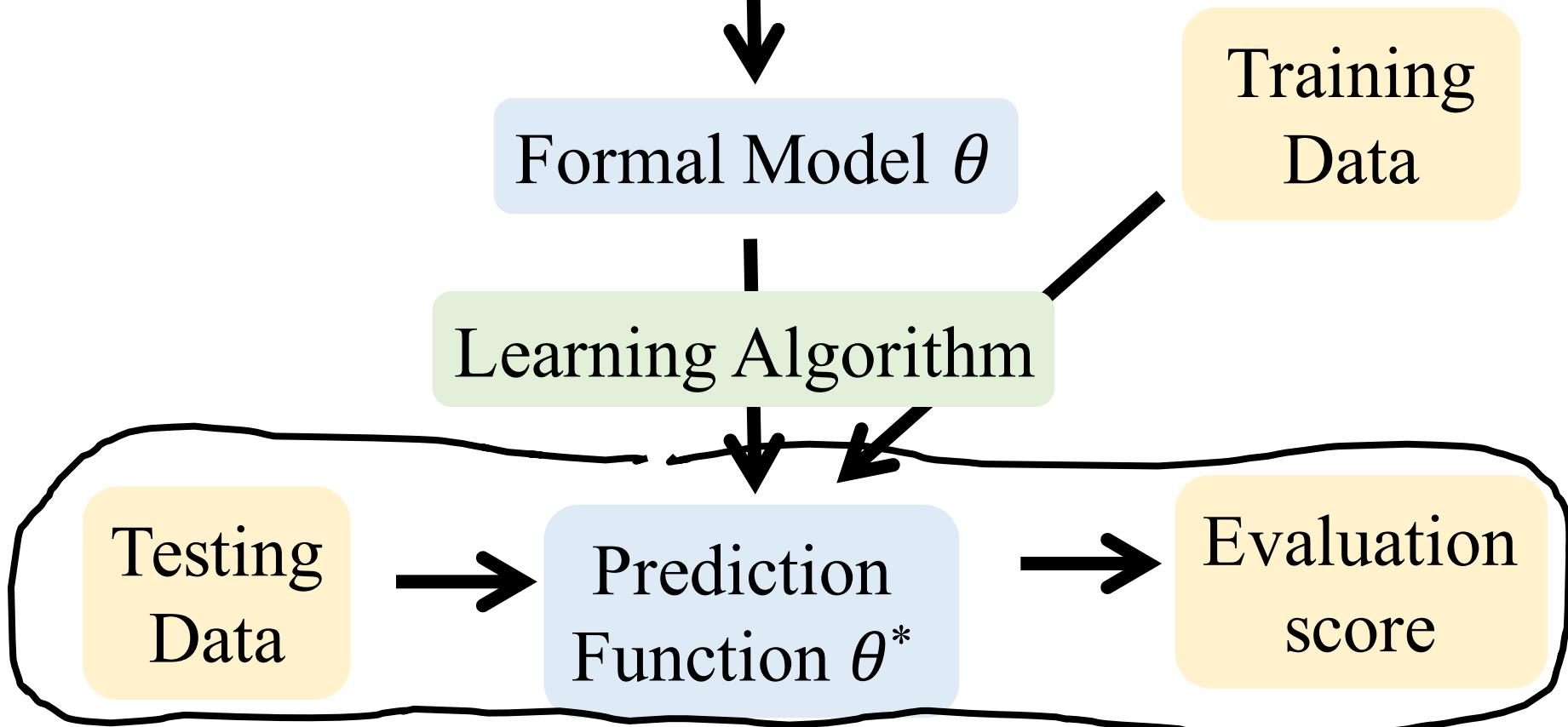
Training
Data

Learning Algorithm

Testing
Data

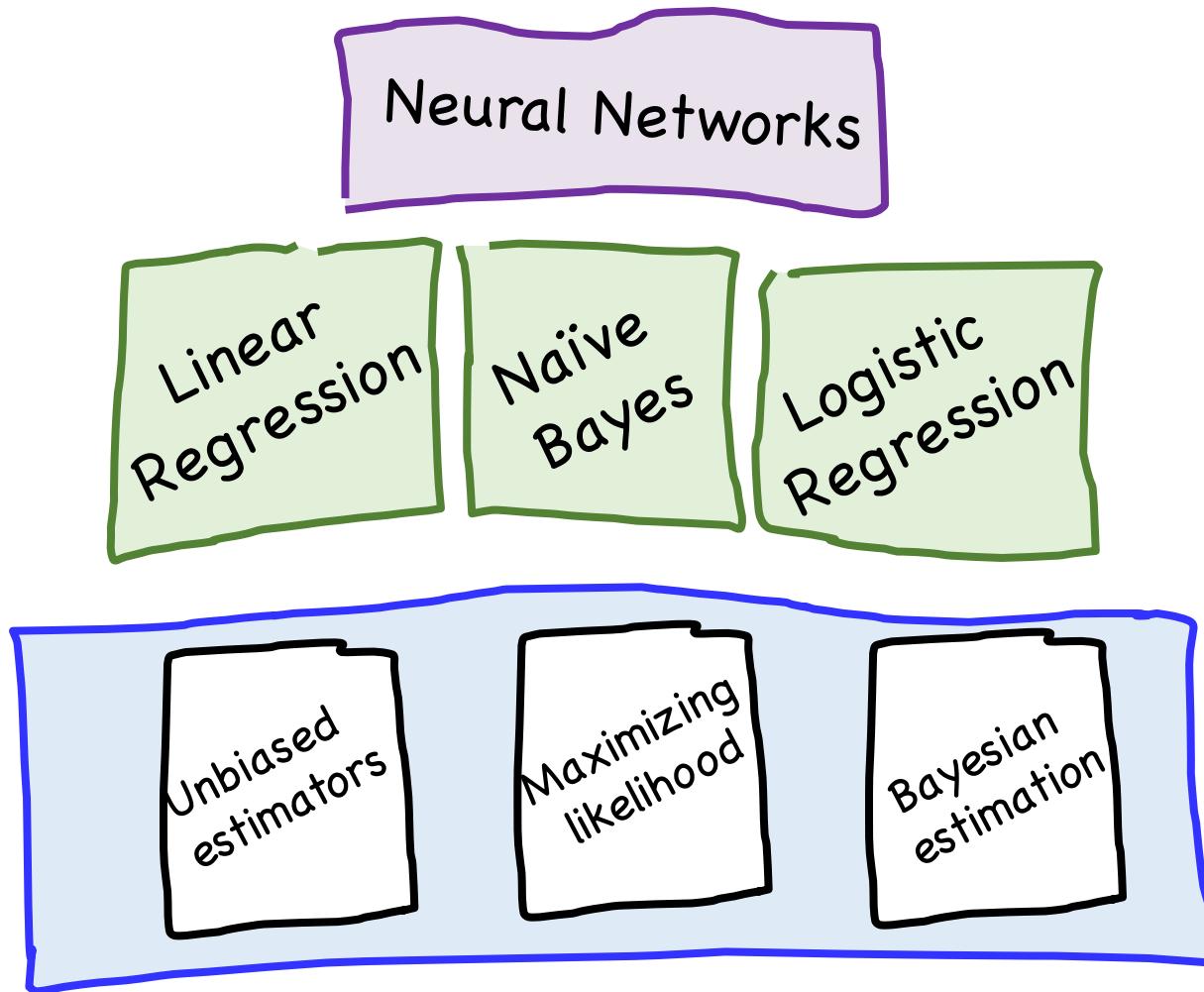
Prediction
Function θ^*

Evaluation
score

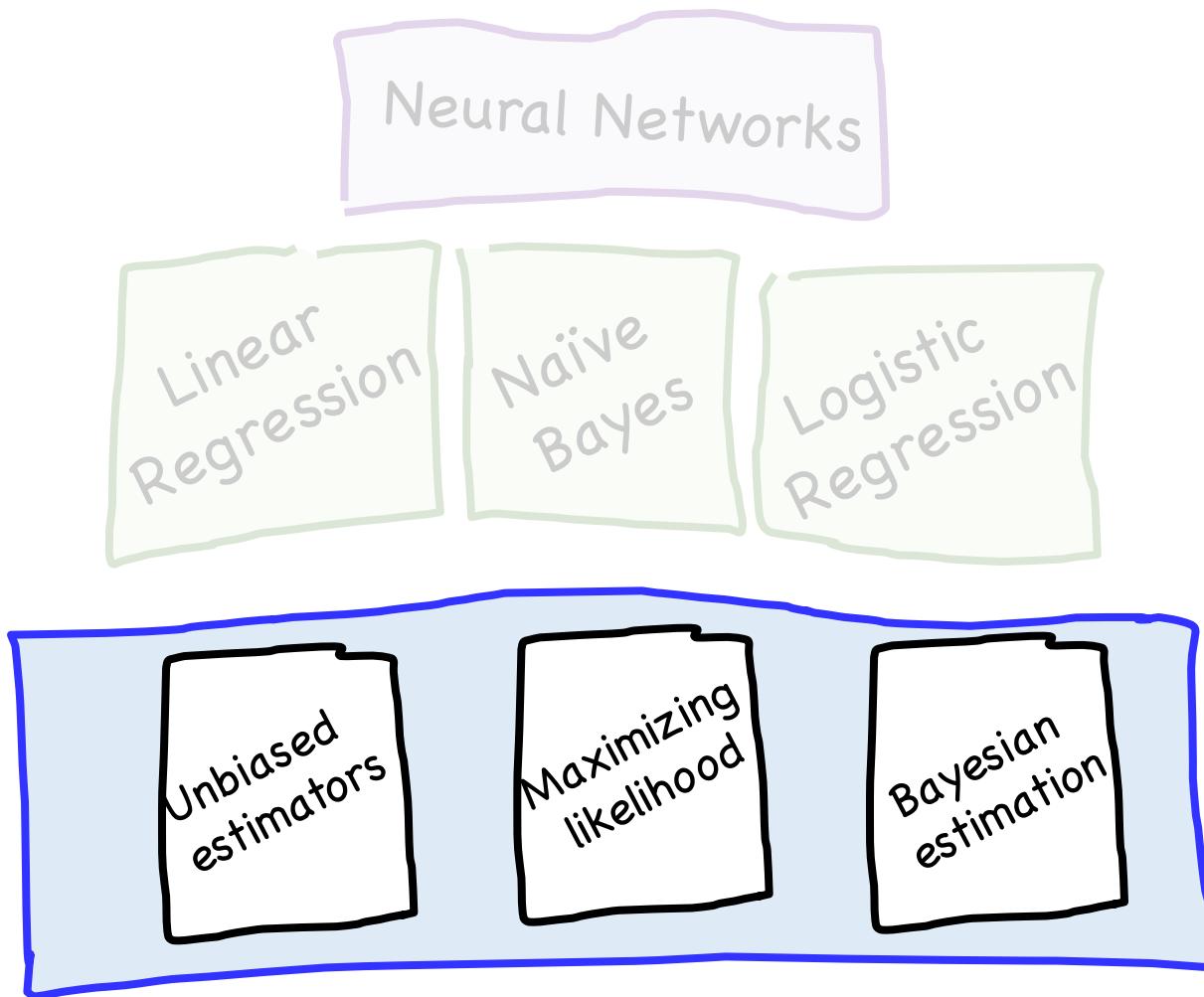


Basis for learning from data

Our Path



Parameter Estimation



Recall Sample Mean + Variance?

- Consider n I.I.D. random variables X_1, X_2, \dots, X_n
 - X_i have distribution F with $E[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2$
 - We call sequence of X_i a **sample** from distribution F
 - Recall sample mean: $\bar{X} = \sum_{i=1}^n \frac{X_i}{n}$ where $E[\bar{X}] = \mu$
$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \text{ as } n \rightarrow \infty$$
 - Recall sample variance:

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

Estimate parameters for
Bernoulli and Normal

Limited tool: how could we use that for
fitting a beta?

Great idea in Machine Learning

Likelihood of Data

- Consider n I.I.D. random variables X_1, X_2, \dots, X_n
 - X_i is a sample from density function $f(X_i | \theta)$
 - Note: now explicitly specify parameter θ of distribution



Likelihood question:
How likely is the data given the samples?

$$\text{Likelihood}(\theta) = f(\text{Samples}|\theta)$$

[Demo](#)



Likelihood of Data

- Consider n I.I.D. random variables X_1, X_2, \dots, X_n
 - X_i is a sample from density function $f(X_i | \theta)$
 - Note: now explicitly specify parameter θ of distribution
 - We want to determine how “likely” the observed data (x_1, x_2, \dots, x_n) is based on density $f(X_i | \theta)$
 - Define the Likelihood function, $L(\theta)$:

$$L(\theta) = \prod_{i=1}^n f(X_i | \theta)$$

- This is just a product since X_i are I.I.D.
- Intuitively: what is probability of observed data using density function $f(X_i | \theta)$, for some choice of θ

Maximum Likelihood Estimator

- The **Maximum Likelihood Estimator** (MLE) of θ , is the value of θ that maximizes $L(\theta)$
 - More formally: $\theta_{MLE} = \arg \max_{\theta} L(\theta)$

Argmax

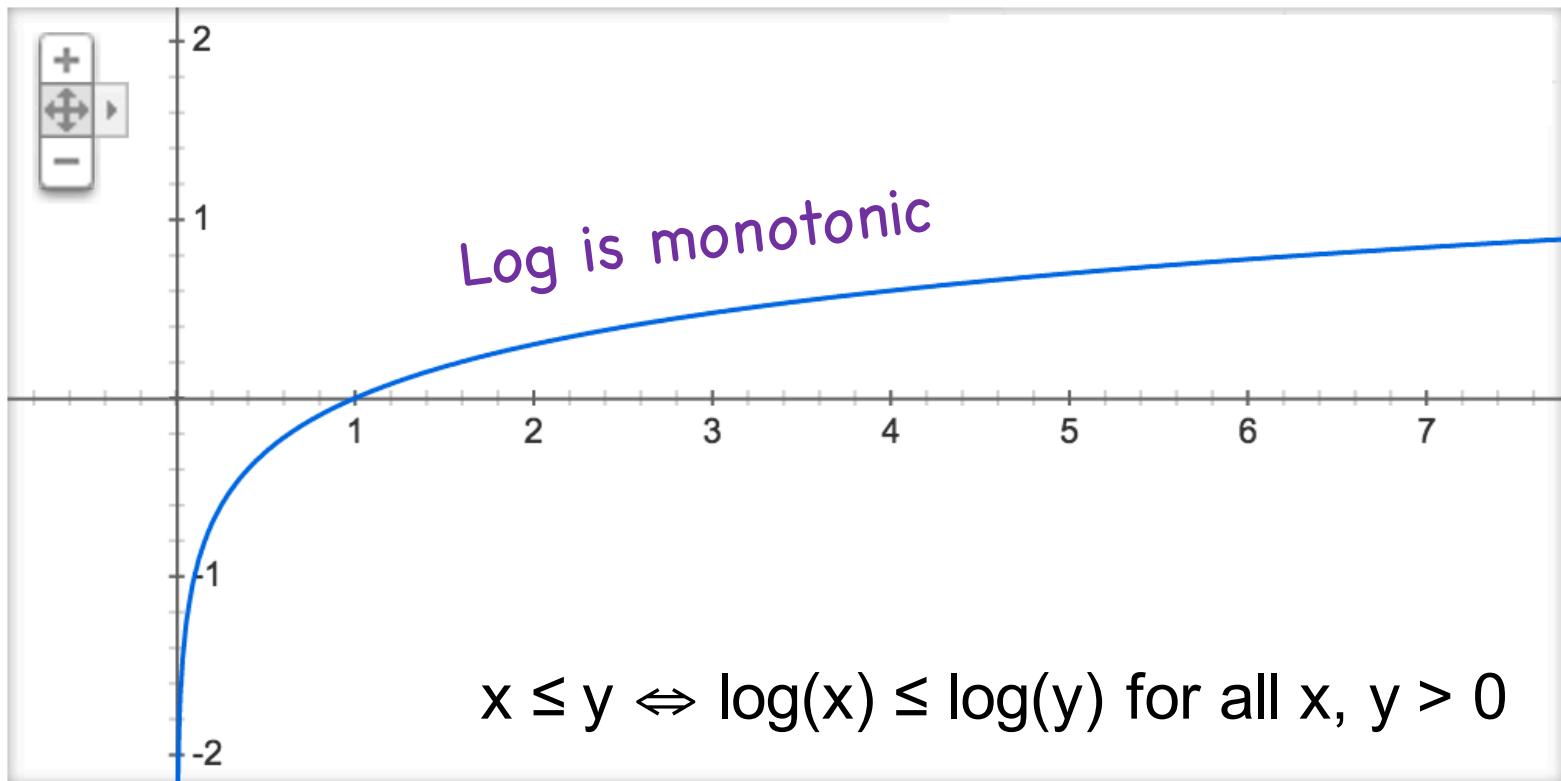
$$f(x) = -x^2 + 5$$

$$\max_x -x^2 + 5 = 5$$

$$\operatorname{argmax}_x -x^2 + 5 = 0$$

Argmax of Log

Graph for $\log(x)$



Claim: $\operatorname{argmax}_x f(x) = \operatorname{argmax}_x \log f(x)$

Argmax of Log

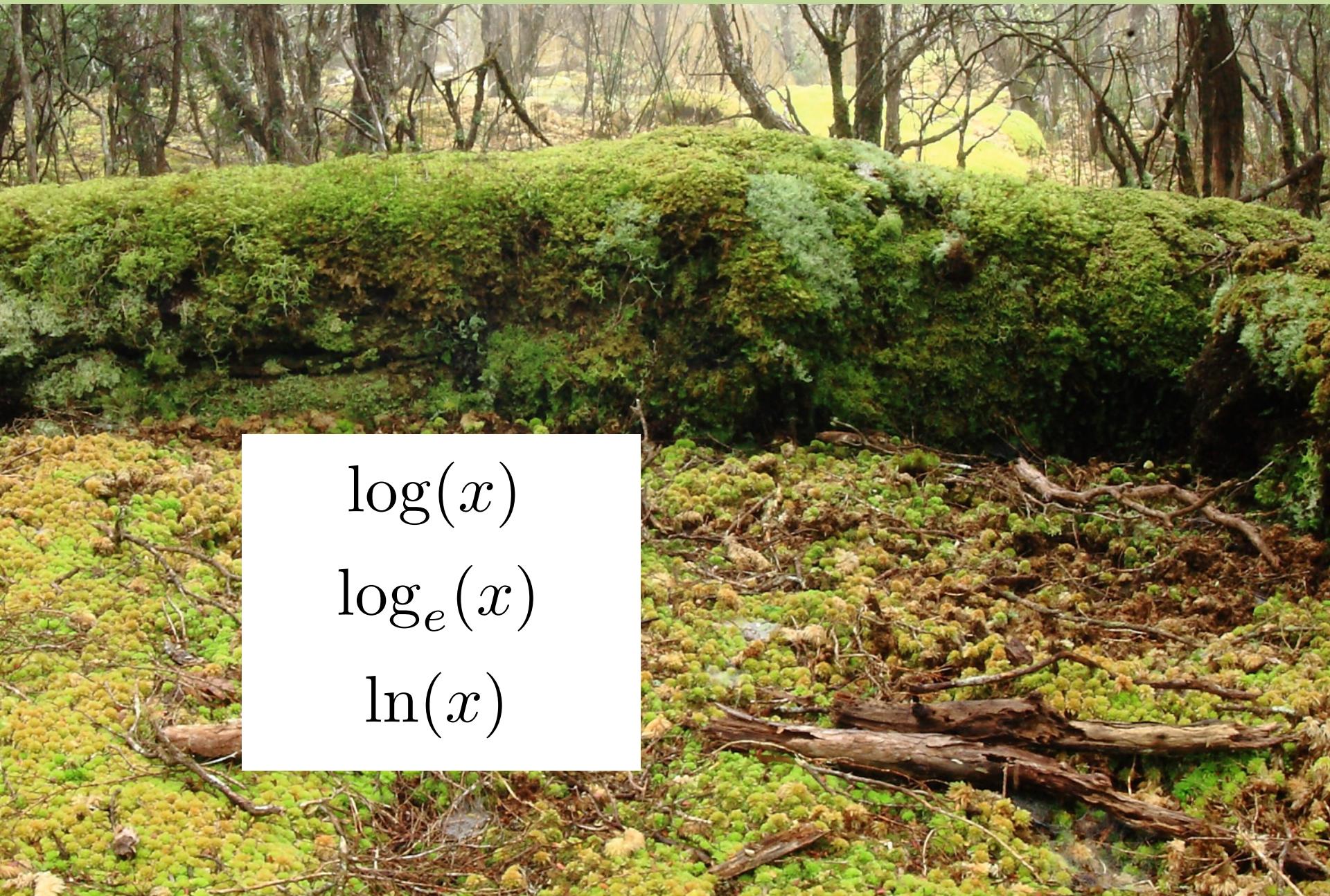


$$\operatorname{argmax}_x f(x) = \operatorname{argmax}_x \log f(x)$$

Log I Love You

$$\log(ab) = \log(a) + \log(b)$$

Natural Log



$\log(x)$

$\log_e(x)$

$\ln(x)$

Maximum Likelihood Estimator

- The **Maximum Likelihood Estimator** (MLE) of θ , is the value of θ that maximizes $L(\theta)$
 - More formally: $\theta_{MLE} = \arg \max_{\theta} L(\theta)$
 - More convenient to use **log-likelihood function**, $LL(\theta)$:

$$LL(\theta) = \log L(\theta) = \log \prod_{i=1}^n f(X_i | \theta) = \sum_{i=1}^n \log f(X_i | \theta)$$

- Note that *log* function is “monotone” for positive values
 - Formally: $x \leq y \Leftrightarrow \log(x) \leq \log(y)$ for all $x, y > 0$
- So, θ that maximizes $LL(\theta)$ also maximizes $L(\theta)$
 - Formally: $\arg \max_{\theta} LL(\theta) = \arg \max_{\theta} L(\theta)$
 - Similarly, for any positive constant c (not dependent on θ):
$$\arg \max_{\theta} (c \cdot LL(\theta)) = \arg \max_{\theta} LL(\theta) = \arg \max_{\theta} L(\theta)$$

Story so far: We can chose parameters by
finding the argmax of the log likelihood of our
data



But how do we compute argmax?

Option #1: Straight optimization

Computing the MLE

- General approach for finding MLE of θ
 - Determine formula for $LL(\theta)$
 - Differentiate $LL(\theta)$ w.r.t. (each) θ : $\frac{\partial LL(\theta)}{\partial \theta}$
 - To maximize, set $\frac{\partial LL(\theta)}{\partial \theta} = 0$
 - Solve resulting (simultaneous) equations to get θ_{MLE}
 - Make sure that derived $\hat{\theta}_{MLE}$ is actually a maximum (and not a minimum or saddle point). E.g., check $LL(\theta_{MLE} \pm \varepsilon) < LL(\theta_{MLE})$
 - This step often ignored in expository derivations
 - So, we'll ignore it here too (and won't require it in this class)

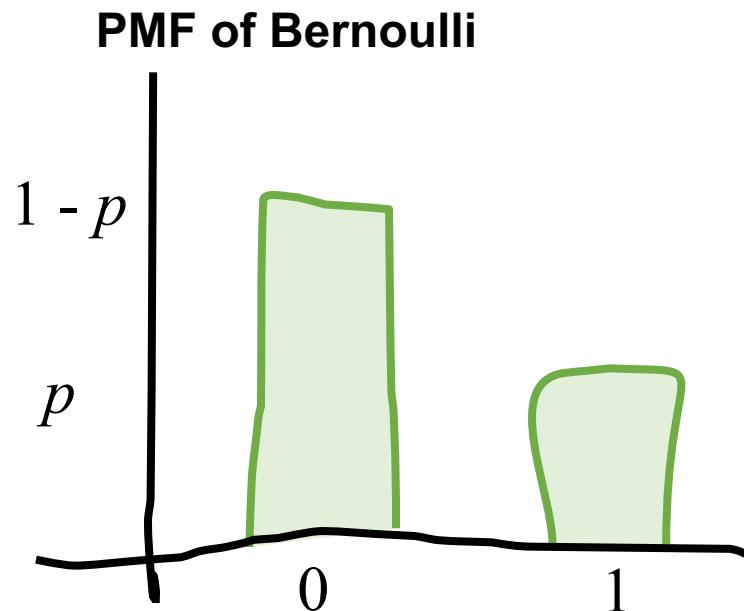
Maximizing Likelihood with Bernoulli

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Ber}(p)$
 - Probability mass function, $f(X_i | p)$:

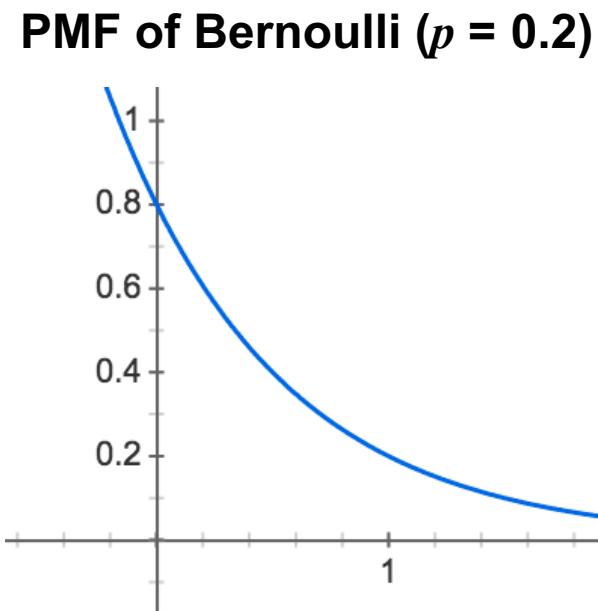


Maximizing Likelihood with Bernoulli

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Ber}(p)$
 - Probability mass function, $f(X_i | p)$:



$$f(X_i | p) = p^{x_i} (1-p)^{1-x_i}$$



$$f(x) = 0.2^x (1 - 0.2)^{1-x}$$

Bernoulli PMF

$$X \sim \text{Ber}(p)$$



$$f(X = x|p) = p^x(1 - p)^{1-x}$$

Maximizing Likelihood with Bernoulli

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Ber}(p)$
 - Probability mass function, $f(X_i | p)$, can be written as:

$$f(X_i | p) = p^{x_i} (1 - p)^{1-x_i} \quad \text{where } x_i = 0 \text{ or } 1$$

- Likelihood: $L(\theta) = \prod_{i=1}^n p^{X_i} (1 - p)^{1-X_i}$

- Log-likelihood:

$$\begin{aligned} LL(\theta) &= \sum_{i=1}^n \log(p^{X_i} (1 - p)^{1-X_i}) = \sum_{i=1}^n [X_i(\log p) + (1 - X_i)\log(1 - p)] \\ &= Y(\log p) + (n - Y)\log(1 - p) \quad \text{where } Y = \sum_{i=1}^n X_i \end{aligned}$$

- Differentiate w.r.t. p , and set to 0:

$$\frac{\partial LL(p)}{\partial p} = Y \frac{1}{p} + (n - Y) \frac{-1}{1 - p} = 0 \quad \Rightarrow \quad p_{MLE} = \frac{Y}{n} = \frac{1}{n} \sum_{i=1}^n X_i$$

Isn't that the same as
unbiased estimator?

Yes. For Bernoulli.



Maximum Likelihood Algorithm

1. Decide on a model for the distribution of your samples. Define the PMF / PDF for your sample.
2. Write out the log likelihood function.
3. State that the optimal parameters are the argmax of the log likelihood function.
4. Use an optimization algorithm to calculate argmax



Maximizing Likelihood with Poisson

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Poi}(\lambda)$
 - PMF: $f(X_i | \lambda) = \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$ Likelihood: $L(\theta) = \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{X_i}}{X_i!}$
 - Log-likelihood:
$$\begin{aligned} LL(\theta) &= \sum_{i=1}^n \log\left(\frac{e^{-\lambda} \lambda^{X_i}}{X_i!}\right) = \sum_{i=1}^n [-\lambda \log(e) + X_i \log(\lambda) - \log(X_i!)] \\ &= -n\lambda + \log(\lambda) \sum_{i=1}^n X_i - \sum_{i=1}^n \log(X_i!) \end{aligned}$$
 - Differentiate w.r.t. λ , and set to 0:

$$\frac{\partial LL(\lambda)}{\partial \lambda} = -n + \frac{1}{\lambda} \sum_{i=1}^n X_i = 0 \quad \Rightarrow \quad \lambda_{MLE} = \frac{1}{n} \sum_{i=1}^n X_i$$

Its so general!

Maximizing Likelihood with Normal

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim N(\mu, \sigma^2)$
 - PDF: $f(X_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(X_i - \mu)^2 / (2\sigma^2)}$
 - Log-likelihood:

$$LL(\theta) = \sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi}\sigma} e^{-(X_i - \mu)^2 / (2\sigma^2)}\right) = \sum_{i=1}^n \left[-\log(\sqrt{2\pi}\sigma) - (X_i - \mu)^2 / (2\sigma^2) \right]$$

- First, differentiate w.r.t. μ , and set to 0:

$$\frac{\partial LL(\mu, \sigma^2)}{\partial \mu} = \sum_{i=1}^n 2(X_i - \mu) / (2\sigma^2) = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = 0$$

- Then, differentiate w.r.t. σ , and set to 0:

$$\frac{\partial LL(\mu, \sigma^2)}{\partial \sigma} = \sum_{i=1}^n -\frac{1}{\sigma} + 2(X_i - \mu)^2 / (2\sigma^3) = -\frac{n}{\sigma} + \sum_{i=1}^n (X_i - \mu)^2 / (\sigma^3) = 0$$

Being Normal, Simultaneously

- Now have two equations, two unknowns:

$$\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = 0 \quad -\frac{n}{\sigma} + \sum_{i=1}^n (X_i - \mu)^2 / (\sigma^3) = 0$$

- First, solve for μ_{MLE} :

$$\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = 0 \Rightarrow \sum_{i=1}^n X_i = n\mu \Rightarrow \mu_{MLE} = \frac{1}{n} \sum_{i=1}^n X_i$$

- Then, solve for σ^2_{MLE} :

$$-\frac{n}{\sigma} + \sum_{i=1}^n (X_i - \mu)^2 / (\sigma^3) = 0 \Rightarrow n\sigma^2 = \sum_{i=1}^n (X_i - \mu)^2$$

$$\sigma^2_{MLE} = \frac{1}{n} \sum_{i=1}^n (X_i - \mu_{MLE})^2$$

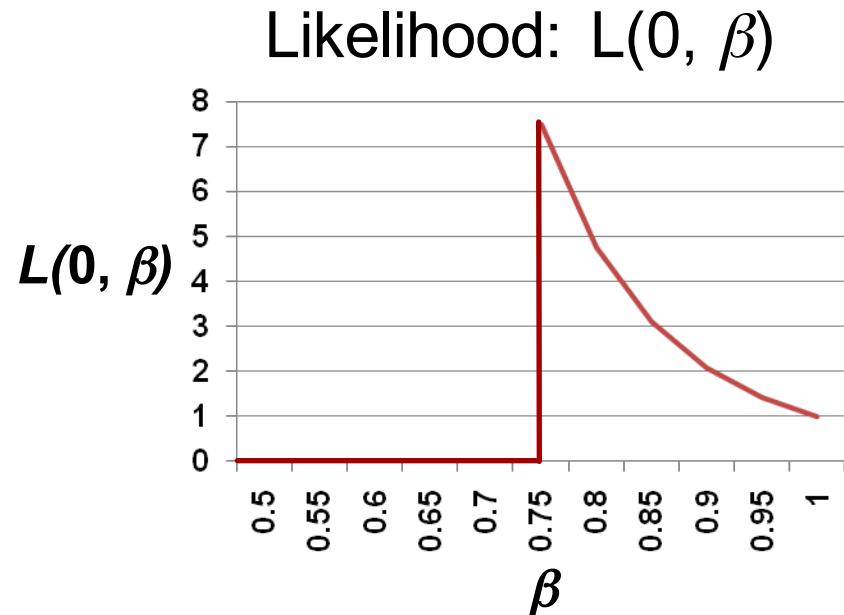
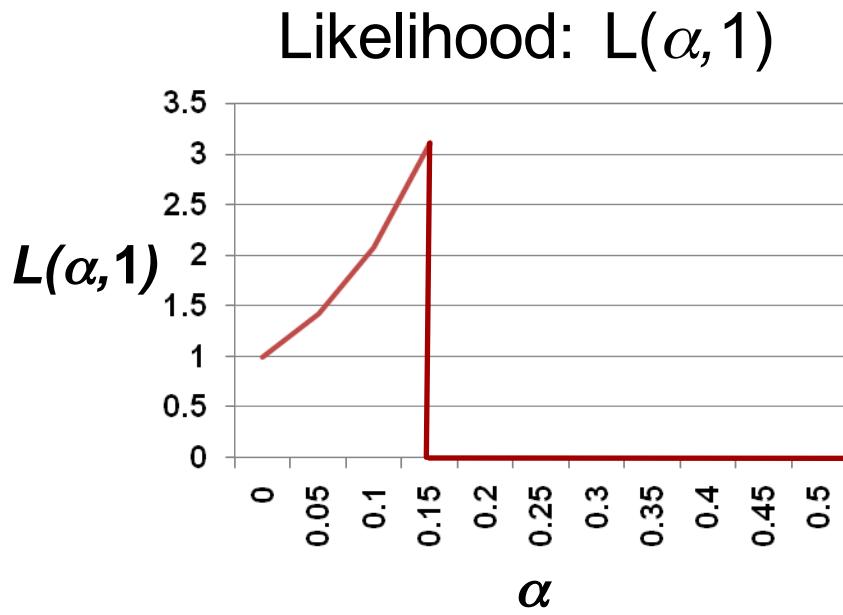
- Note: μ_{MLE} unbiased, but σ^2_{MLE} biased

Maximizing Likelihood with Uniform

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Uni}(\alpha, \beta)$
 - PDF: $f(X_i | \alpha, \beta) = \begin{cases} \frac{1}{\beta-\alpha} & \alpha \leq x_i \leq \beta \\ 0 & \text{otherwise} \end{cases}$
 - Likelihood: $L(\theta) = \begin{cases} \left(\frac{1}{\beta-\alpha}\right)^n & \alpha \leq x_1, x_2, \dots, x_n \leq \beta \\ 0 & \text{otherwise} \end{cases}$
 - Constraint $\alpha \leq x_1, x_2, \dots, x_n \leq \beta$ makes differentiation tricky
 - Intuition: want interval size $(\beta - \alpha)$ to be as small as possible to maximize likelihood function for each data point
 - But need to make sure all observed data contained in interval
 - If all observed data not in interval, then $L(\theta) = 0$
 - Solution: $\alpha_{MLE} = \min(x_1, \dots, x_n)$ $\beta_{MLE} = \max(x_1, \dots, x_n)$

Understanding MLE with Uniform

- Consider I.I.D. random variables X_1, X_2, \dots, X_n
 - $X_i \sim \text{Uni}(0, 1)$
 - Observe data:
 - 0.15, 0.20, 0.30, 0.40, 0.65, 0.70, 0.75



Small Samples = Problems

- How do small samples affect MLE?
 - In many cases, $\mu_{MLE} = \frac{1}{n} \sum_{i=1}^n X_i$ = sample mean
 - Unbiased. Not too shabby...
 - As seen with Normal, $\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu_{MLE})^2$
 - Biased. Underestimates for small n (e.g., 0 for $n = 1$)
 - As seen with Uniform, $\alpha_{MLE} \geq \alpha$ and $\beta_{MLE} \leq \beta$
 - Biased. Problematic for small n (e.g., $\alpha = \beta$ when $n = 1$)
 - Small sample phenomena intuitively make sense:
 - Maximum likelihood \Rightarrow best explain data we've seen
 - Does not attempt to generalize to unseen data

Properties of MLE

- Maximum Likelihood Estimators are generally:
 - Consistent: $\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| < \varepsilon) = 1$ for $\varepsilon > 0$
 - Potentially biased (though asymptotically less so)
 - Asymptotically optimal
 - Has smallest variance of “good” estimators for large samples
 - Often used in practice where sample size is large relative to parameter space
 - But be careful, there are some very large parameter spaces

[if time, MLE of line]