



# Parameter Estimation

Noah Arthurs

CS109, Stanford University

Review

# Bootstrapping for p values

# Null Hypothesis Test

<b>Nepal Happiness</b>	<b>Bhutan Happiness</b>
4.44	2.15
3.36	3.01
5.87	2.02
2.31	1.43
...	...
3.70	1.83

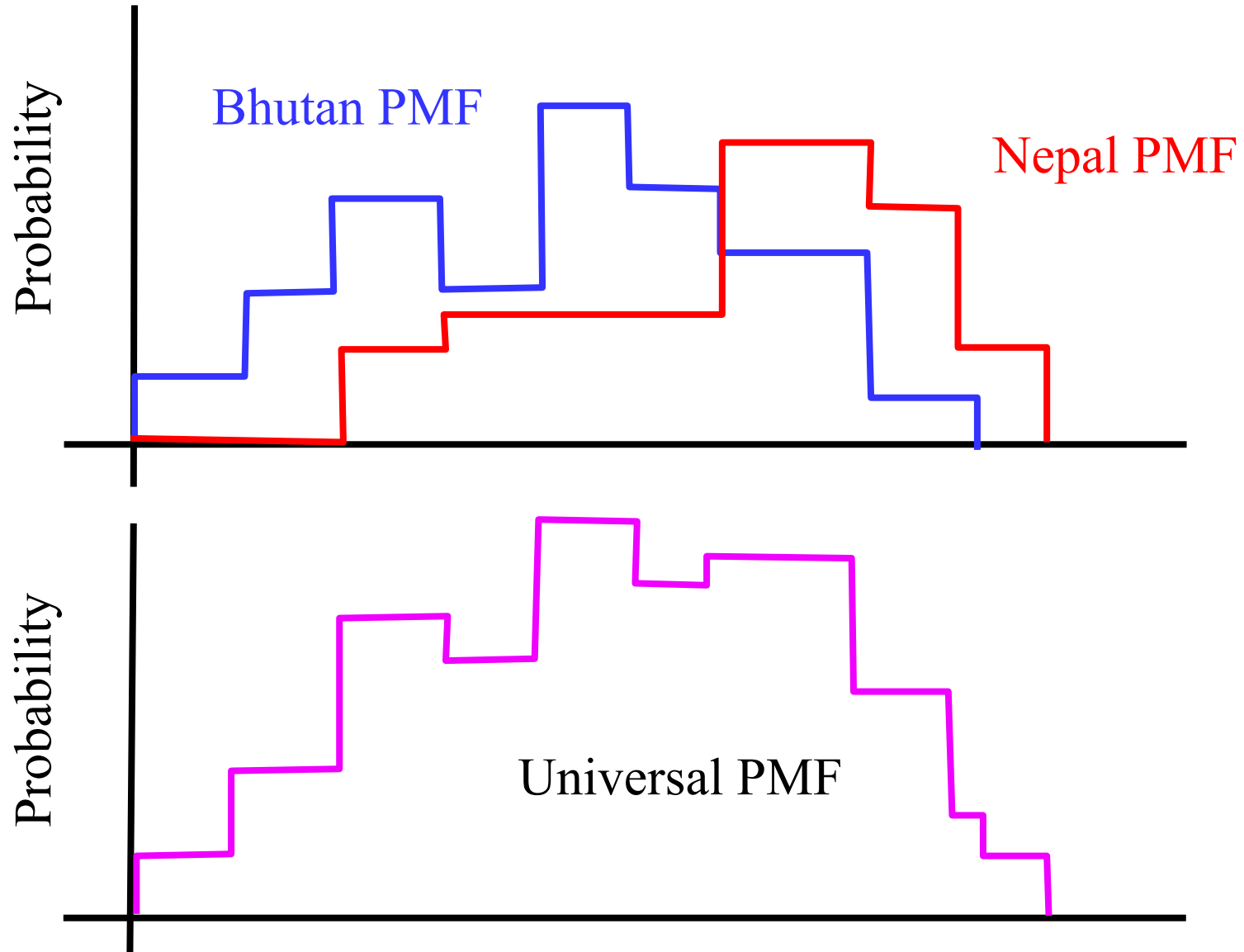
$\mu_1 = 3.1$                        $\mu_2 = 2.4$

Claim: The difference in happiness between Nepal and Bhutan is 0.7 happiness points.



**Null hypothesis:** even if **there is no pattern** (ie the two samples are identically distributed) your claim might have arisen by **chance**.

# Universal Sample



# Bootstrap for P Values

```
def pvalueBootstrap(bhutanSample, nepalSample):  
    N = size of the bhutanSample  
    M = size of the nepalSample  
  
    uniSamples = combine bhutanSamples and nepalSamples  
    count = 0  
  
    repeat 10,000 times:  
        bhutanResample = draw N resamples from the uniSamples  
        nepalResample = draw M resamples from the uniSamples  
        muBhutan = sample mean of the bhutanResample  
        muNepal = sample mean of the nepalResample  
        meanDiff = |muNepal - muBhutan|  
        if meanDiff > observedDifference:  
            count += 1  
  
    pValue = count / 10,000
```

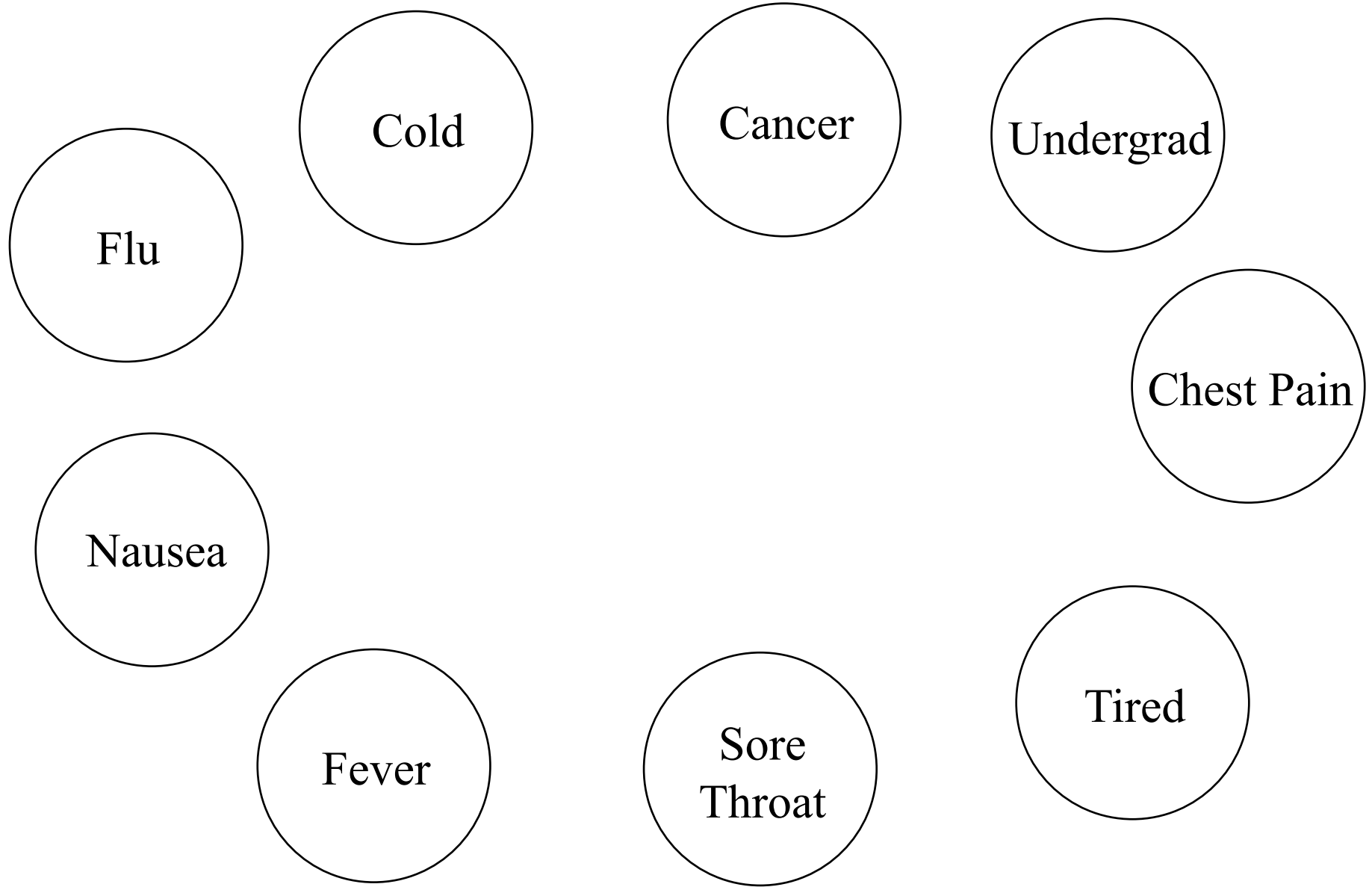
# Null Hypothesis Test

<u>Nepal Happiness</u>	<u>Bhutan Happiness</u>
4.44	2.15
3.36	3.01
5.87	2.02
2.31	1.43
...	...
3.70	1.83

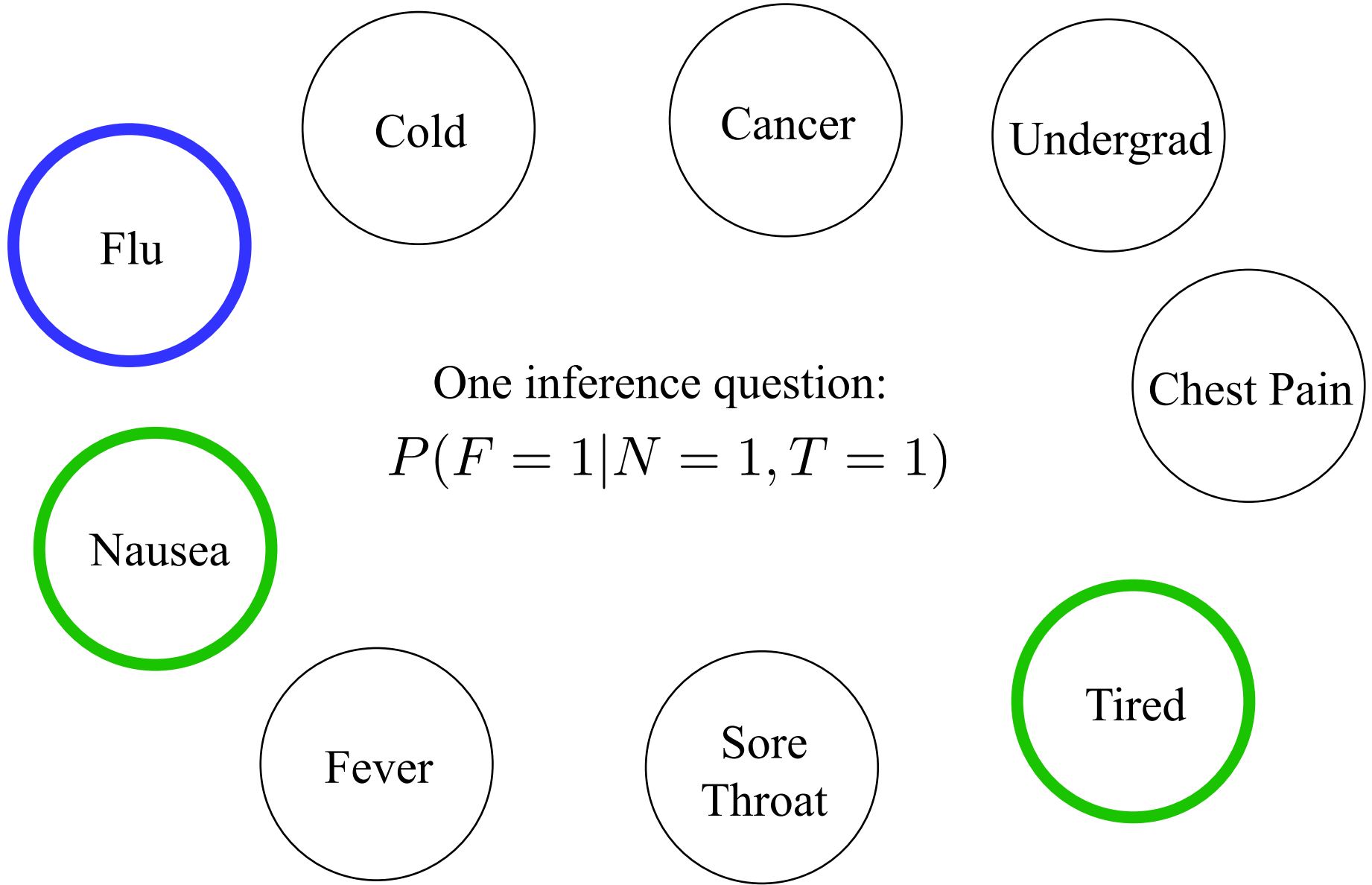
$\mu_1 = 3.1$                        $\mu_2 = 2.4$

**Claim: The difference in happiness between Nepal and Bhutan is 0.7 happiness points ( $p < 0.008$ ).**

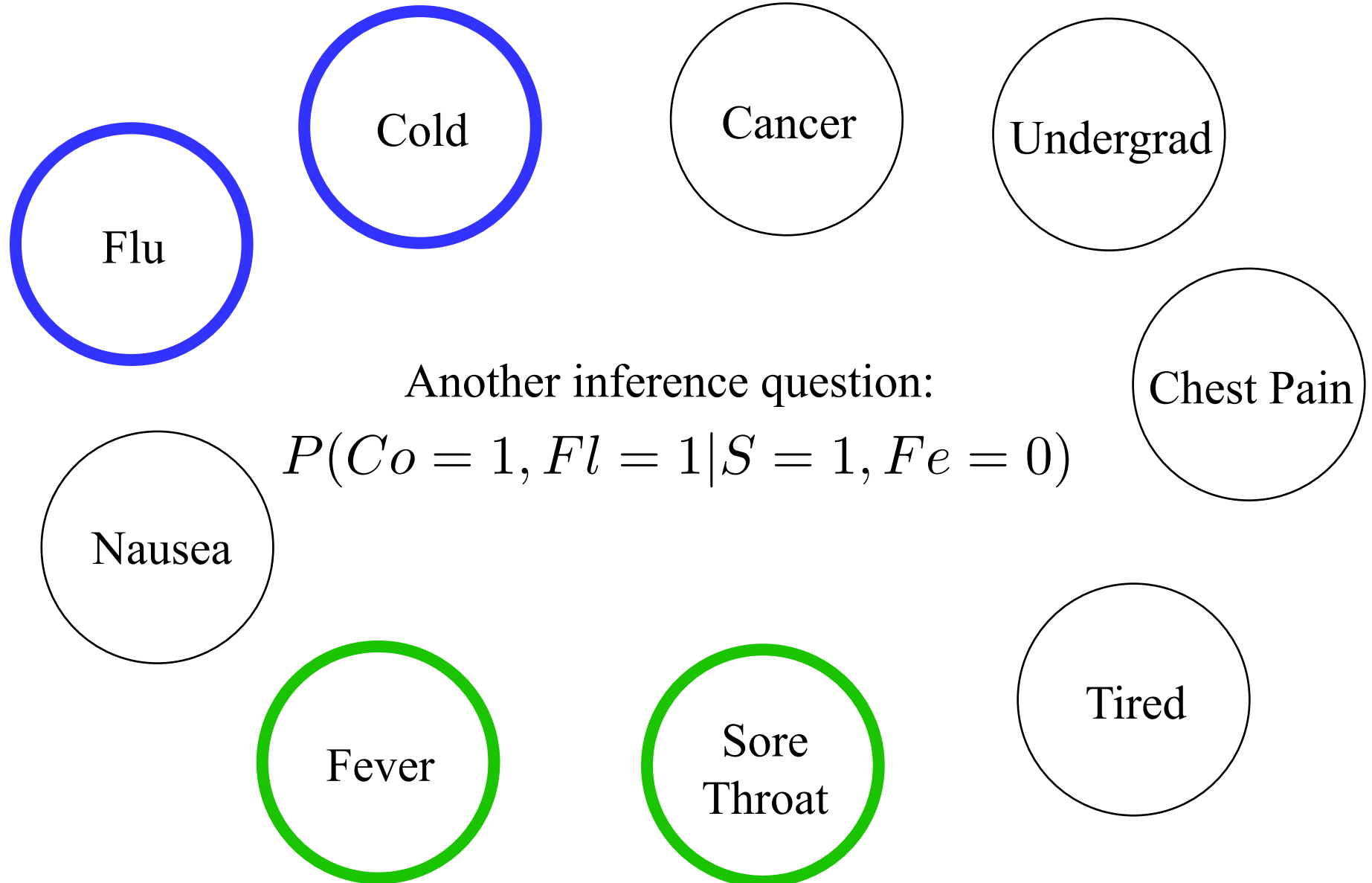
# General “Inference”



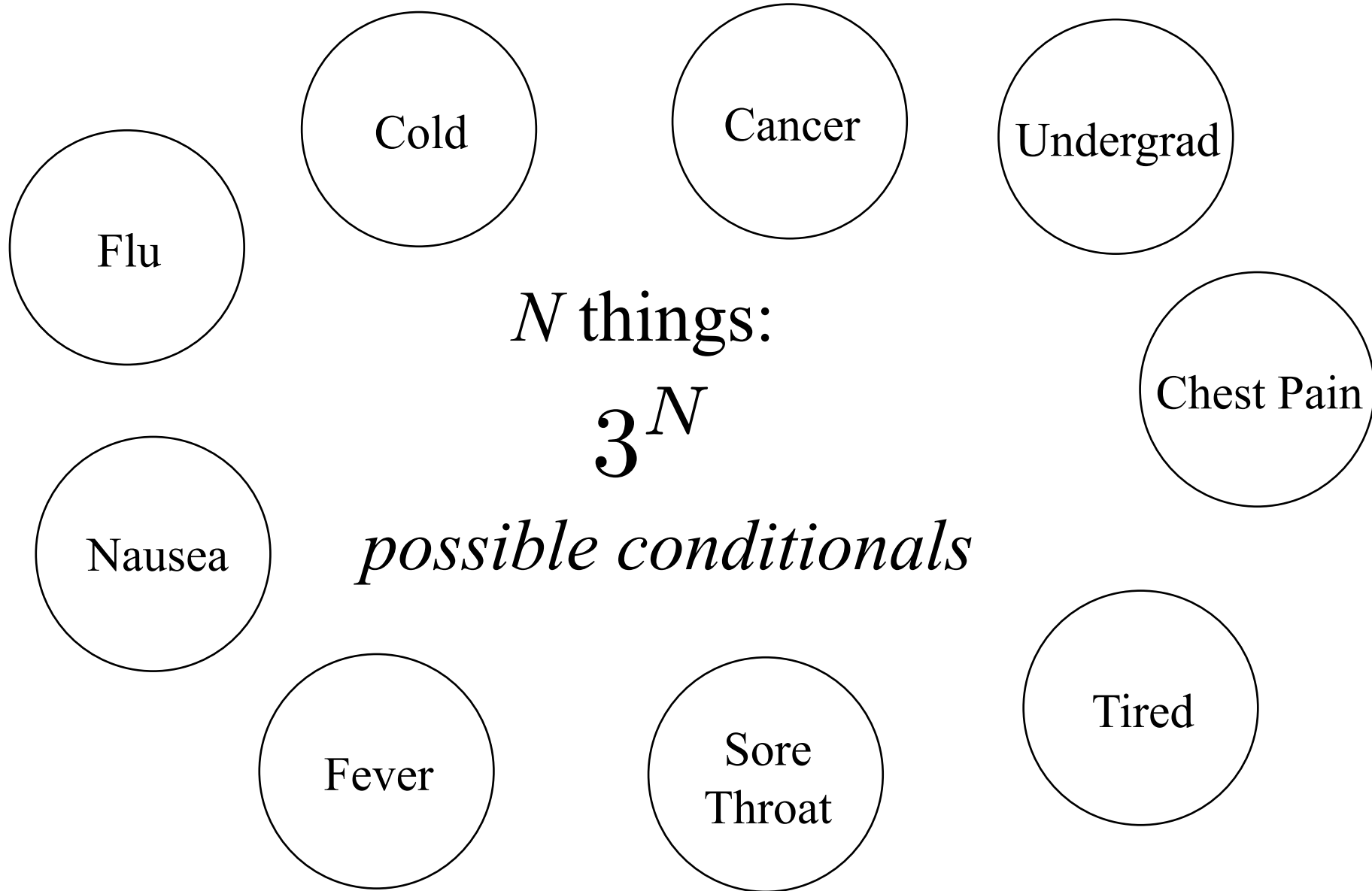
# General “Inference”



# General “Inference”



# How Many Things Can You Condition On?



Flu

Cold

Cancer

Undergrad

$N$  things:

$$3^N$$

Chest Pain

Nausea

*possible conditionals*

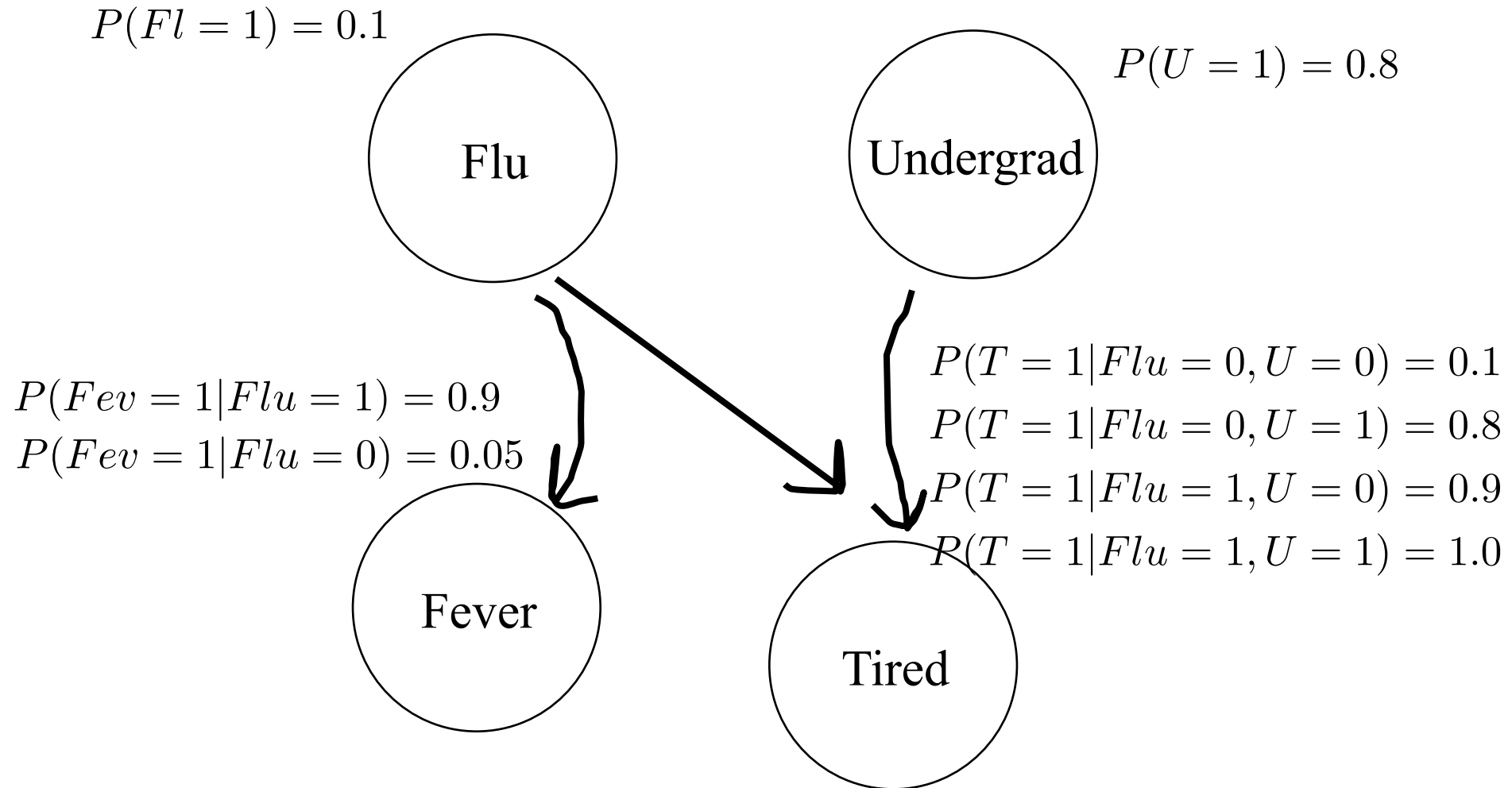
Tired

Fever

Sore  
Throat

# Bayesian Network

Describe the joint using causality!



$$P(Fl = a, Fe = b, U = c, T = d)?$$

## Bayesian Network:

*If you know causality,*



Make a network of assumed direct causality for your random vars.

You simply need to give:

$P(\text{values} \mid \text{parents})$

for each random variable.

Prob can be a conditional probability table or  
an equation!

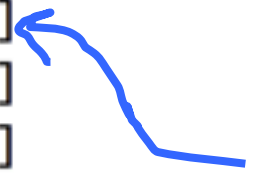
# Alg #1: Joint Sampling

```
3 N_SAMPLES = 100000
4
5 # Program: Joint Sample
6 # -----
7 # we can answer any pro
8 # with multivariate sam
9 # where conditioned var
10 def main():
11     obs = getObservatio
12     print 'Observation
13
14     samples = sampleATo
15     prob = probFluGiven
16     print 'Pr(Flu) = ',
```

```
webMd -- -bash -- 30x20
[0, 1, 0, 1]
[1, 1, 1, 1]
[0, 1, 0, 1]
[0, 1, 0, 0]
[0, 1, 0, 0]
[0, 1, 0, 1]
[0, 1, 0, 1]
[0, 0, 0, 0]
[0, 0, 0, 0]
[0, 1, 0, 1]
[0, 1, 0, 1]
[0, 1, 0, 1]
[0, 1, 0, 1]
[0, 1, 0, 1]
[0, 1, 0, 1]
[0, 1, 0, 0]
[0, 0, 0, 0]
[0, 1, 0, 1]
[0, 1, 0, 1]
[0, 1, 0, 1]
[1, 1, 0, 1]
```

```
webMd — -bash — 39x20
[0, 1, 1, 0]
[1, 0, 1, 1]
[0, 1, 0, 1]
[0, 1, 0, 0]
[0, 1, 0, 0]
[0, 1, 1, 0]
[1, 1, 1, 1]
[0, 1, 0, 0]
[0, 0, 0, 1]
[0, 1, 0, 1]
[0, 1, 0, 1]
[0, 1, 0, 1]
[0, 1, 0, 0]
[0, 1, 0, 1]
[0, 1, 0, 0]
[0, 0, 0, 0]
[0, 0, 0, 1]
Observation = [None, None, None, None]
Pr(Flu | Obs) = 0.10164
>
```

If you can sample enough from the joint distribution, you can answer any probability question



Each one of these is one joint sample:  
[Flu, Undergrad, Fever, Tired]



# BAE's Theorem?

$$P(A | B E) = \frac{P(B | A E) P(A | E)}{P(B | E)}$$



# $P(F = 1 \mid \text{all other rvs})$

Know:  $P(\text{symptom} \mid \text{flu}, \text{undergrad})$        $P(\text{flu})$        $P(\text{undergrad})$

Flu is independent of undergrad

Tired and fever are conditionally independent given flu, undergrad

---

$$\begin{aligned} &P(F = 1 \mid \text{symptoms}, U = u) \\ &= \frac{P(\text{symptoms} \mid F = 1, U = u)P(F = 1 \mid U = u)}{P(\text{symptoms} \mid U = u)} \\ &\propto P(\text{symptoms} \mid F = 1, U = u)P(F = 1 \mid U = u) \\ &\propto P(F = 1)P(\text{symptoms} \mid F = 1, U = u) \\ &\propto P(F = 1) \prod_i P(\text{symptom}_i \mid F = 1, U = u) \end{aligned}$$

# General Inference Summary

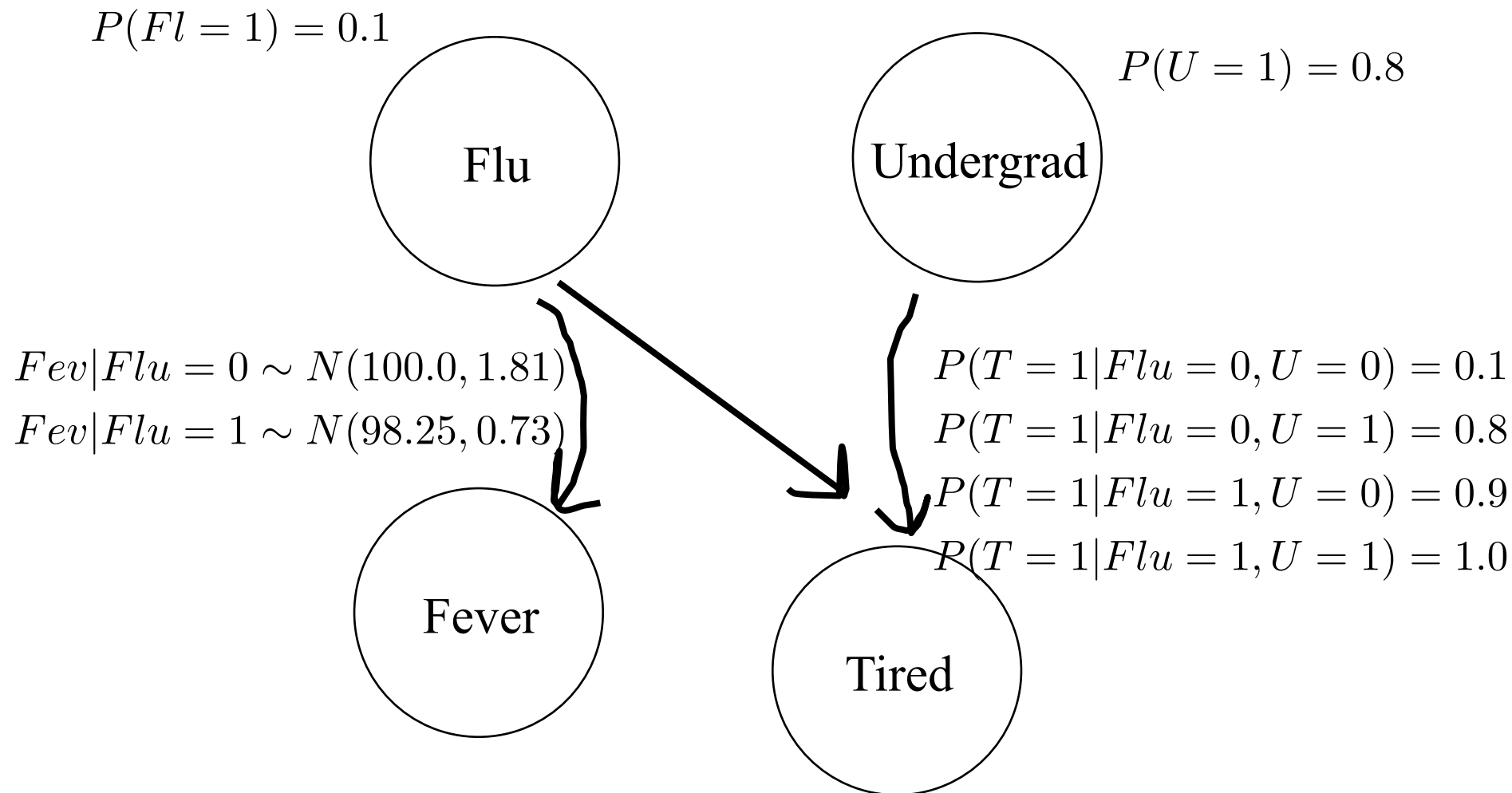


- **Straight Math** is fast, but can be prohibitively hard for complex models (see hw).
- **Joint Sampling** is really easy to program but fails for continuous variables (and when what you are conditioning on is rare)
- **Take CS228** to learn about many, many more algorithms!

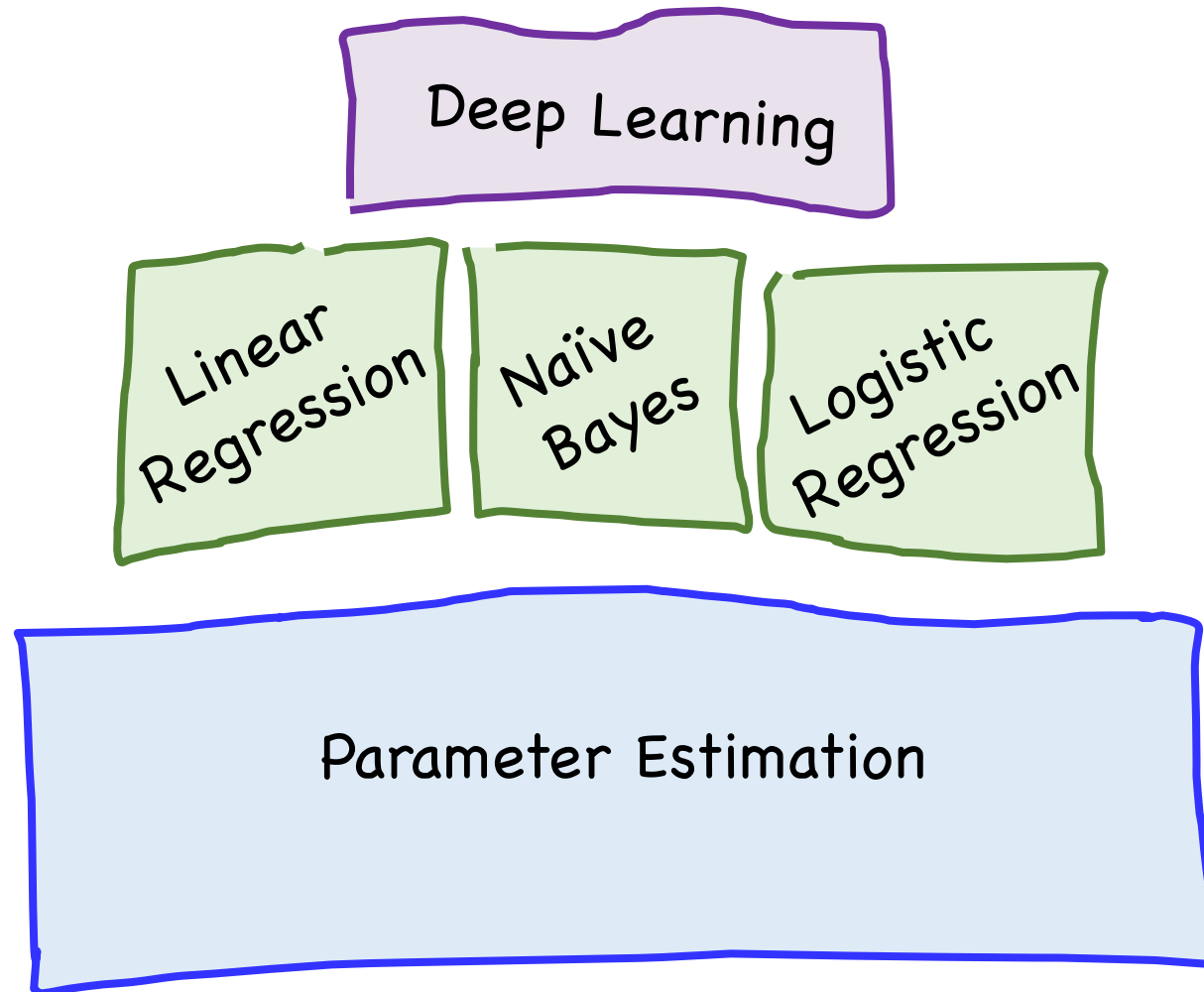
# MACHINE LEARNING



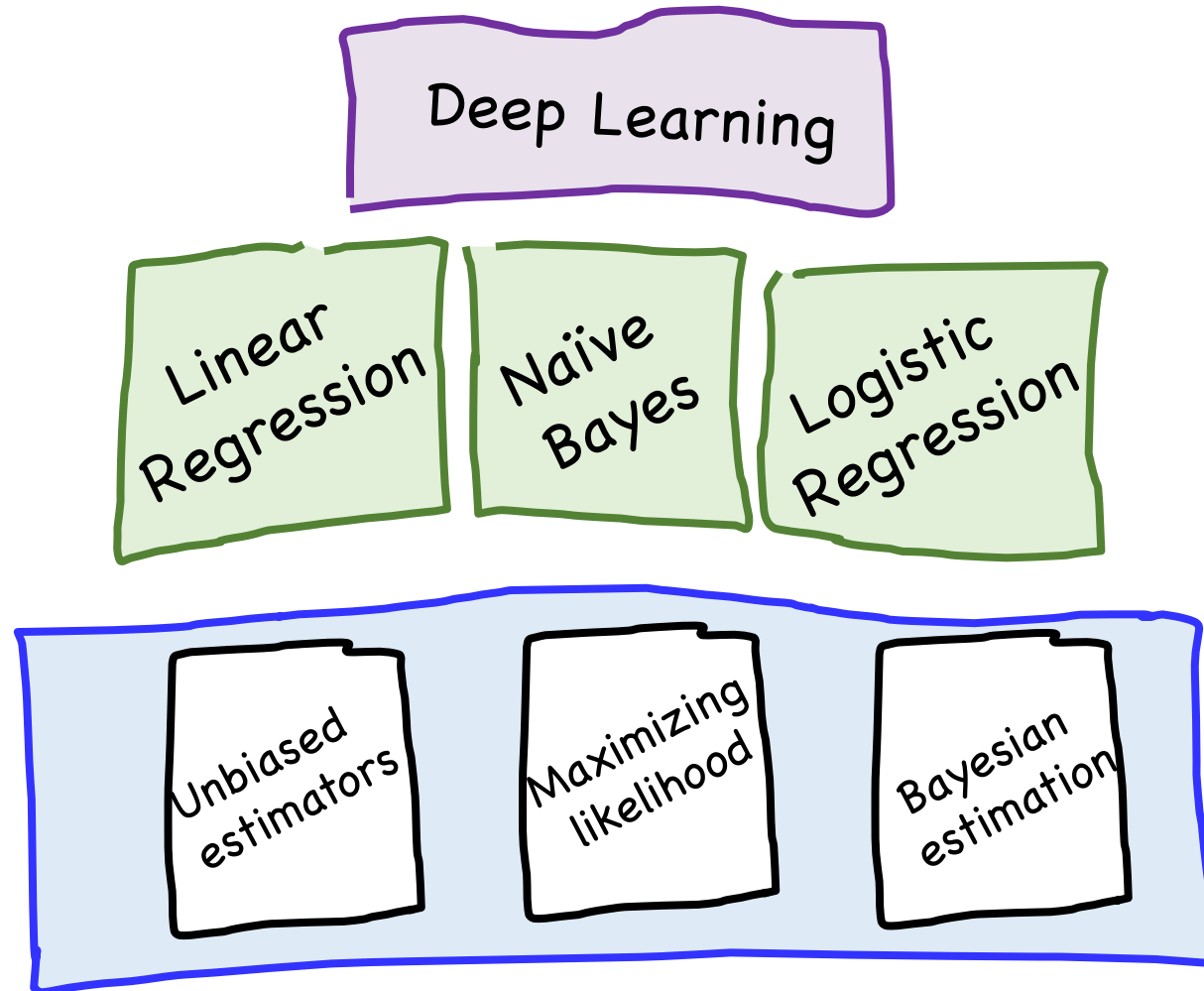
# Where Do The Numbers Come From?



# Our Path



# Our Path



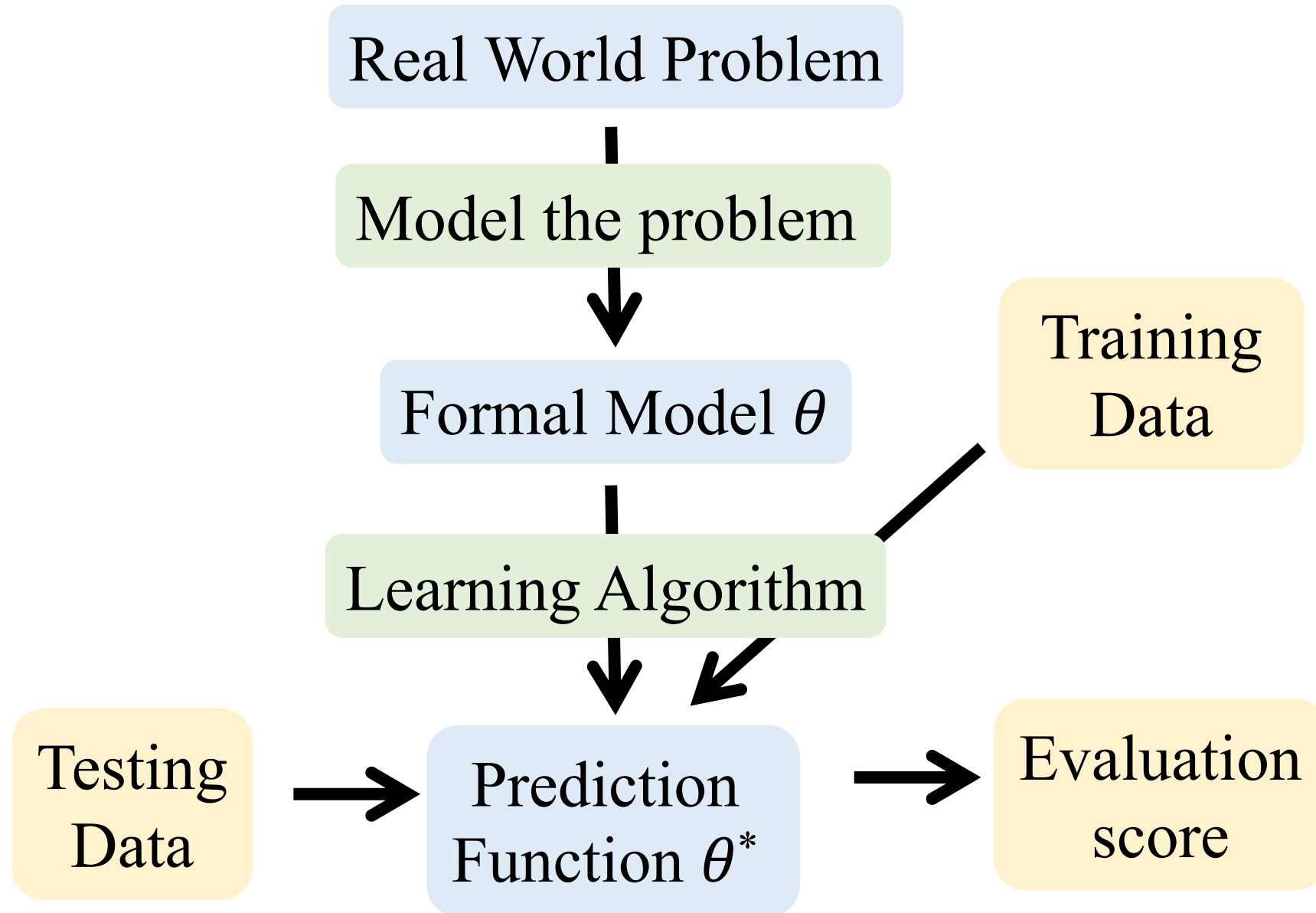
# What are Parameters?

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

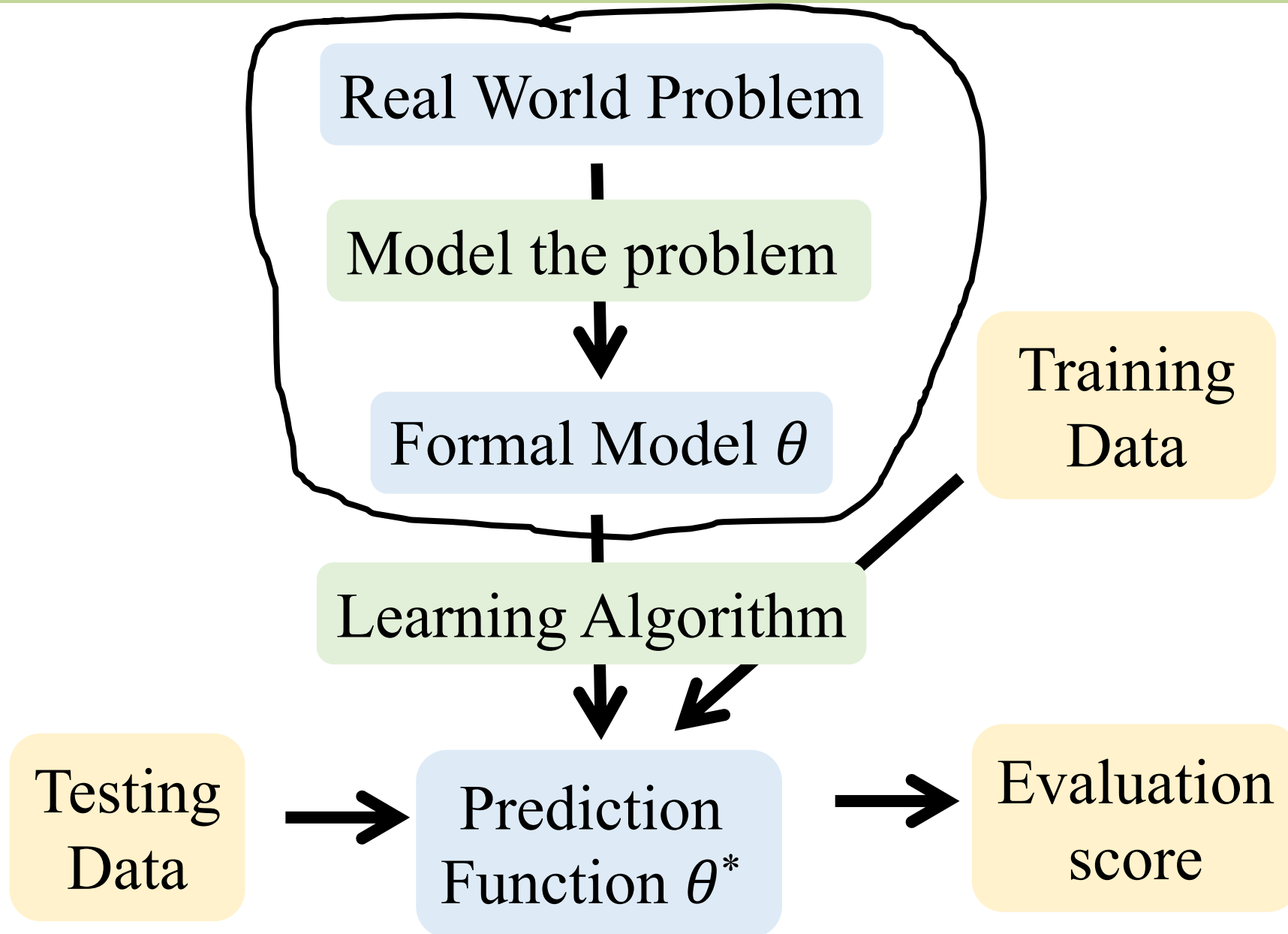
# Why Do We Care?

- 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

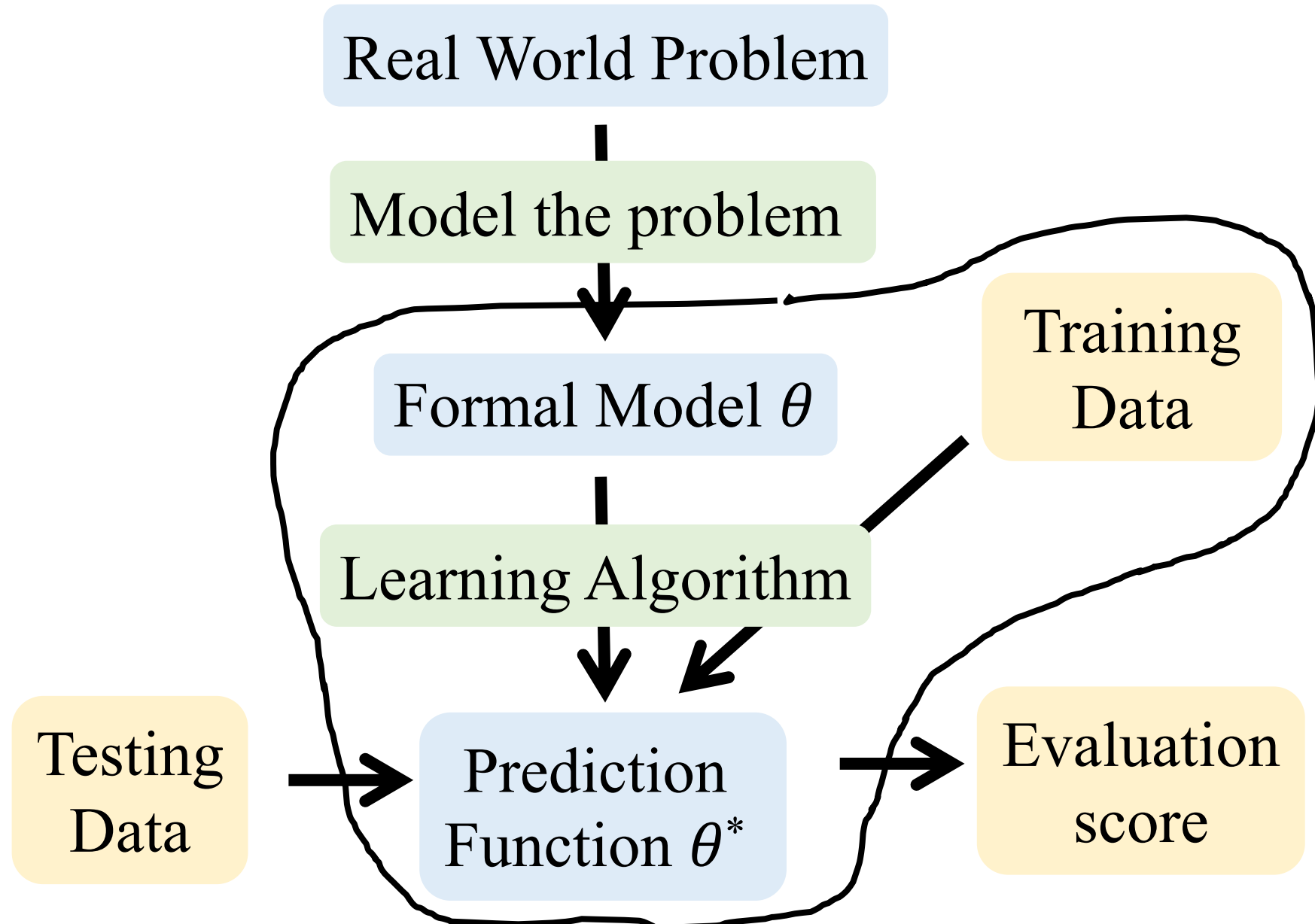
# Supervised Learning



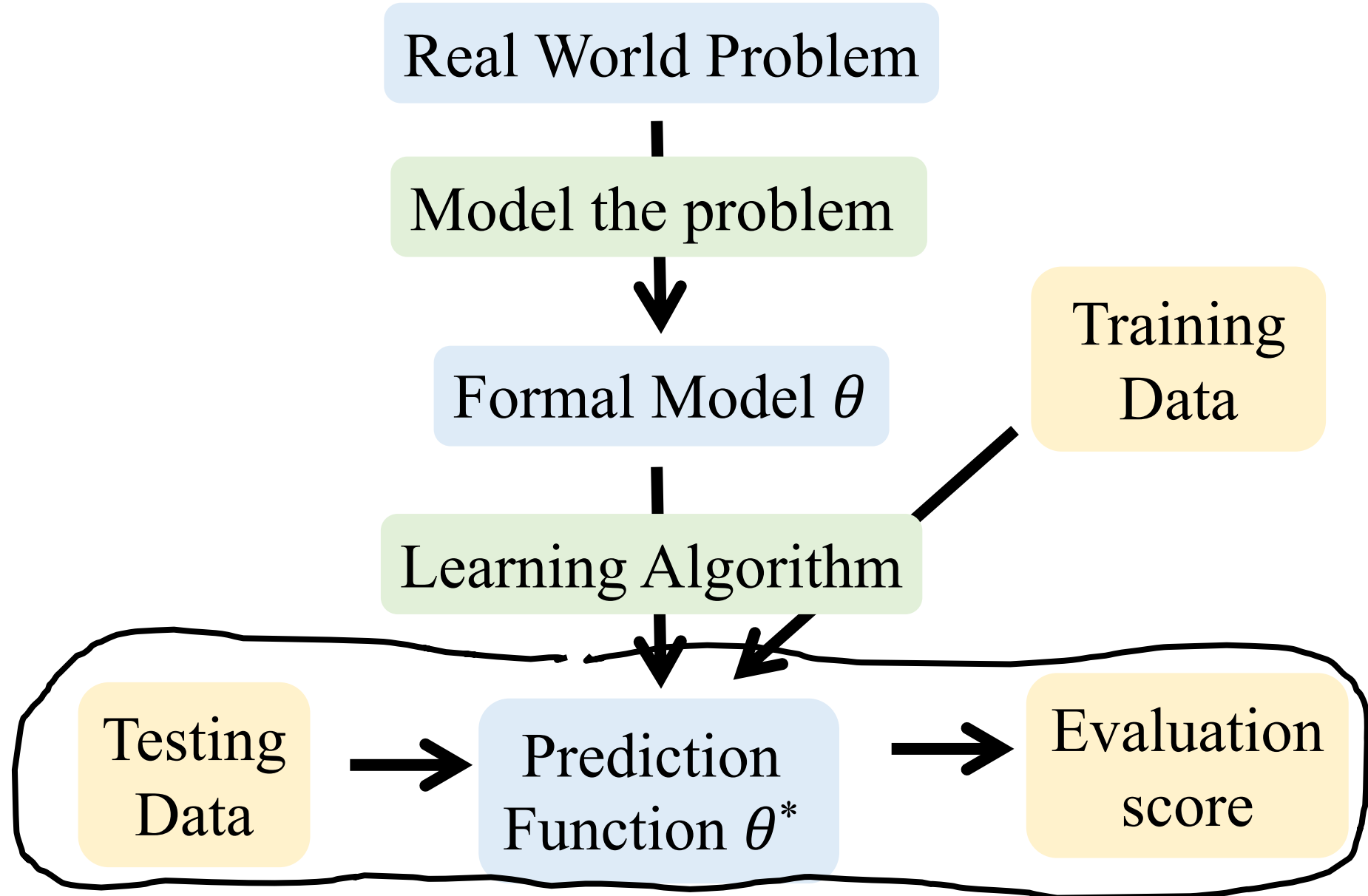
# Modelling



# Training



# Testing



End Review

Basis for learning from data

# Our Path

Neural Networks

Linear  
Regression

Naive  
Bayes

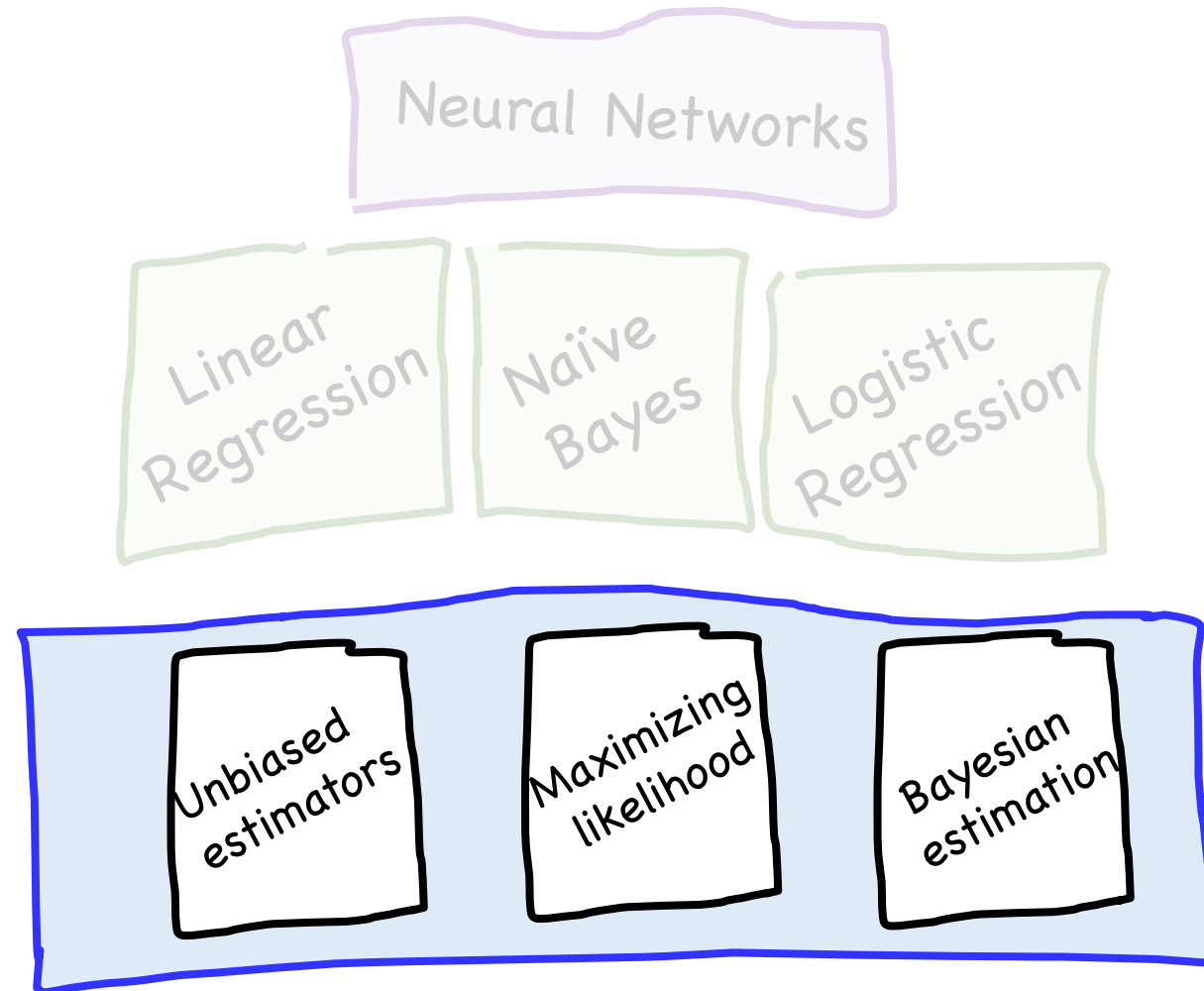
Logistic  
Regression

Unbiased  
estimators

Maximizing  
likelihood

Bayesian  
estimation

# 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 “Mixture of Gaussians”?

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  $\theta$  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  $\theta$

# Maximum Likelihood Estimator

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



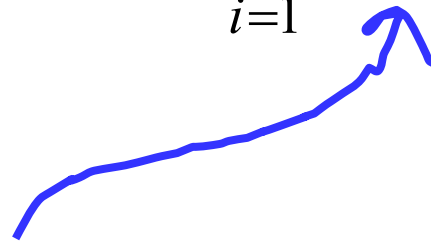
Likelihood (of data given parameters):

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



Likelihood (of data given parameters):

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



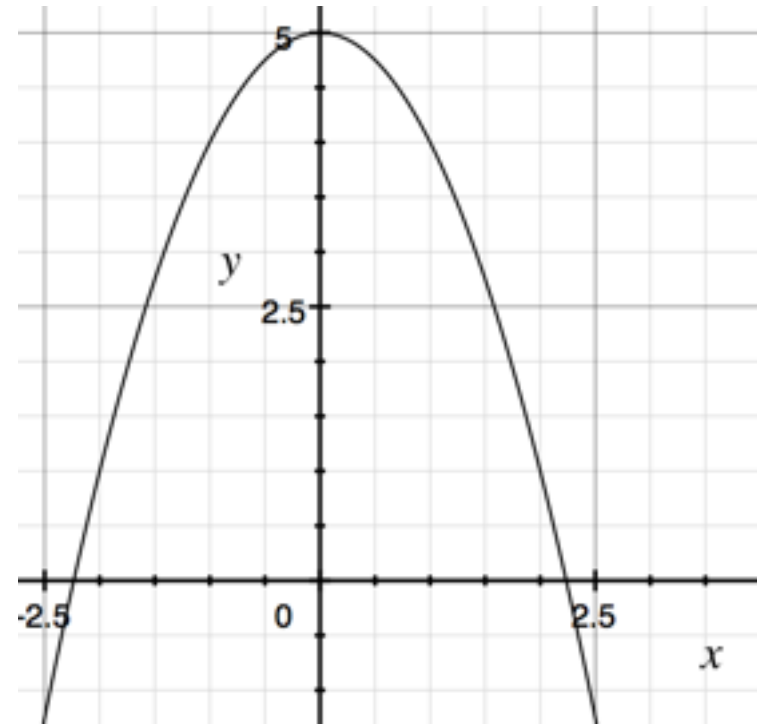
Either the  
PDF (continuous) or  
PMF (discrete)

# Argmax

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

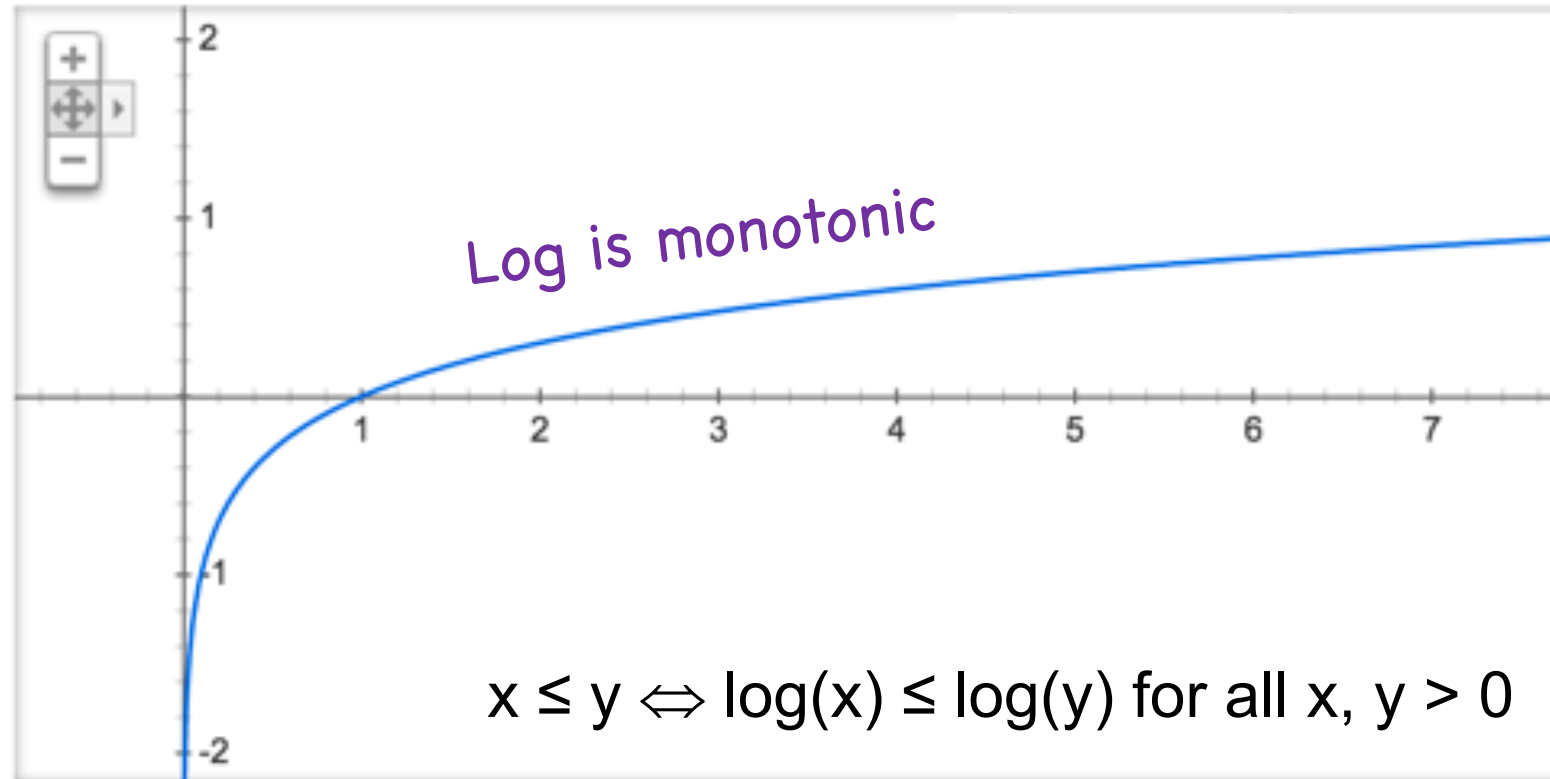
$$\max_x -x^2 + 5 = 5$$

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



# 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



# Maximum Likelihood Estimator

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

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

- $\theta$  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  $\theta$ ):

$$\arg \max_{\theta} (c \cdot LL(\theta)) = \arg \max_{\theta} LL(\theta) = \arg \max_{\theta} L(\theta)$$

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



## Maximum Likelihood

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

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

$$\hat{\theta} = \operatorname{argmax}_{\theta} LL(\theta)$$

But how do we compute  $\operatorname{argmax}$ ?

Option #1: Straight optimization

# Computing the MLE

- General approach for finding MLE of  $\theta$ 
  - Determine formula for  $LL(\theta)$
  - Differentiate  $LL(\theta)$  w.r.t. (each)  $\theta$ :  $\frac{\partial LL(\theta)}{\partial \theta}$
  - To maximize, set  $\frac{\partial LL(\theta)}{\partial \theta} = 0$
  - Solve resulting (simultaneous) equations to get  $\theta_{MLE}$ 
    - 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

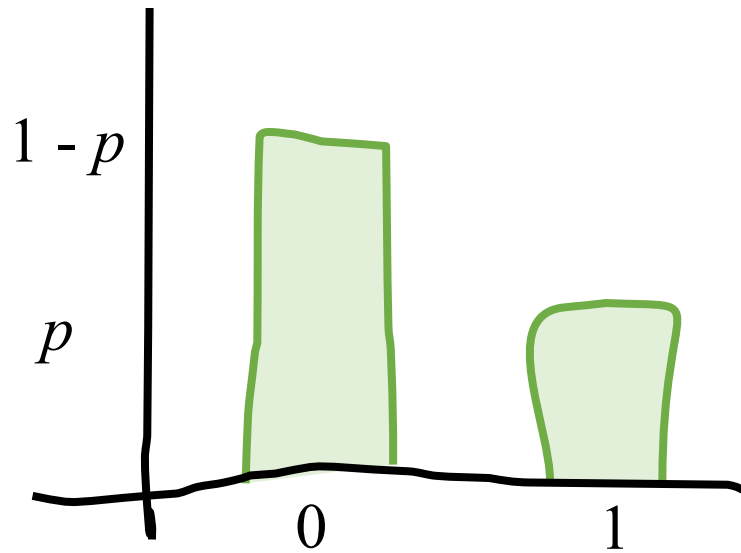
- Consider I.I.D. random variables  $X_1, X_2, \dots, X_n$ 
  - $X_i \sim \text{Ber}(p)$
  - Probability mass function,  $f(X_i | p)$ :



# Maximizing Likelihood with Bernoulli

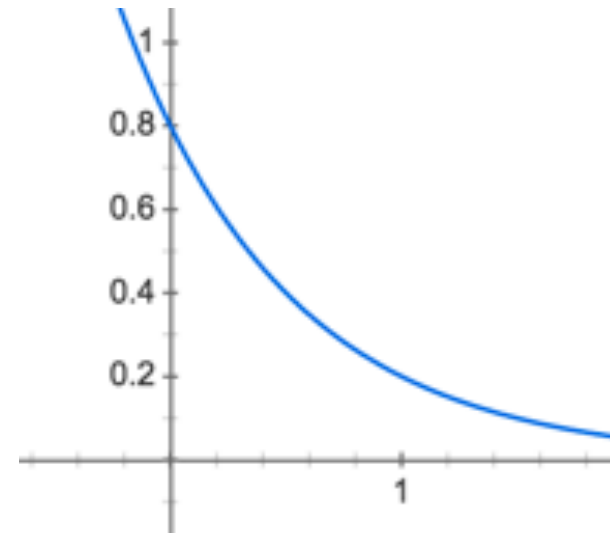
- Consider I.I.D. random variables  $X_1, X_2, \dots, X_n$ 
  - $X_i \sim \text{Ber}(p)$
  - Probability mass function,  $f(X_i | p)$ :

PMF of Bernoulli



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

PMF of Bernoulli ( $p = 0.2$ )



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

# 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.  $\lambda$ , and set to 0:

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

It's 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.  $\mu$ , 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.  $\sigma$ , 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  $\mu_{MLE}$ :

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

- Then, solve for  $\sigma^2_{MLE}$ :

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

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

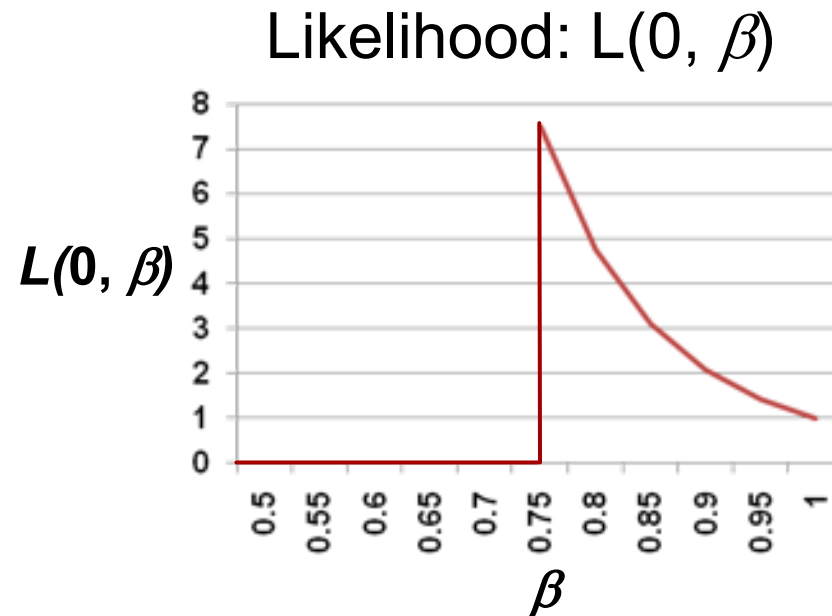
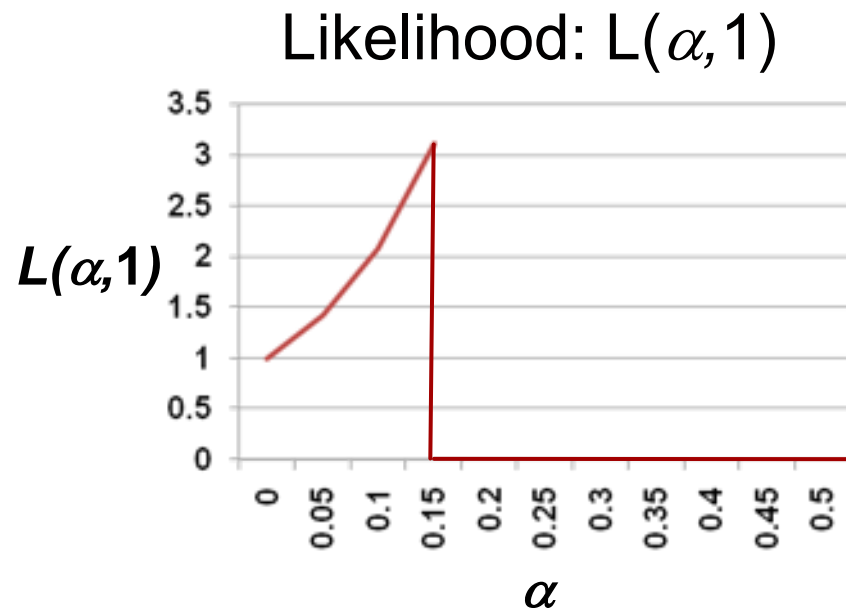
- Note:  $\mu_{MLE}$  unbiased, but  $\sigma^2_{MLE}$  biased

# 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,  $\mu_{MLE} = \frac{1}{n} \sum_{i=1}^n X_i = \text{sample mean}$

- 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

Machine Learning:  
Learn parameters (mostly with MLE) for  
probabilistic models.

Stretch!



SCPD Code:



# Gradient Ascent

Noah Arthurs

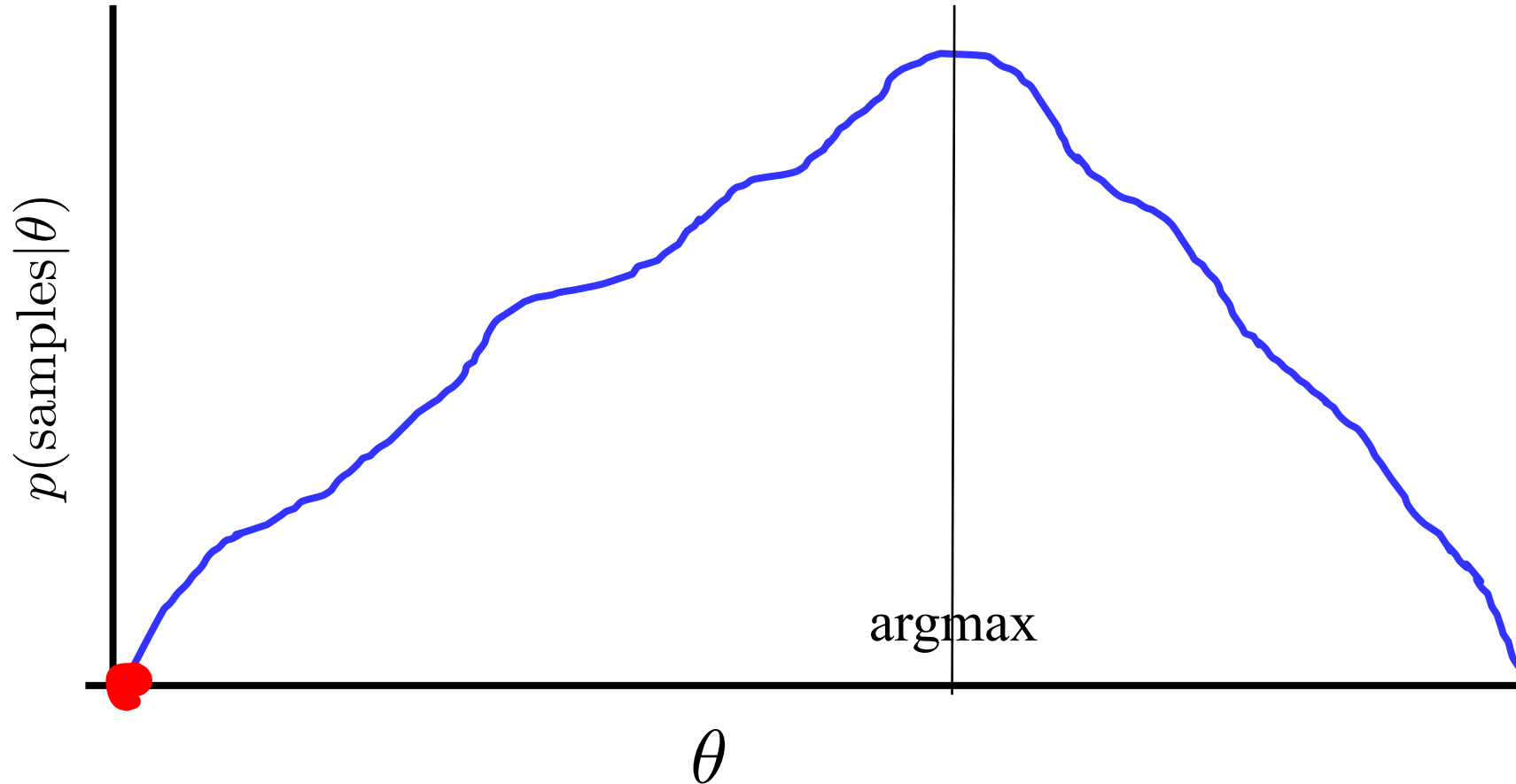
CS109, Stanford University

Argmax

Option #1: Straight optimization

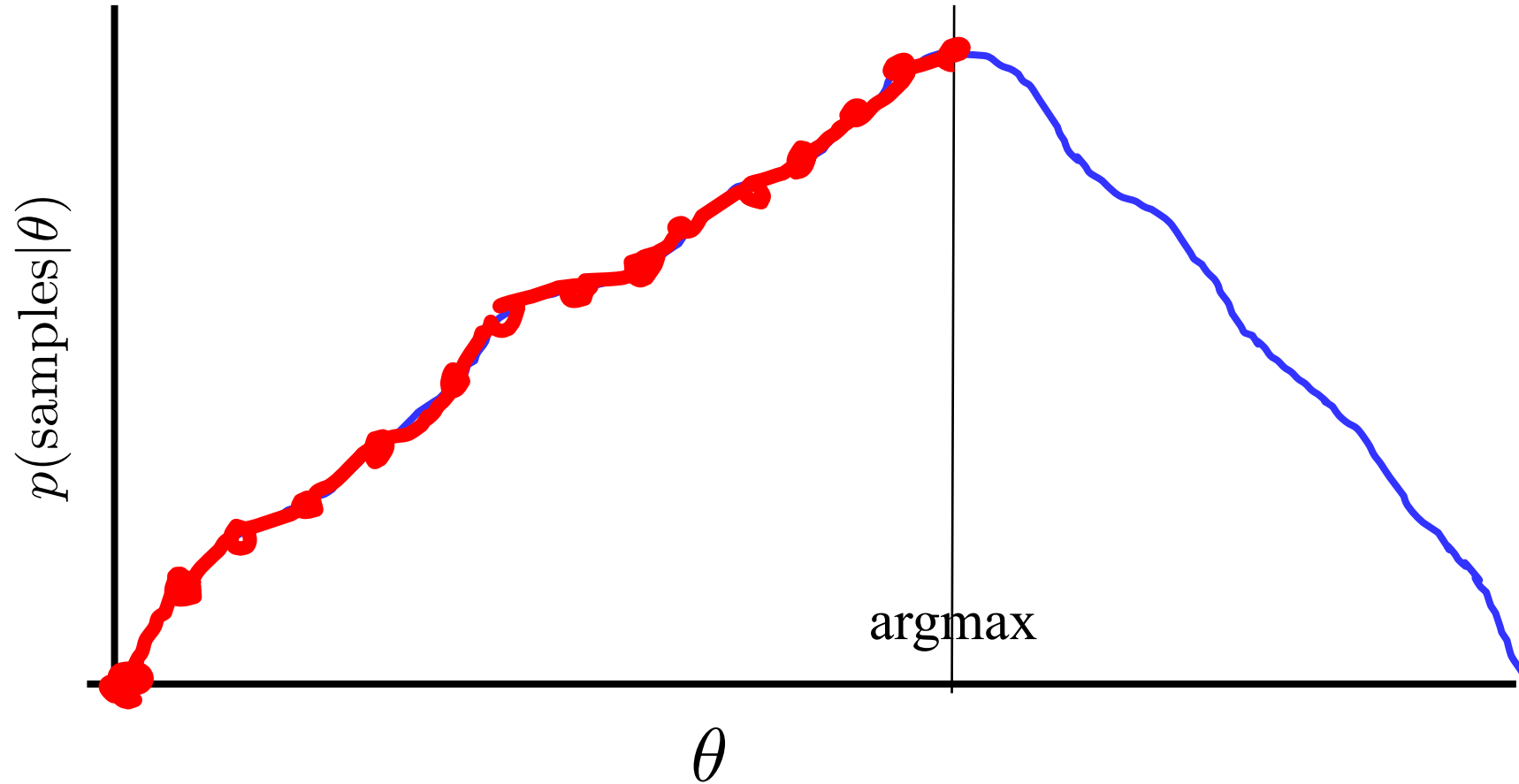
Argmax  
Option #2: Gradient Ascent

# Gradient Ascent



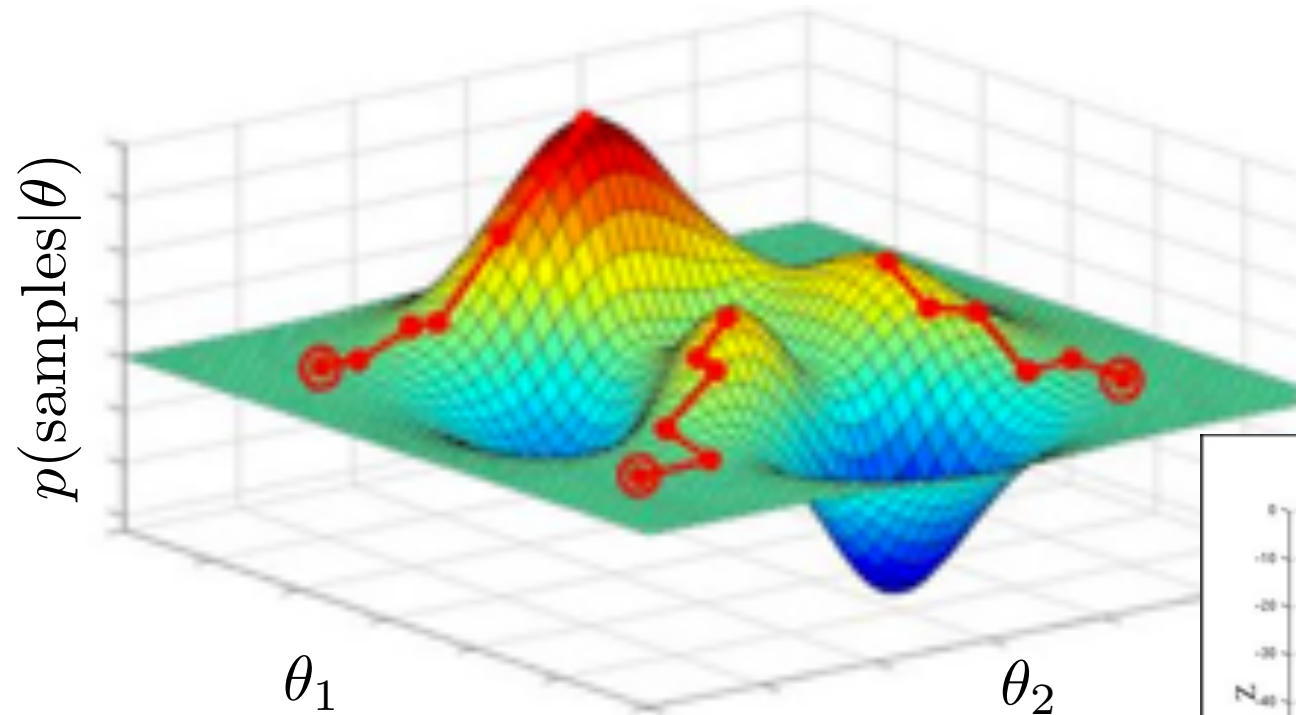
Walk uphill and you will find a local maxima  
(if your step size is small enough)

# Gradient Ascent

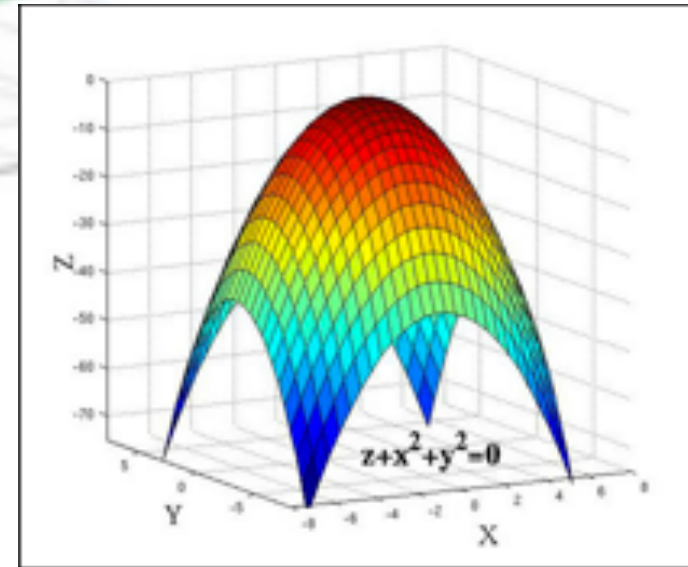


Walk uphill and you will find a local maxima  
(if your step size is small enough)

# Gradient Ascent



Especially good if  
function is convex



Walk uphill and you will find a local maxima  
(if your step size is small enough)

# Gradient Ascent

Repeat many times

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \frac{\partial LL(\theta^{\text{old}})}{\partial \theta_j^{\text{old}}}$$

This is some **profound** life philosophy

Walk uphill and you will find a local maxima  
(if your step size is small enough)



Gradient ascent is your  
bread and butter  
algorithm for optimization  
(eg argmax)