



# Maximum A Posteriori

Noah Arthurs

CS109, Stanford University

# Problem Set 6

**Naïve Bayes:** Need today's lecture

**Logistic Regression:** Need Wednesday

**Neural Net (EC):** Need Friday

**Mini-Project (EC):** Whatever you'd like! Start early!

Review

# Review: Maximum Likelihood Algorithm

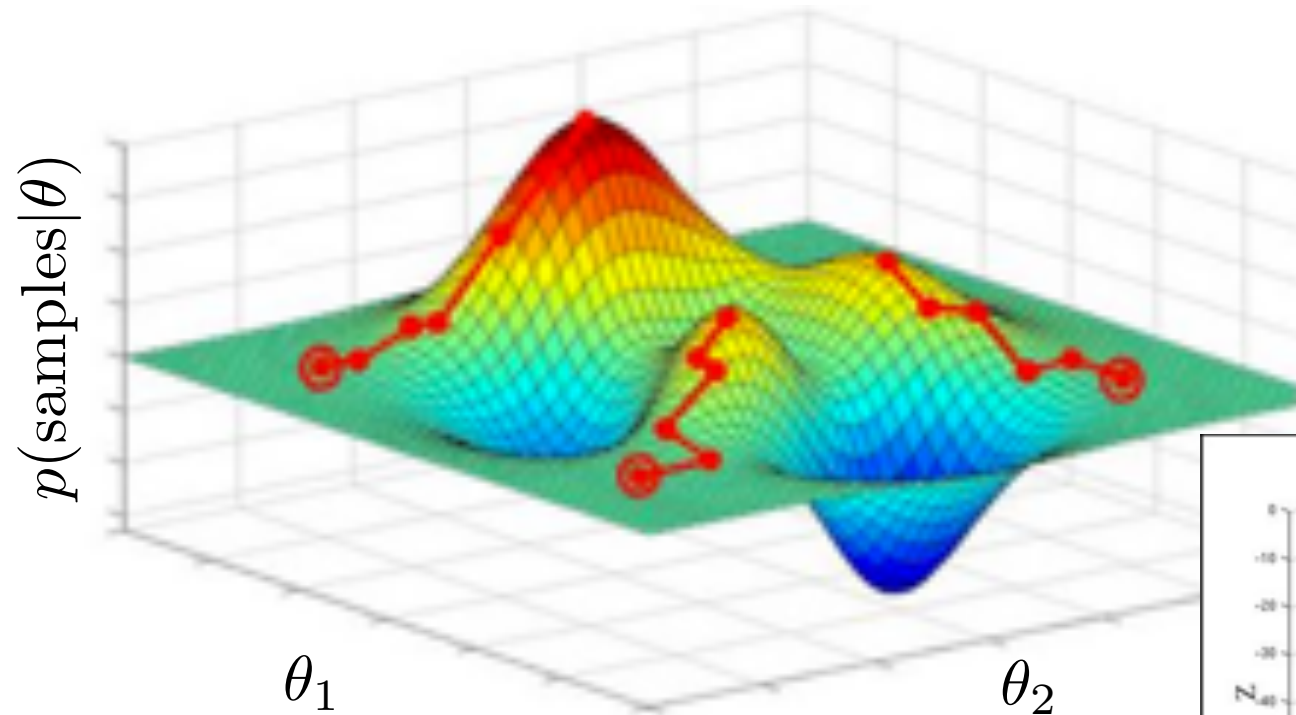
1. Decide on a model for the likelihood of your samples. This is often using a PMF or PDF.

2. Write out the log likelihood function.

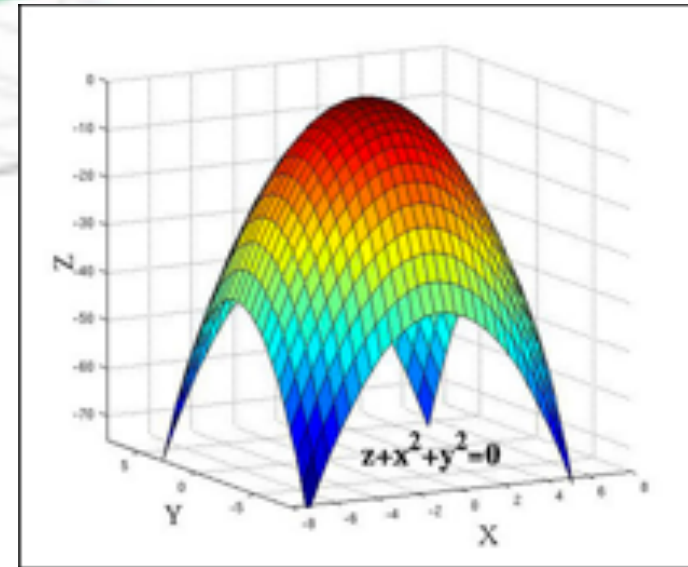
3. State that the optimal parameters are the argmax of the log likelihood function.

4. Use an optimization algorithm to calculate argmax

# 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

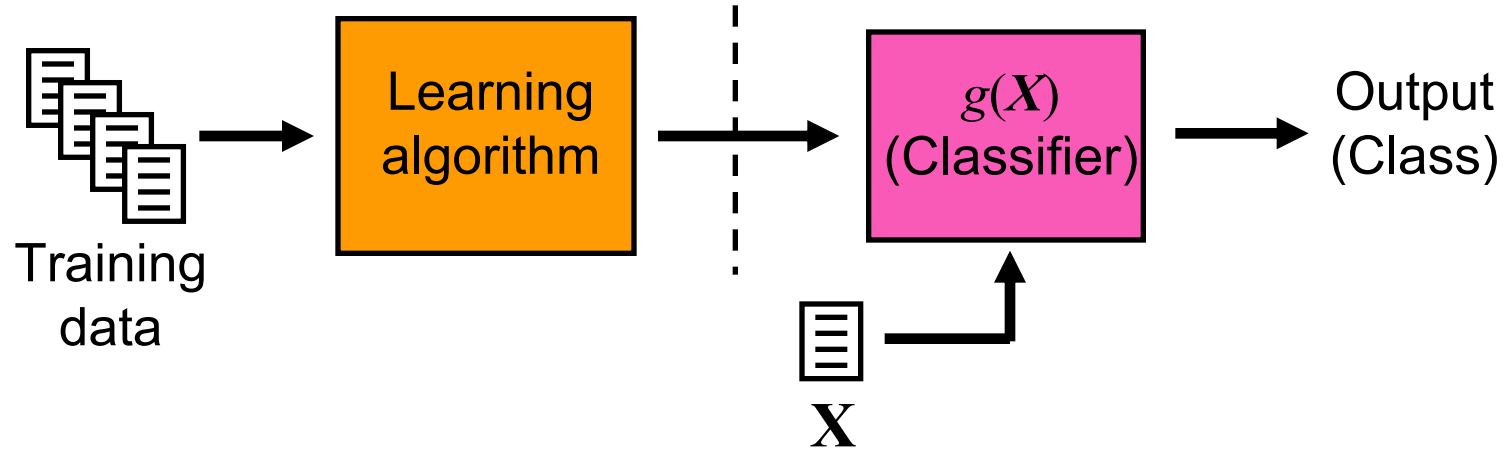
Repeat many times

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

This is some **profound** life philosophy

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

# The Machine Learning Process

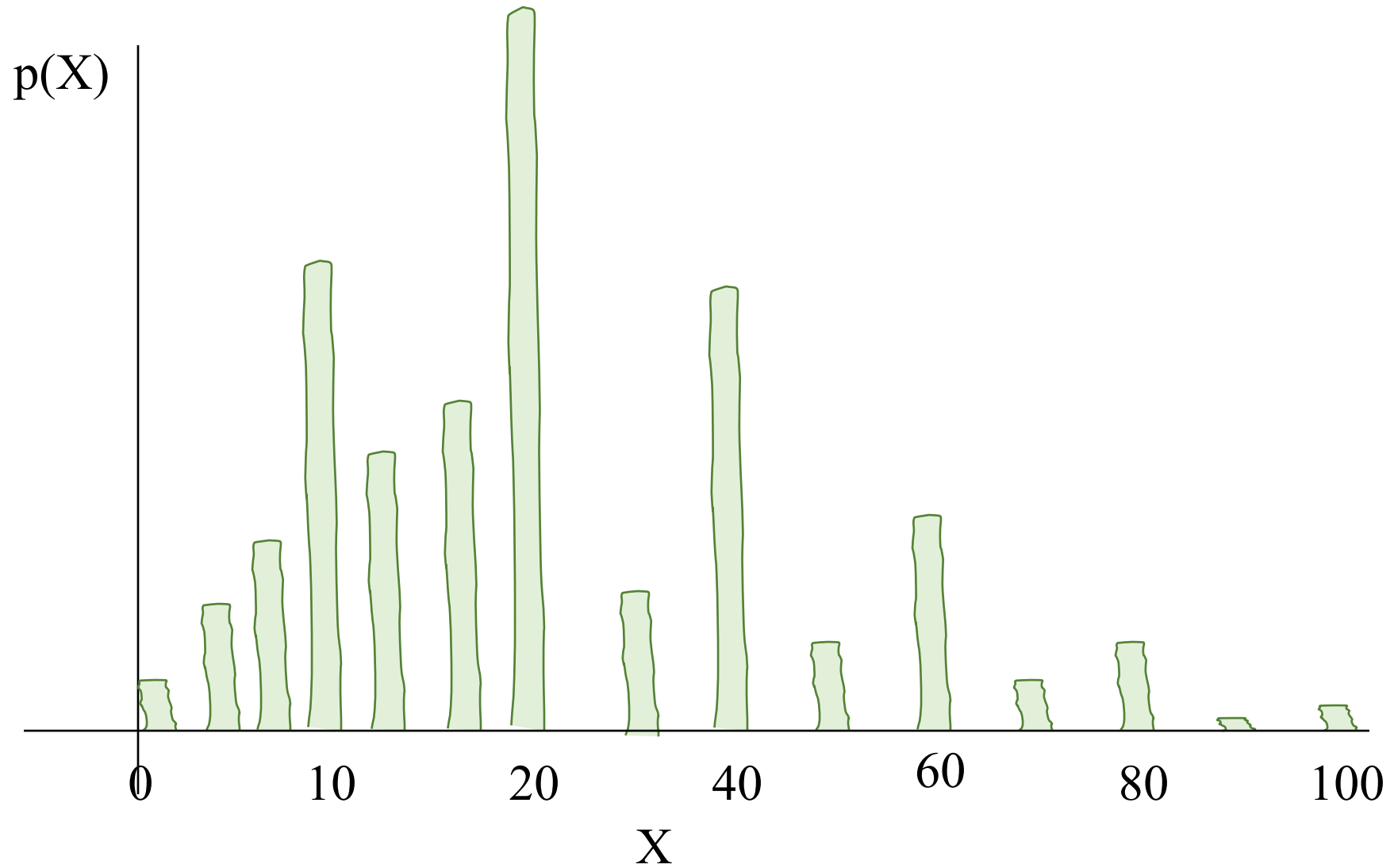


- Training data: set of  $N$  pre-classified data instances
  - $N$  training pairs:  $(\mathbf{x}^{(1)}, y^{(1)})$ ,  $(\mathbf{x}^{(2)}, y^{(2)})$ , ...,  $(\mathbf{x}^{(n)}, y^{(n)})$ 
    - Use superscripts to denote  $i$ -th training instance
- Learning algorithm: method for determining  $g(X)$ 
  - Given a new input observation of  $\mathbf{x} = x_1, x_2, \dots, x_m$
  - Use  $g(\mathbf{x})$  to compute a corresponding output (prediction)

# Two Envelopes

- I have two envelopes, will allow you to have one
  - One contains  $\$X$ , the other contains  $\$2X$
  - Select an envelope
    - Open it!
  - Now, would you like to switch for other envelope?
  - To help you decide, compute  $E[\$ \text{ in other envelope}]$ 
    - Let  $Y = \$$  in envelope you selected
$$E[\$ \text{ in other envelope}] = \frac{1}{2} \cdot \frac{Y}{2} + \frac{1}{2} \cdot 2Y = \frac{5}{4} Y$$
  - Before opening envelope, think either equally good
  - So, what happened by opening envelope?
    - And does it really make sense to switch?

# Subjectivity of Probability



# Maximum A Posteriori

**data:**  $x^{(1)}, \dots, x^{(n)}$        $\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} f(\theta | x^{(1)}, \dots, x^{(n)})$

---

$$\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} \frac{g(\theta) f(x^{(1)}, x^{(2)}, \dots, x^{(n)} | \theta)}{h(x^{(1)}, x^{(2)}, \dots, x^{(n)})}$$

$$= \operatorname{argmax}_{\theta} \frac{g(\theta) \prod_{i=1}^n f(x^{(i)} | \theta)}{h(x^{(1)}, x^{(2)}, \dots, x^{(n)})}$$

$$= \operatorname{argmax}_{\theta} g(\theta) \prod_{i=1}^n f(x^{(i)} | \theta)$$

$$= \operatorname{argmax}_{\theta} \left( \log(g(\theta)) + \sum_{i=1}^n \log(f(x^{(i)} | \theta)) \right)$$



# Maximum A Posteriori



Estimated  
parameter



Log prior



$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} \left( \log(g(\theta)) + \sum_{i=1}^n \log(f(x^{(i)} | \theta)) \right)$$



Chose the value of theta  
that maximizes:

Sum of  
log likelihood



# MLE vs MAP

**Data:**  $x^{(1)}, \dots, x^{(n)}$

## Maximum Likelihood Estimation

$$\begin{aligned}\hat{\theta}_{MLE} &= \operatorname{argmax}_{\theta} f(x^{(1)}, \dots, x^{(n)} | \theta) \\ &= \operatorname{argmax}_{\theta} \left( \sum_i \log f(x^{(i)} | \theta) \right)\end{aligned}$$

## Maximum A Posteriori

$$\begin{aligned}\hat{\theta}_{MAP} &= \operatorname{argmax}_{\theta} f(\theta | x^{(1)}, \dots, x^{(n)}) \\ &= \operatorname{argmax}_{\theta} \left( \log(g(\theta)) + \sum_{i=1}^n \log(f(x^{(i)} | \theta)) \right)\end{aligned}$$

# $P(\theta | D)$ For Bernoulli

- Prior:  $\theta \sim \text{Beta}(a, b)$ ; data =  $\{n \text{ heads}, m \text{ tails}\}$
- Estimate  $p$ , aka  $\theta$

$$\begin{aligned}\hat{\theta}_{MAP} &= \underset{\theta}{\operatorname{argmax}} f(\theta | \text{data}) &&= \underset{\theta}{\operatorname{argmax}} f(\text{data} | \theta) g(\theta) \\ & && \text{This is the beta PDF} \\ &= \underset{\theta}{\operatorname{argmax}} \log g(\theta) + \log f(\text{data} | \theta) && \text{Product of thetas and (1-theta)s} \\ &= \underset{\theta}{\operatorname{argmax}} \log \left[ \frac{1}{\beta} \theta^{a-1} (1 - \theta)^{b-1} \right] \\ & \quad + n \log f(\text{heads} | \theta) \\ & \quad + m \log f(\text{tails} | \theta) \\ &= \underset{\theta}{\operatorname{argmax}} \log \frac{1}{\beta} + (a - 1) \log \theta + (b - 1) \log(1 - \theta) + n \log \theta + m \log(1 - \theta) \\ &= \underset{\theta}{\operatorname{argmax}} (a - 1 + n) \log \theta + (b - 1 + m) \log(1 - \theta)\end{aligned}$$

# $P(\theta | D)$ For Bernoulli

- Prior:  $\theta \sim \text{Beta}(a, b)$ ;  $D = \{n \text{ heads}, m \text{ tails}\}$
- Estimate  $p$ , aka  $\theta$

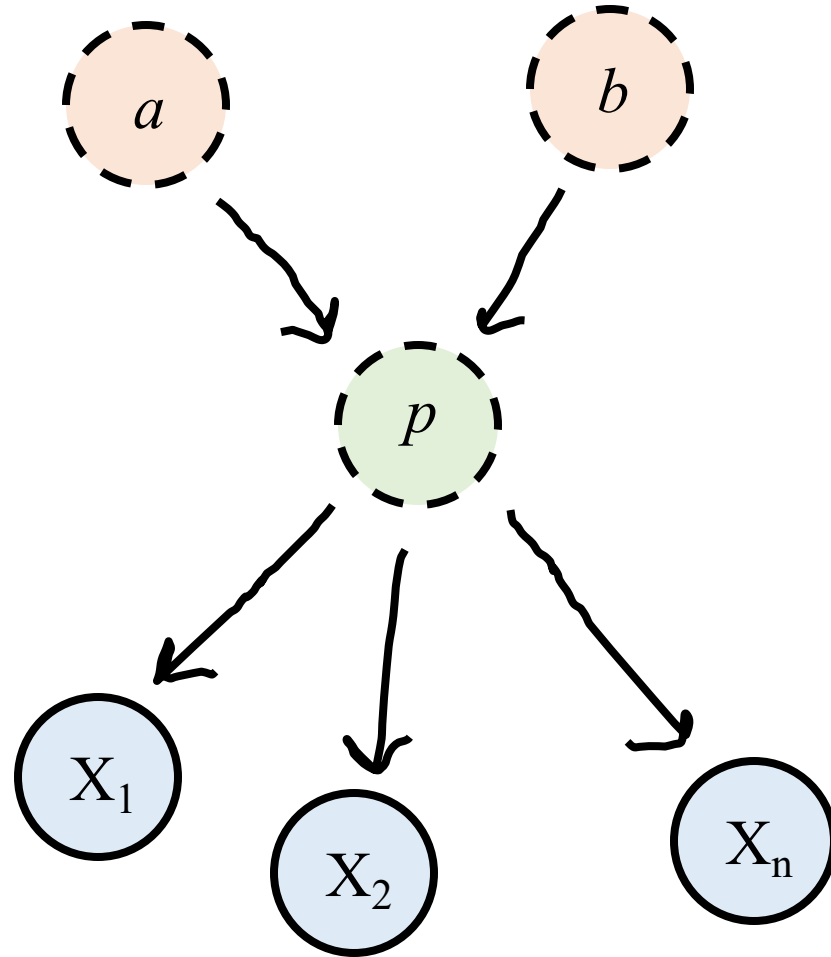
$$\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} f(\theta | \text{data})$$

$$= \operatorname{argmax}_{\theta} (a - 1 + n) \log \theta + (b - 1 + m) \log(1 - \theta)$$

$$= \frac{n + a - 1}{n + m + a + b - 2}$$

That's the mode of the updated beta

# Hyper Parameters



Hyperparameter  
 $a, b$  are fixed

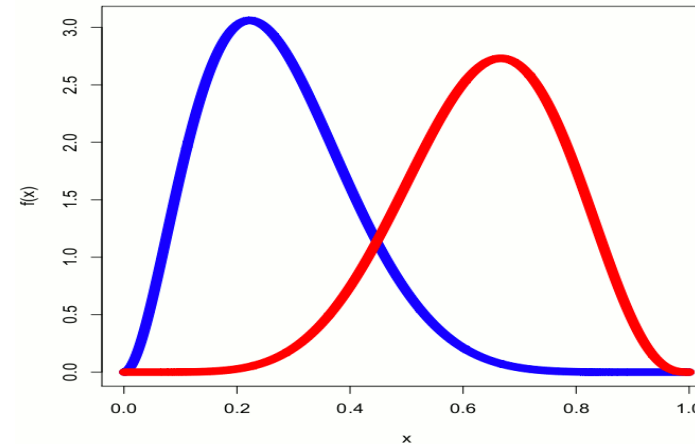
Prior  
 $p \sim \text{Beta}(a, b)$

Data distribution  
 $X_i \sim \text{Bern}(p)$

MAP will estimate the most likely value of  $p$  for this model

# Where'd Ya Get Them $P(\theta)$ ?

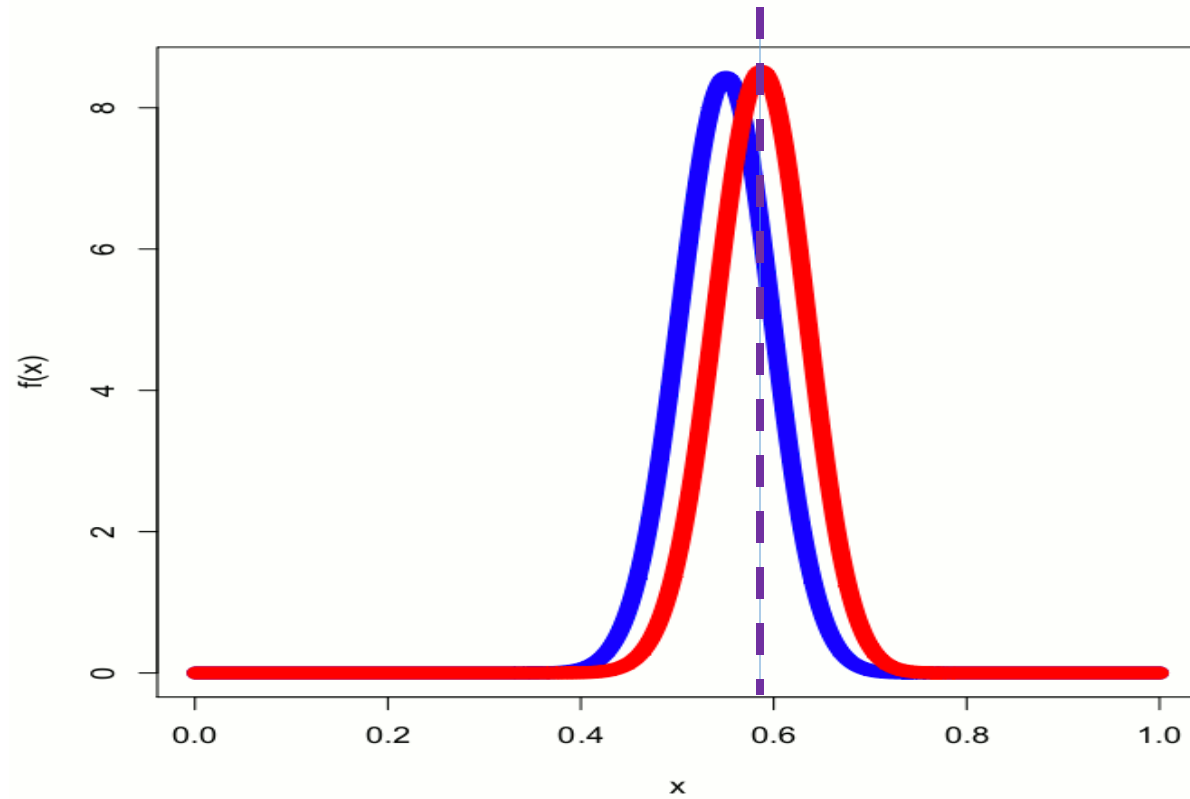
- $\theta$  is the probability a coin turns up heads
- Model  $\theta$  with 2 different priors:
  - $P_1(\theta)$  is Beta(3,8) (blue)
  - $P_2(\theta)$  is Beta(7,4) (red)
- They look pretty different!



- Now flip 100 coins; get 58 heads and 42 tails
  - What do posteriors look like?

# It's Like Having Twins

argmax returns the mode



- As long as we collect enough data, posteriors will converge to the true value!

# Conjugate Distributions Without Tears

- Just for review...
- Have coin with unknown probability  $\theta$  of heads
  - Our prior (subjective) belief is that  $\theta \sim \text{Beta}(a, b)$
  - Now flip coin  $k = n + m$  times, getting  $n$  heads,  $m$  tails
  - Posterior density:  $(\theta \mid n \text{ heads}, m \text{ tails}) \sim \text{Beta}(a+n, b+m)$ 
    - Beta is conjugate for Bernoulli, Binomial, Geometric, and Negative Binomial
  - $a$  and  $b$  are called “hyperparameters”
    - Saw  $(a + b - 2)$  imaginary trials, of those  $(a - 1)$  are “successes”
  - For a coin you never flipped before, use  $\text{Beta}(x, x)$  to denote you think coin likely to be fair
    - How strongly you feel coin is fair is a function of  $x$

End Review

# Gonna Need Priors

Parameter

Distribution for Parameter

Bernoulli  $p$

Beta

Binomial  $p$

Beta

Poisson  $\lambda$

Gamma

Exponential  $\lambda$

Gamma

Multinomial  $p_i$

Dirichlet

Normal  $\mu$

Normal

Normal  $\sigma^2$

Inverse Gamma

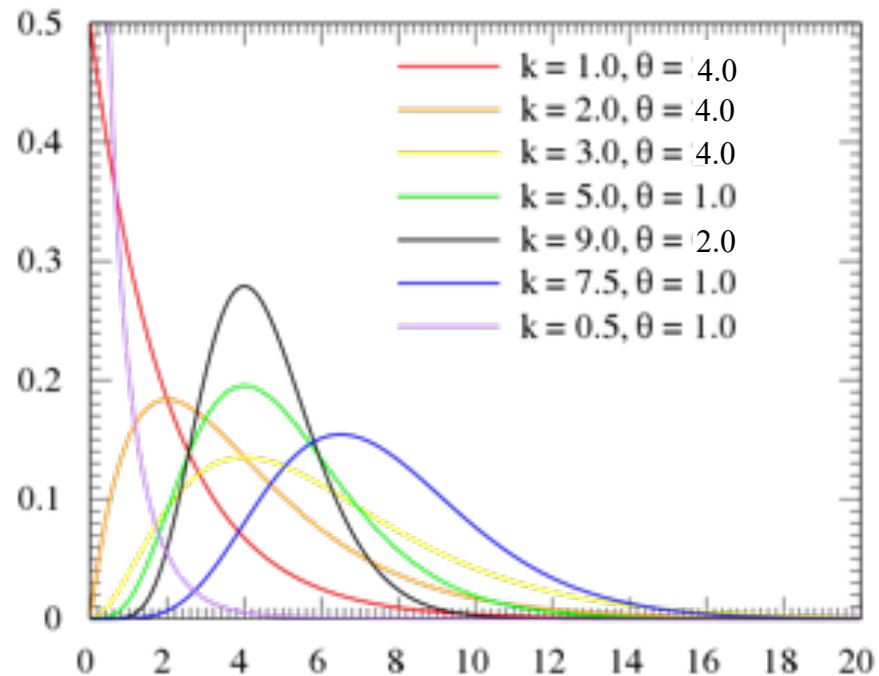
Know

Familiarity

Not necessary but  
good to know of

# Good Times with Gamma

- Gamma( $k, \theta$ ) distribution
  - Conjugate for Poisson Rate
    - Also conjugate for Exponential, but we won't delve into that
  - Intuitive understanding of hyperparameters:
    - Saw  $k$  total imaginary events during  $\theta$  prior time periods



# Good Times with Gamma

- Gamma( $k, \theta$ ) distribution
  - Conjugate for Poisson Rate
    - Also conjugate for Exponential, but we won't delve into that
  - Intuitive understanding of hyperparameters:
    - Saw  $k$  total imaginary events during  $\theta$  prior time periods
  - Updating with observations
    - After observing  $n$  events during next  $t$  time periods...
    - ... posterior distribution is Gamma( $k + n, \theta + t$ )
    - ...MAP estimator for Poisson with Gamma prior is  $(k+n)/(\theta + t)$
    - Example: Prior for rate is Gamma(10, 5)
    - Saw 10 events in 5 time periods. Like observing at rate = 2
    - Now see 11 events in next 2 time periods  $\rightarrow$  Gamma(21, 7)
    - MAP rate = 3

# Reviving an Old Story Line

The Multinomial Distribution  $\text{Mult}(p_1, \dots, p_k)$

$$p(x_1, \dots, x_k) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}$$

# Multinomial is Multiple Times the Fun

- Dirichlet( $a_1, a_2, \dots, a_m$ ) distribution
  - Conjugate for Multinomial
    - Dirichlet generalizes Beta in same way Multinomial generalizes Bernoulli

$$f(X_1 = x_1, X_2 = x_2, \dots, X_m = x_m) = K \prod_{i=1}^m x_i^{a_i - 1}$$

- Intuitive understanding of hyperparameters:
  - Saw  $\sum_{i=1}^m a_i - m$  imaginary trials, with  $(a_i - 1)$  of outcome  $i$
- Updating to get the posterior distribution
  - After observing  $n_1 + n_2 + \dots + n_m$ , new trials with  $n_i$  of outcome  $i$ ...
  - ... posterior distribution is Dirichlet( $a_1 + n_1, a_2 + n_2, \dots, a_m + n_m$ )

# Example: Estimating Die Parameters



# Your Happy Laplace

- Recall example of 6-sides die rolls:
  - $X \sim \text{Multinomial}(p_1, p_2, p_3, p_4, p_5, p_6)$
  - Roll  $n = 12$  times
  - Result: 3 ones, 2 twos, 0 threes, 3 fours, 1 fives, 3 sixes
    - MLE:  $p_1=3/12$ ,  $p_2=2/12$ ,  $p_3=0/12$ ,  $p_4=3/12$ ,  $p_5=1/12$ ,  $p_6=3/12$
  - Dirichlet prior allows us to pretend we saw each outcome  $k$  times before. MAP estimate:  $p_i = \frac{X_i + k}{n + mk}$ 
    - Laplace's "law of succession": idea above with  $k = 1$
    - Laplace estimate:  $p_i = \frac{X_i + 1}{n + m}$
    - Laplace:  $p_1=4/18$ ,  $p_2=3/18$ ,  $p_3=1/18$ ,  $p_4=4/18$ ,  $p_5=2/18$ ,  $p_6=4/18$
    - No longer have 0 probability of rolling a three!

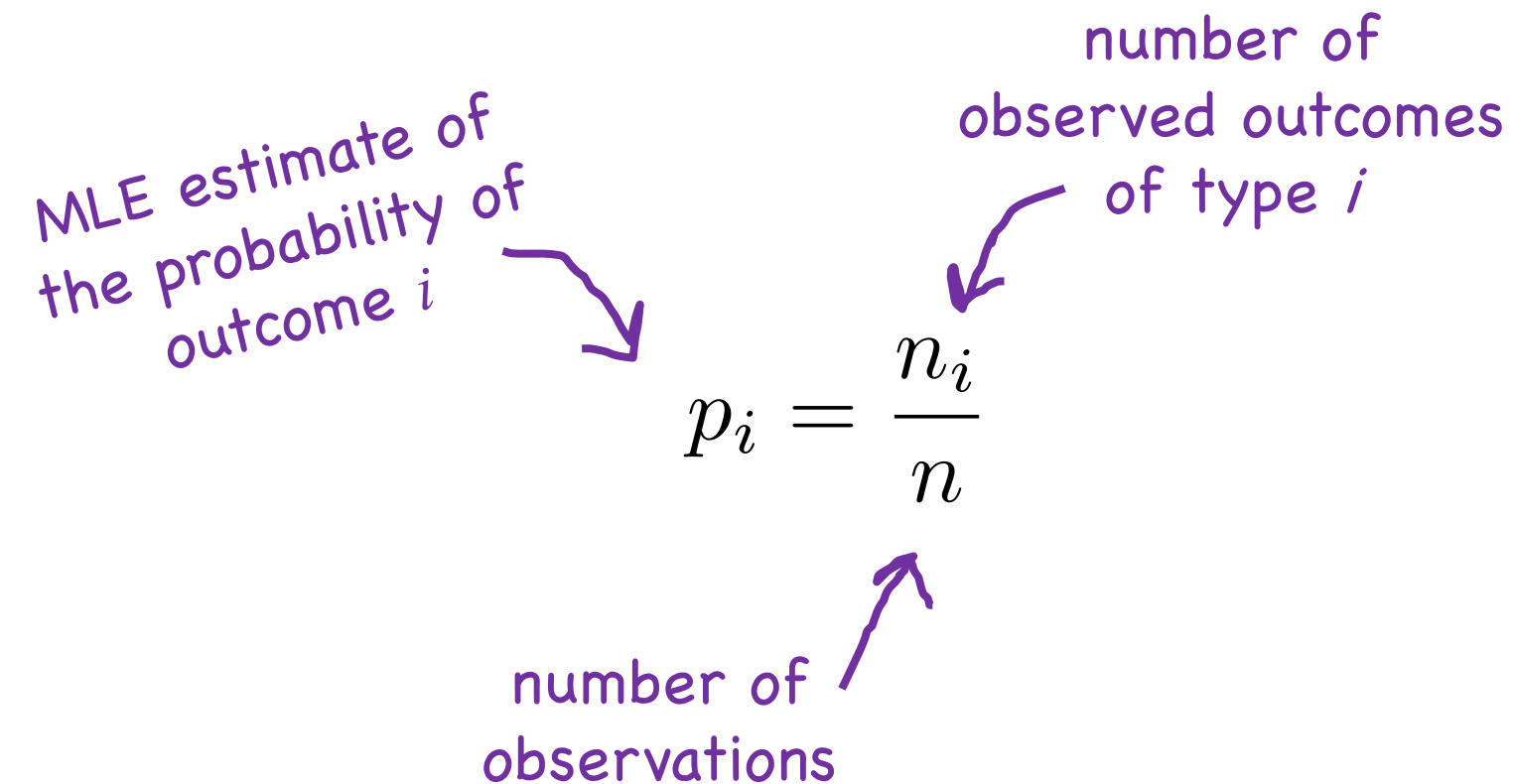
# MLE for Multinomial

MLE estimate of the probability of outcome  $i$

number of observed outcomes of type  $i$

$$p_i = \frac{n_i}{n}$$

number of observations



$\theta$  is  $p$ .  
For a multinomial

# MAP for Multinomial, Laplace Prior

MAP estimate of the probability of outcome  $i$

number of observed outcomes of type  $i$

$$p_i = \frac{n_i + 1}{n + m}$$

number of observations

number of outcome types

$\theta$  is  $p$ .  
For a multinomial

Stretch!



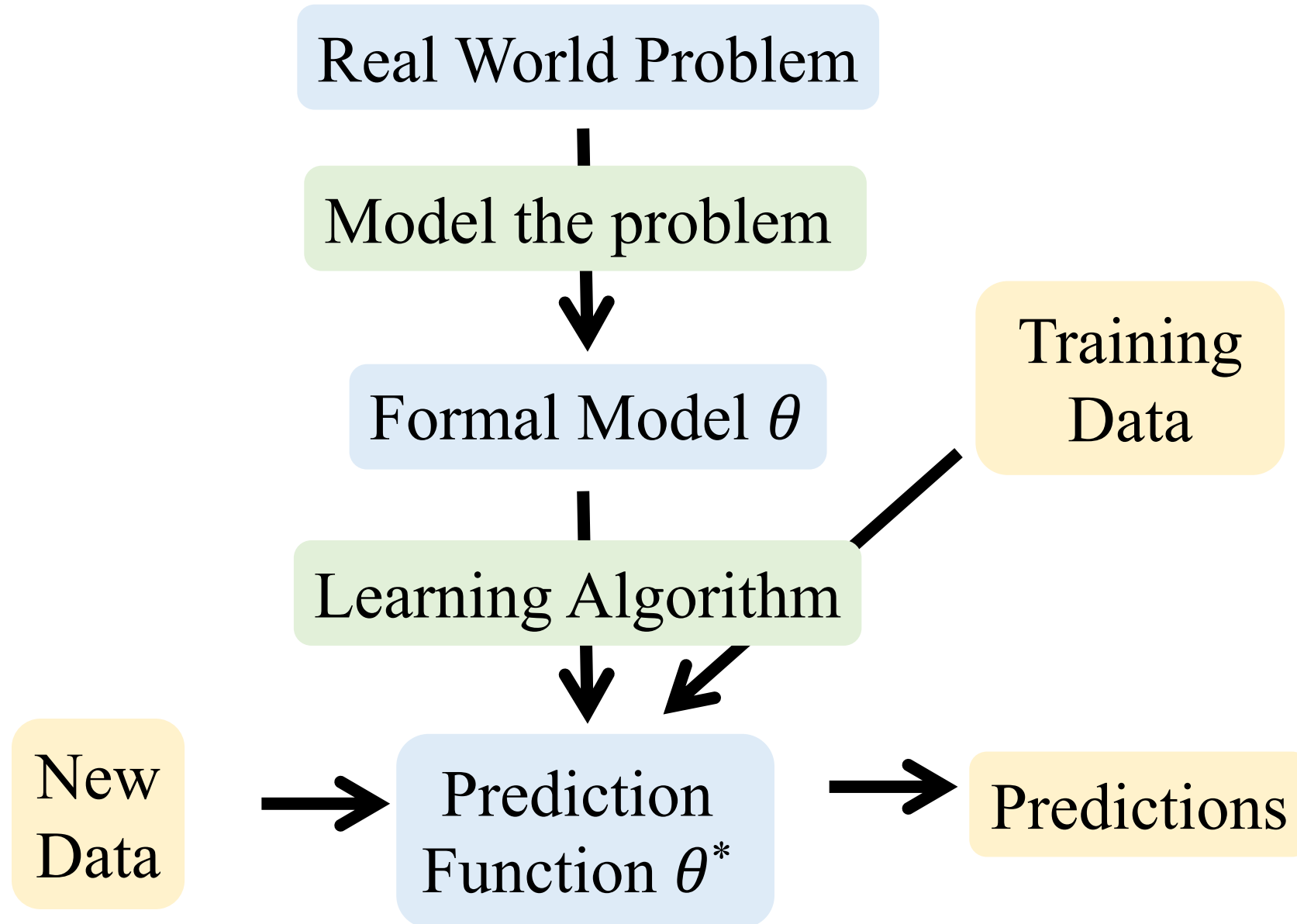
SCPD Code:  
Purple Car

A woman with a gold coin headband and a dark, patterned top is looking intently into a glowing white crystal ball. Her hands are positioned around the ball, as if holding it. The background is dark, and the lighting is focused on the crystal ball and her face.

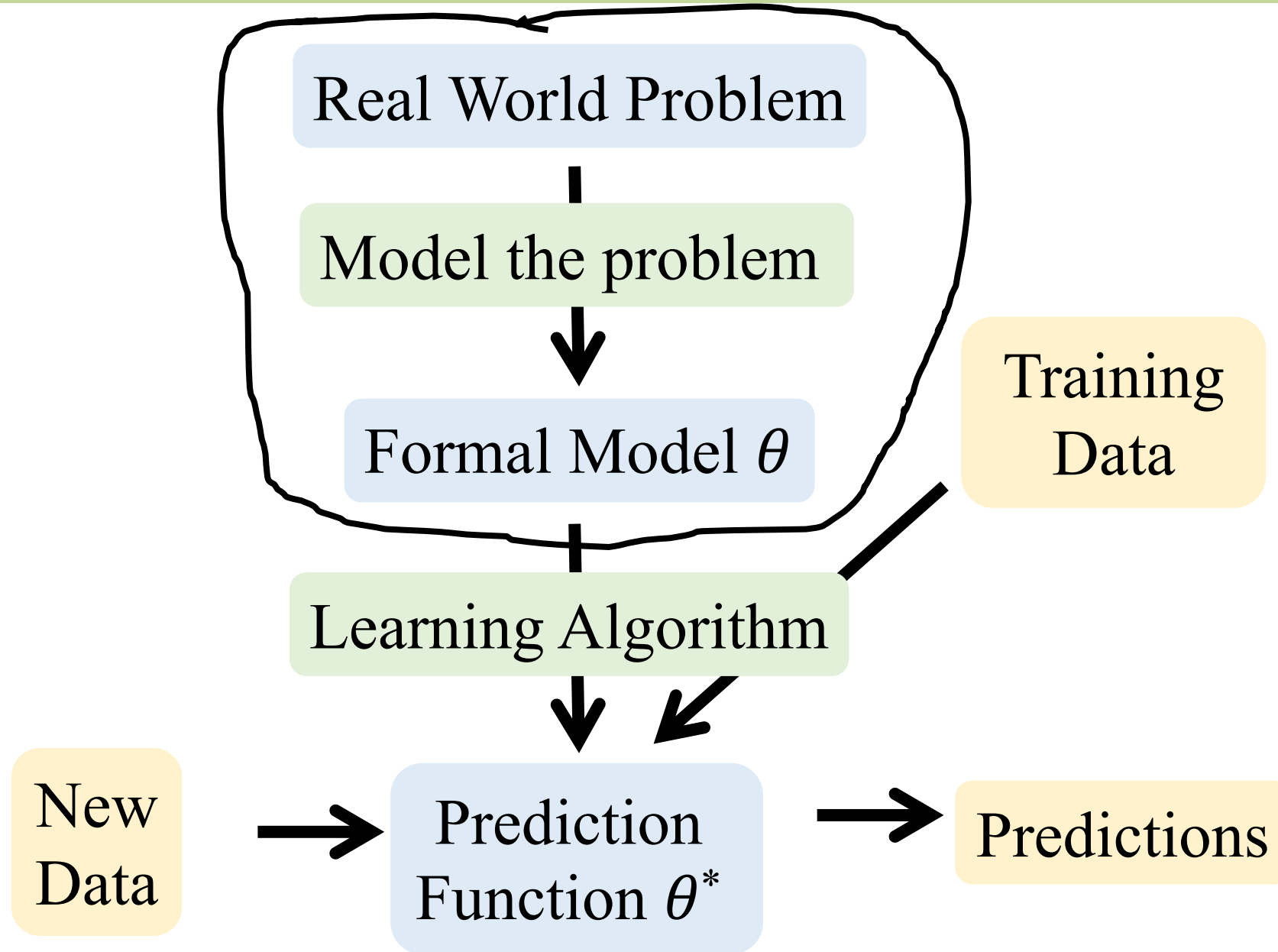
# Naïve Bayes

Noah Arthurs  
CS109, Stanford University

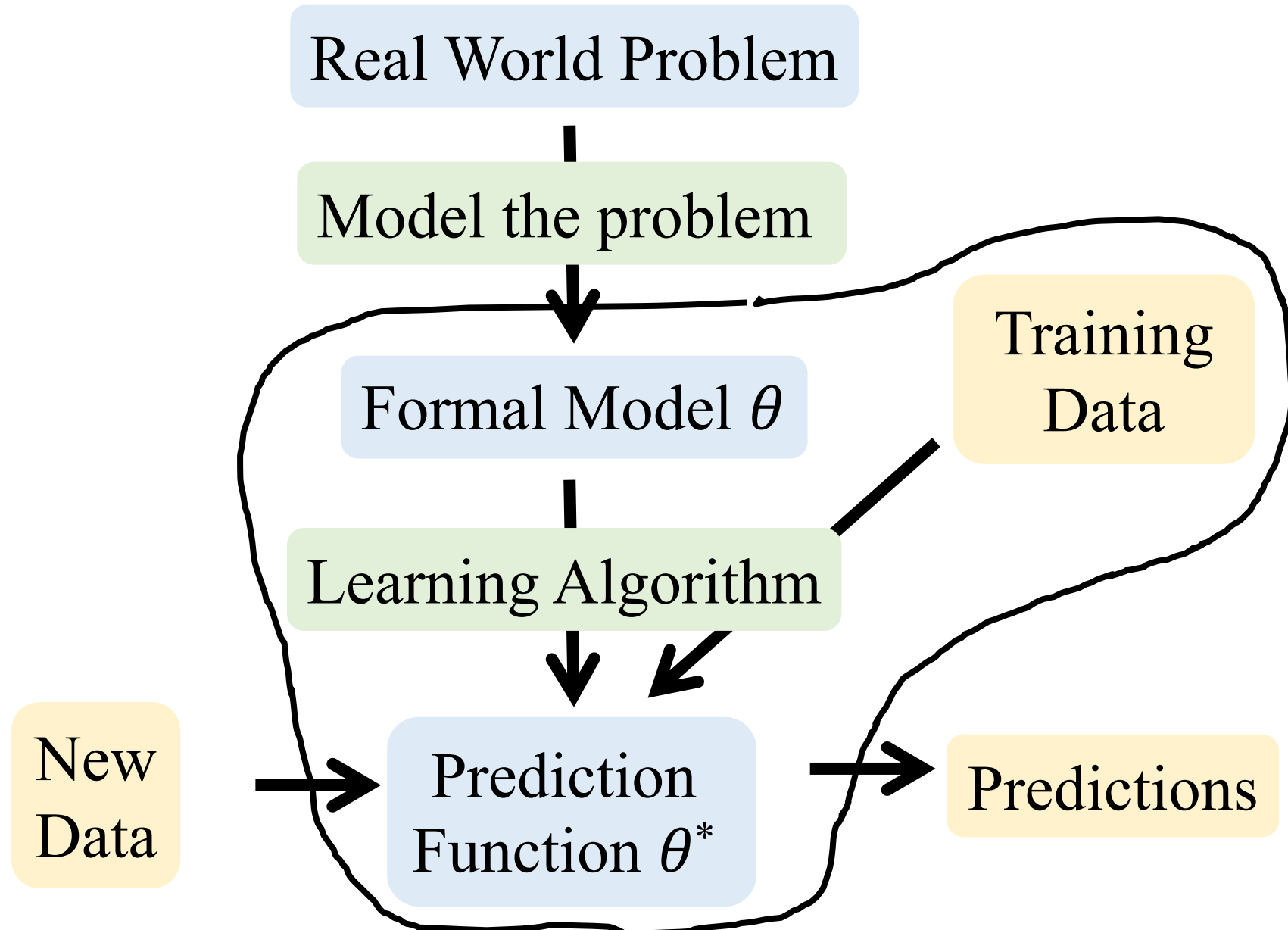
# Supervised Learning



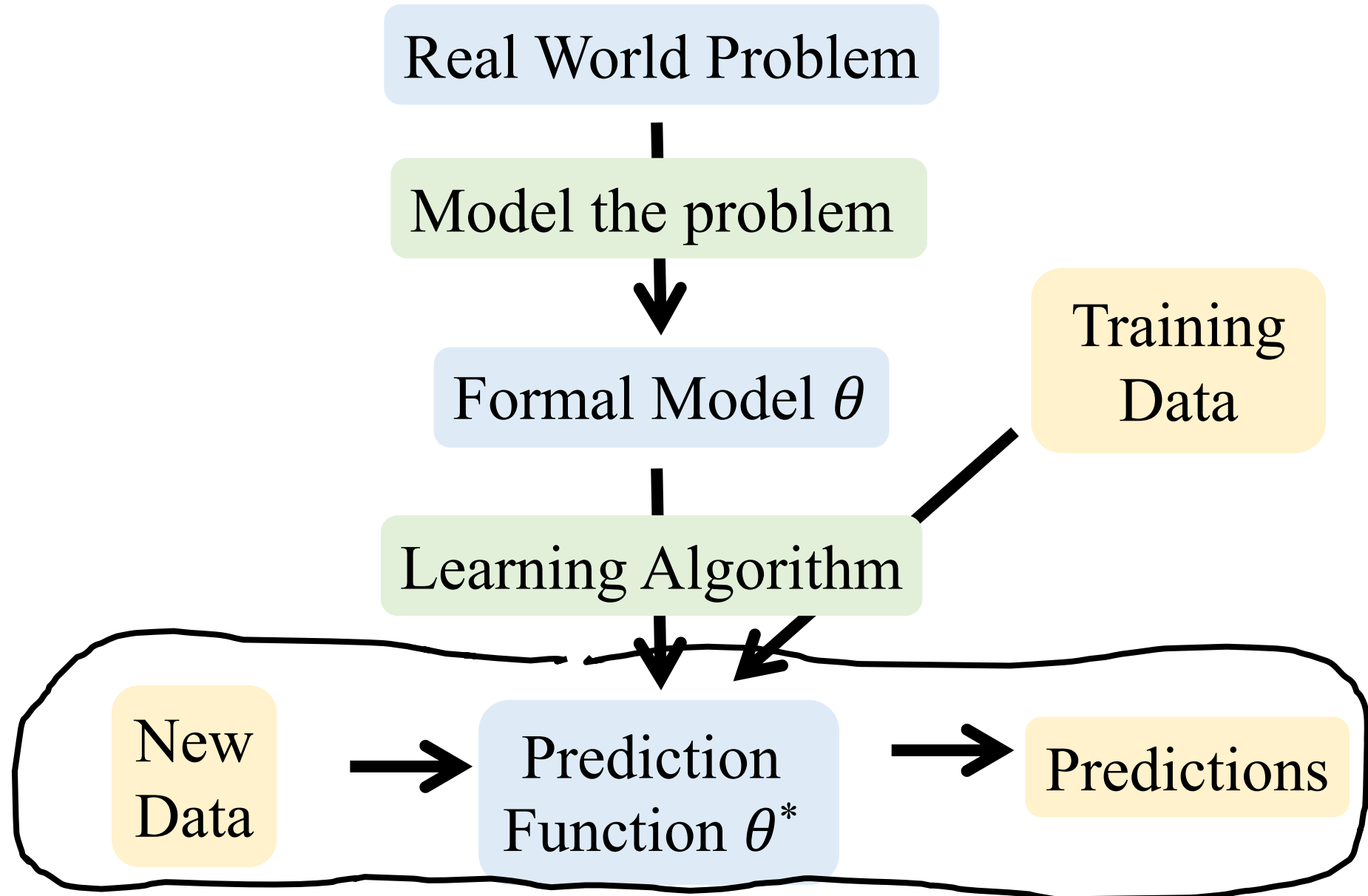
# Modelling



# Training\*



# Make Predictions\*



# Machine Learning: Formally

- Many different forms of “Machine Learning”
  - We focus on the problem of *prediction*
- Want to make a prediction based on observations
  - Vector  $\mathbf{X}$  of  $m$  observed variables:  $\mathbf{X} = [X_1 \dots X_m]$
  - Based on observed  $\mathbf{X}$ , want to predict unseen variable  $Y$ 
    - $Y$  called “output feature/variable” (or the “dependent variable”)
  - Seek to “learn” a function  $g(\mathbf{X})$  to predict  $Y$ :
    - $\hat{Y} = g(\mathbf{X})$
    - When  $Y$  is discrete, prediction of  $Y$  is called “classification”
    - When  $Y$  is continuous, prediction of  $Y$  is called “regression”

# Training Data

Training Data: assignments all random variables  $\mathbf{X}$  and  $Y$

Assume IID data:

*n training datapoints*

$$(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots (\mathbf{x}^{(n)}, y^{(n)})$$

$$m = |\mathbf{x}^{(i)}|$$

Each datapoint has  $m$  features and a single output

# Regression

# Regression: Predicting Real Numbers

Opposing team  
ELO



Points in  
last game



...

At Home?



Output



# Points

Game 1

84

105

1

120

Game 2

90

102

0

95

⋮

⋮

Game  $n$

74

120

0

115

# Linear Regression

Opposing team  
ELO



Points in  
last game



...

At Home?



Output



# Points

Game 1

84

105

1

120

Game 2

90

102

0

95

⋮

⋮

Game  $n$

74

120

0

115

# Linear Regression

$X_1 =$  Opposing team ELO

$X_2 =$  Points in last game

$X_3 =$  Curry playing?

$X_4 =$  Playing at home?

---

$Y =$  Warriors points

# Linear Regression

$Y =$  Warriors points

$$\hat{Y} = \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_{n-1} X_{n-1} + \theta_n 1$$
$$= \theta^T \mathbf{X}$$

---

$X_1 =$  Opposing team ELO

$$\theta_1 = -2.3$$

$X_2 =$  Points in last game

$$\theta_2 = +1.2$$

$X_3 =$  Curry playing?

$$\theta_3 = +10.2$$

$X_4 =$  Playing at home?

$$\theta_4 = +3.3$$

$X_5 = 1$

$$\theta_5 = +95.4$$

# Classification

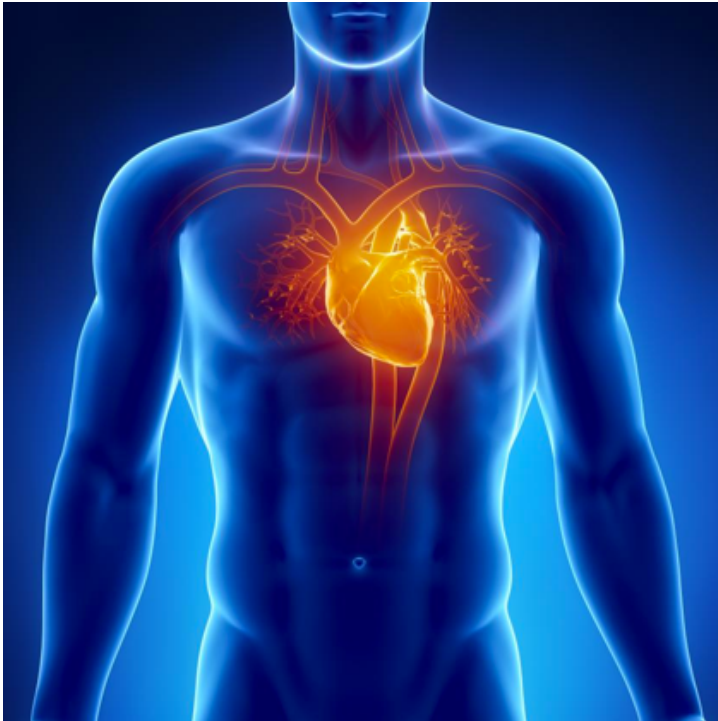
# Classification is Building a Harry Potter Hat



$$\mathbf{x} = [0, 1, \dots, 1]$$

# Example Datasets

Heart



Ancestry

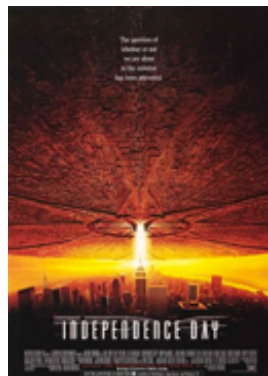


Netflix

The Netflix logo, consisting of the word 'NETFLIX' in a white, bold, sans-serif font with a slight 3D effect, set against a solid red rectangular background.

# Target Movie "Like" Classification

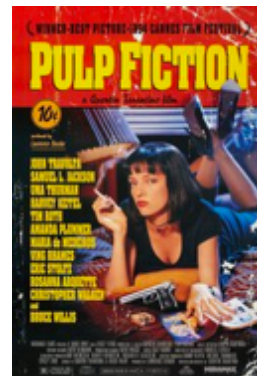
Movie 1



Movie 2



Movie  $m$



...

Output



User 1

1

0

1

1

User 2

1

1

0

0

⋮

⋮

User  $n$


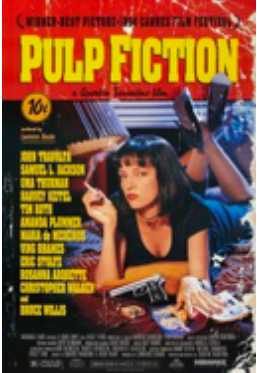

0

0

1

1

# Single Instance

	Movie 1	Movie 2	Movie $m$	Output
				
User 1	1	0	1	1
User 2	1	1	0	0
		⋮		⋮
User $n$	0	0	1	1


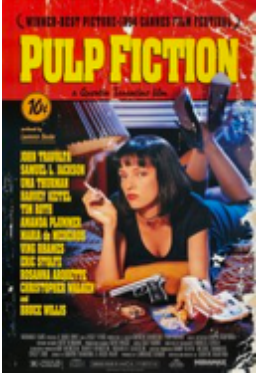
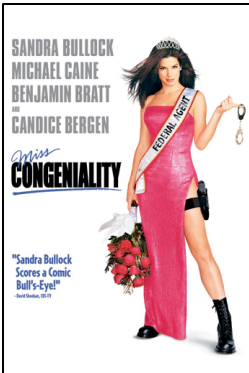
$(\mathbf{x}^{(i)}, y^{(i)})$  such that  $1 \leq i \leq n$

# Feature Vector

	Movie 1	Movie 2	...	Movie $m$	Output
User 1			...		
User 2	1	1		0	0
		⋮			⋮
User $n$	0	0		1	1


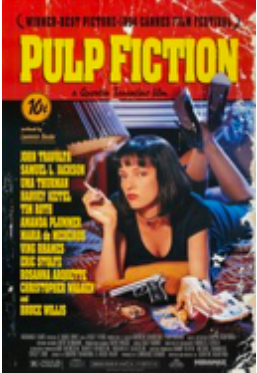
$(\mathbf{x}^{(i)}, y^{(i)})$  such that  $1 \leq i \leq n$

# Output Value

	Movie 1	Movie 2	...	Movie $m$	Output
User 1			...		
User 2	1	1		0	0
			⋮		⋮
User $n$	0	0		1	1

$(\mathbf{x}^{(i)} \quad y^{(i)})$  such that  $1 \leq i \leq n$

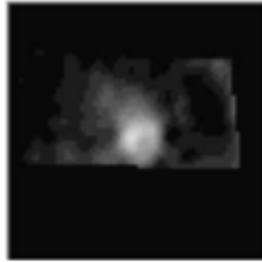
# Single Feature Value

	Movie 1	Movie 2	Movie $m$	Output
				
User 1	1	0	1	1
User 2	1	1	0	0
		⋮		⋮
User $n$	0	0	1	1

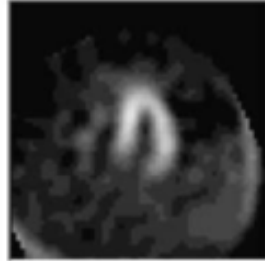
In general:  $x_j^{(i)}$       In this case:  $x_m^{(2)}$

# Healthy Heart Classifier

ROI 1

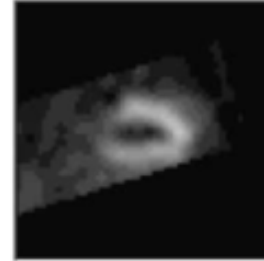


ROI 2



...

ROI  $m$



Output



Heart 1

0

1

1

0

Heart 2

1

1

1

0

⋮

⋮

Heart  $n$

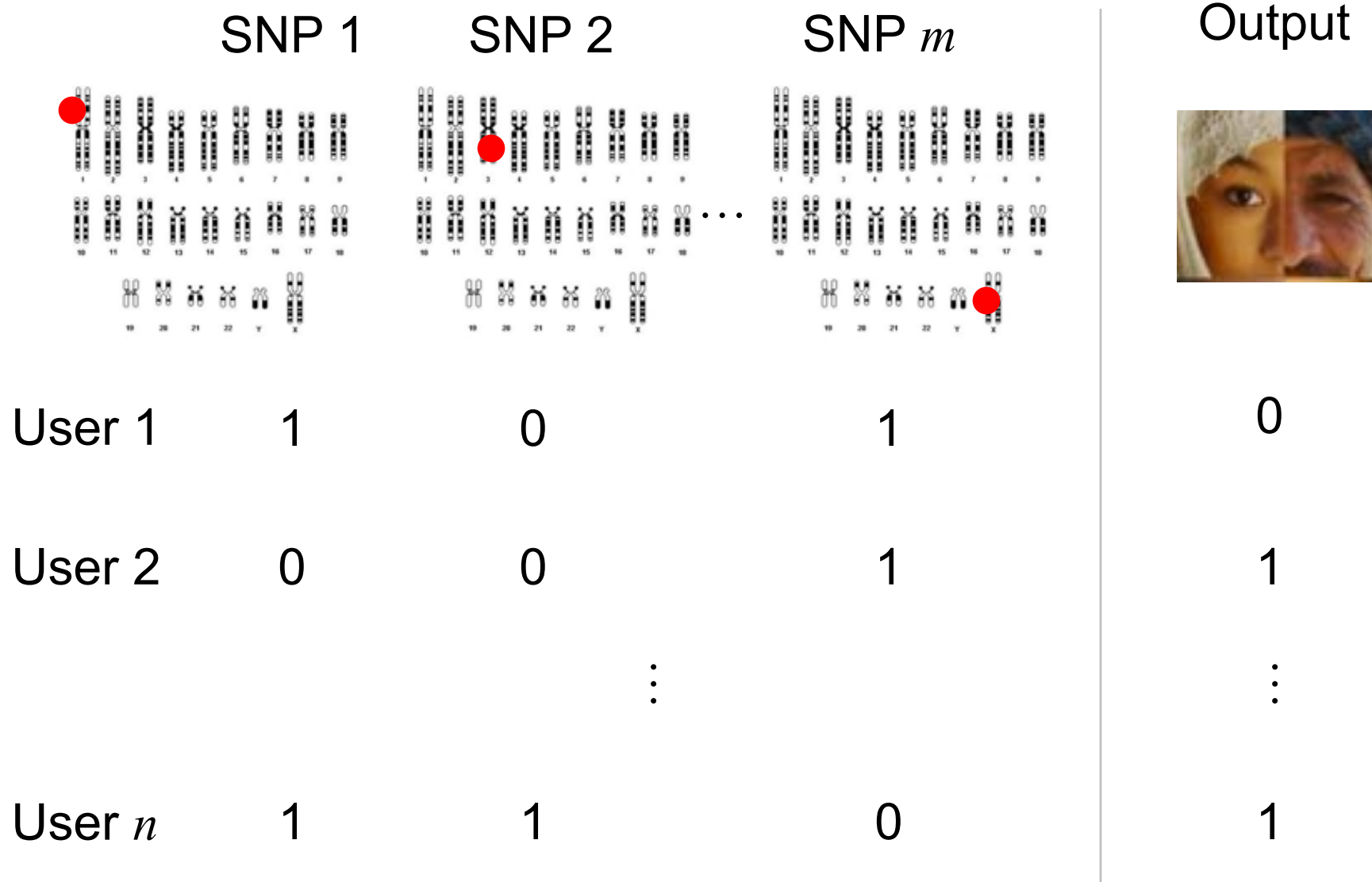
0

0

0

1

# Ancestry Classifier



Stretch!



A woman with a gold coin headband and a dark, patterned top is looking intently into a glowing white crystal ball. Her hands are positioned around the ball, as if holding it. The background is dark, and the lighting is focused on the crystal ball and her face.

# Naïve Bayes

Noah Arthurs  
CS109, Stanford University

**NETFLIX**

**And Learn**

# Target Movie “Like” Classification

Feature 1



Output



User 1

1

1

User 2

1

0

⋮

User  $n$

0

1

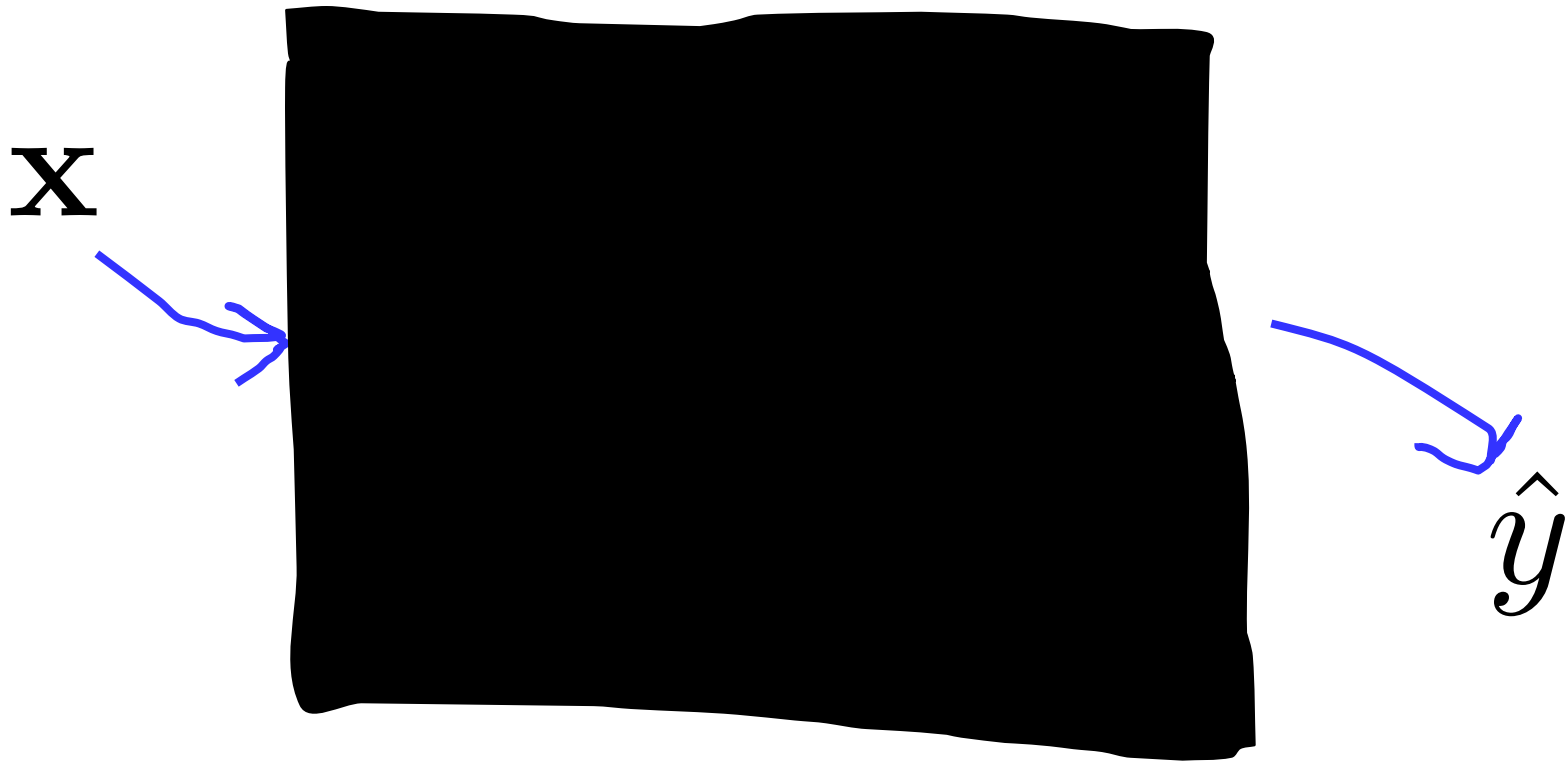
$$x_j^{(i)} \in \{0, 1\}$$

$$y^{(i)} \in \{0, 1\}$$

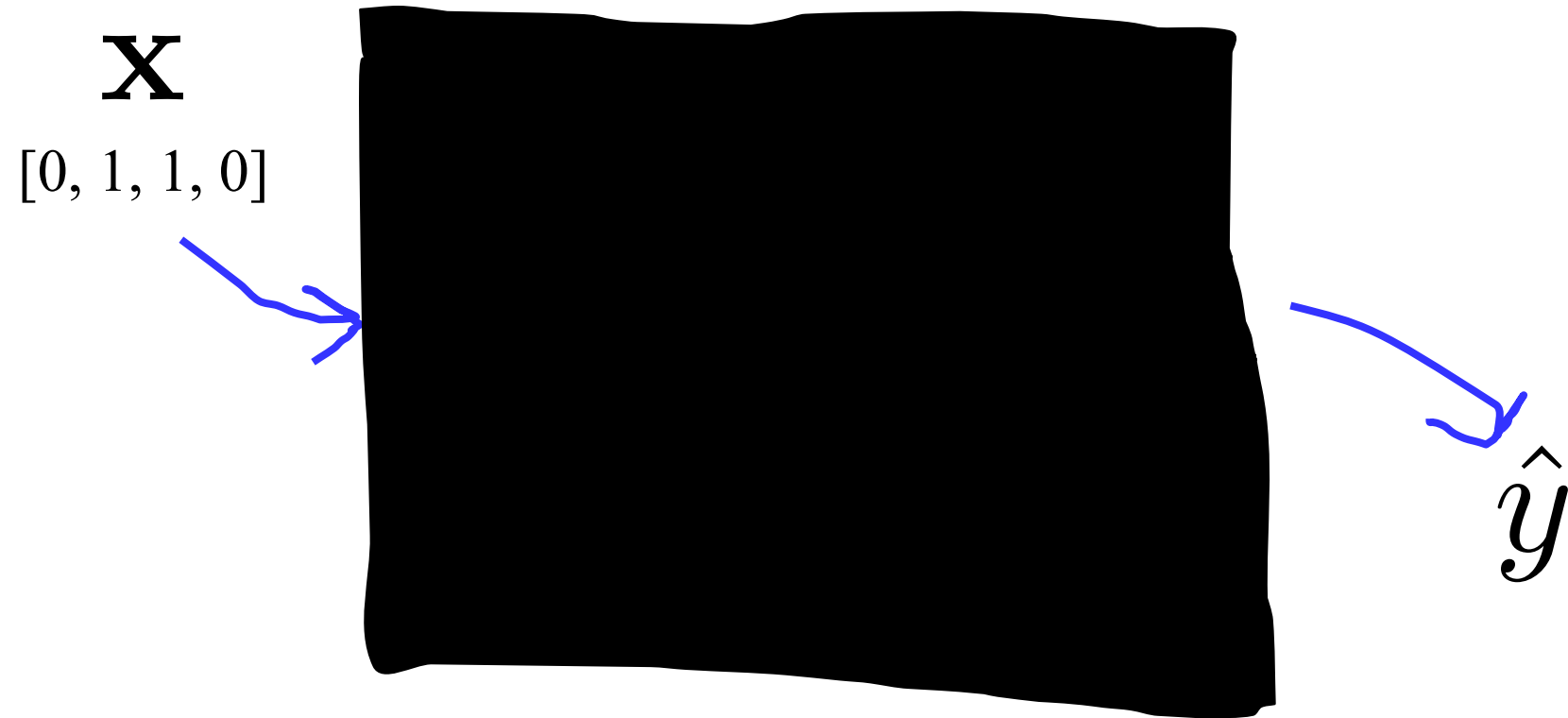
How could we predict the class label:  
will the user like life is beautiful?

# Fake Algorithm: Brute Bayes Classifier

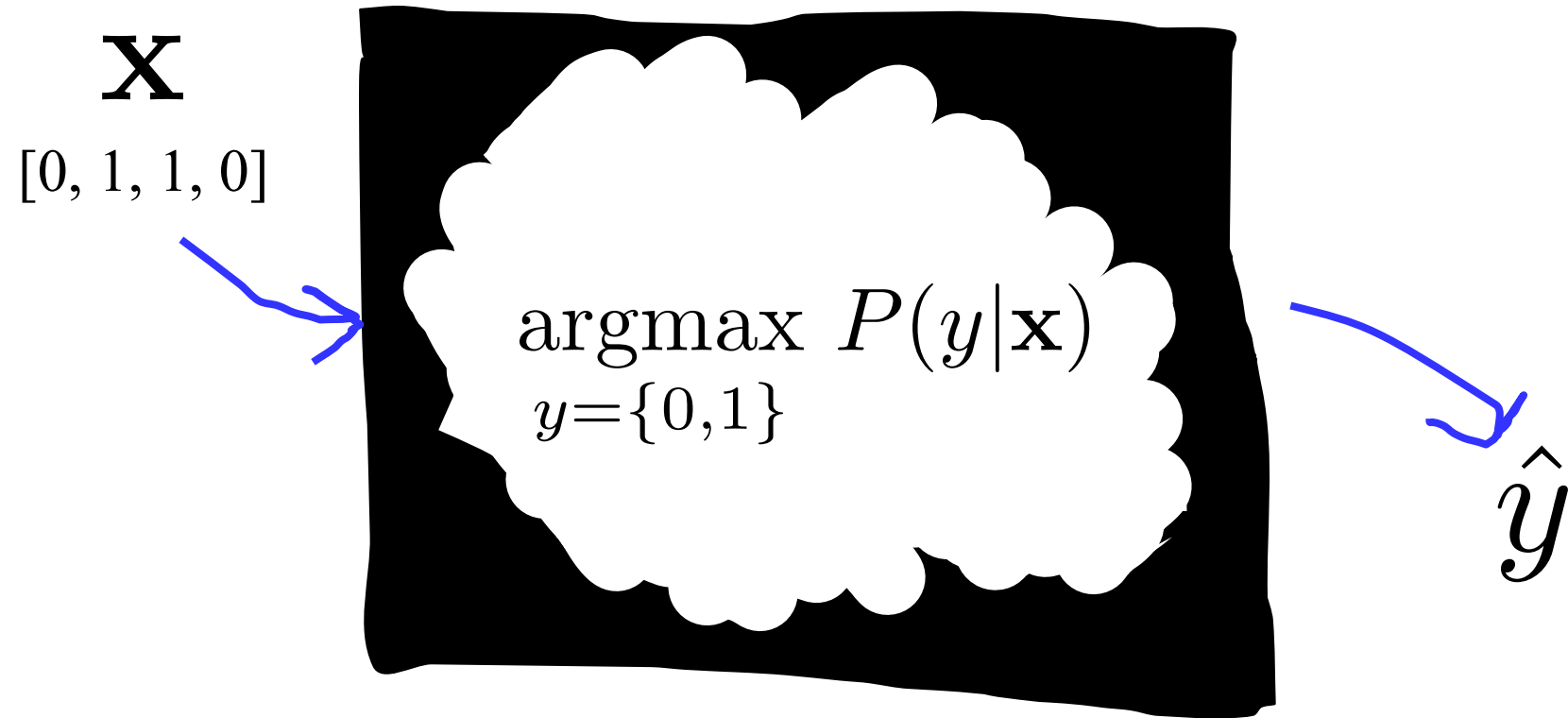
# Brute Force Bayes



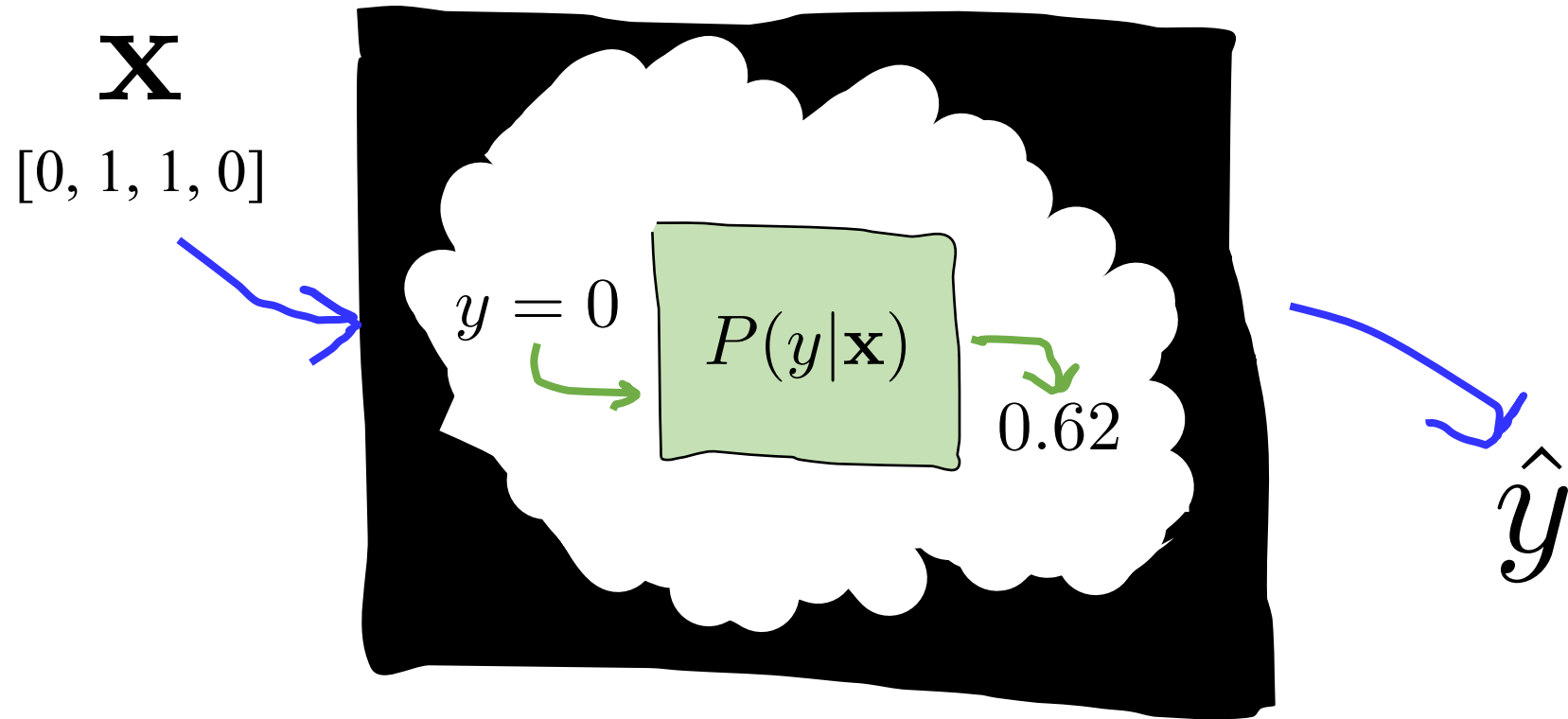
# Brute Force Bayes



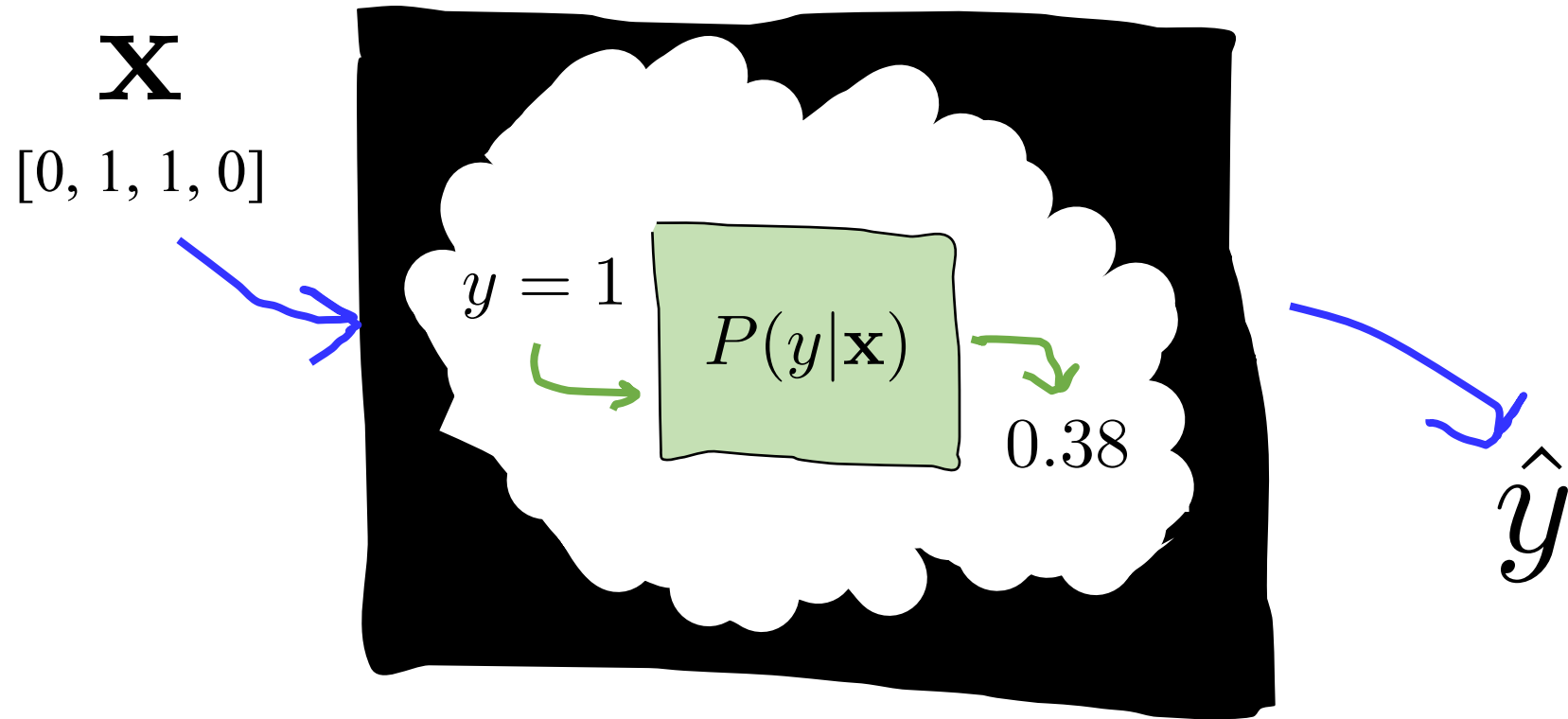
# Brute Force Bayes



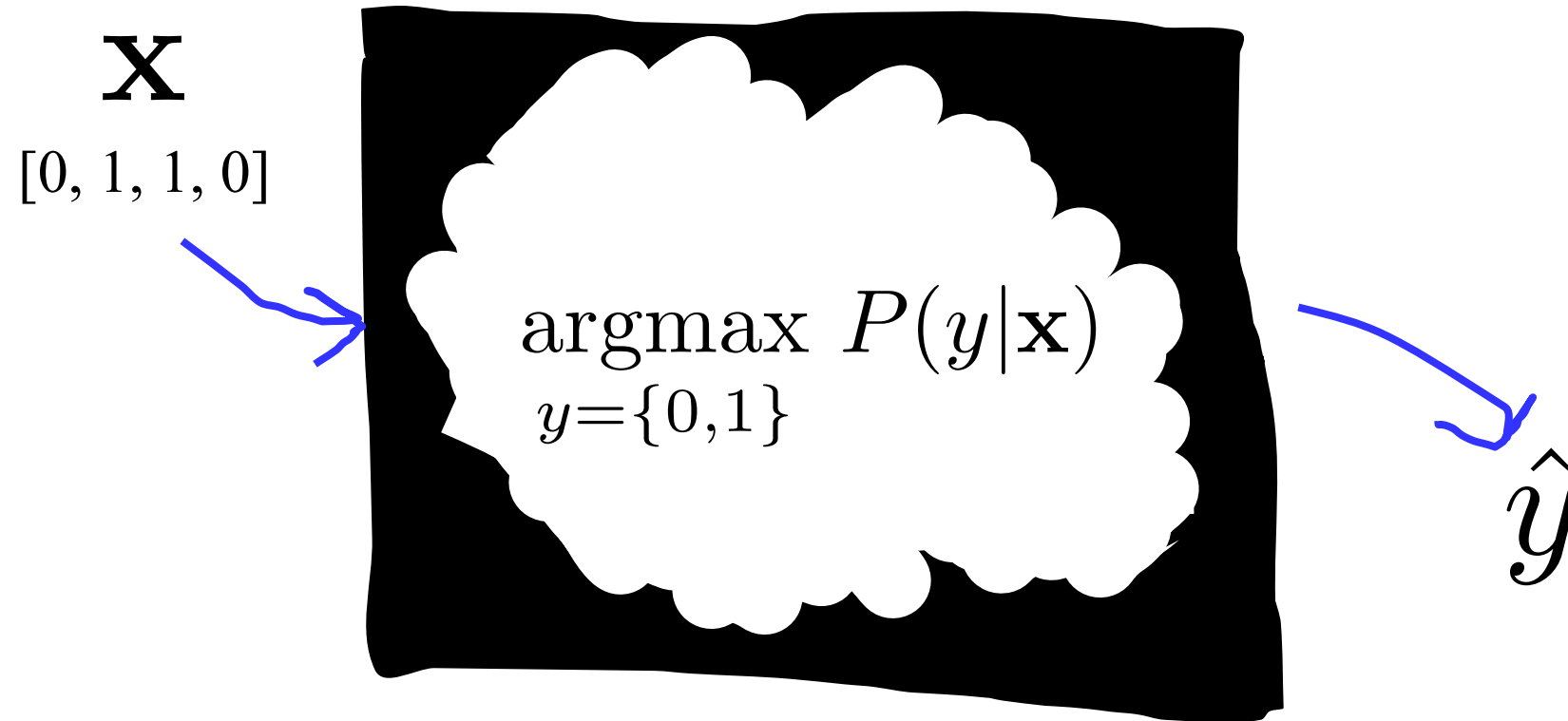
# Brute Force Bayes



# Brute Force Bayes

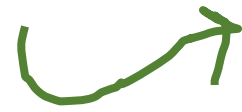


# Brute Force Bayes



# Brute Force Bayes

Prediction: will they like L.I.B.?



$$\hat{y} = \operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x})$$

If  $y = 1$ , they like L.I.B.?



Whether or not they liked Independence day



Simply chose the class label that is the most likely given the data

This is for one user

# Brute Force Bayes

$$\hat{y} = \operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x})$$

Simply chose the class label that is the most likely given the data

This is for one user

# Brute Force Bayes

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x}) \\ &= \operatorname{argmax}_{y=\{0,1\}} \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} \\ &= \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)\end{aligned}$$

Simply chose the class label that is the most likely given the data

This is for one user

\* Note how similar this is to Hamilton example ☺

What are the Parameters?

# Brute Force Bayes

$$\hat{y} = \operatorname{argmax}_{y=\{0,1\}} \underline{P(\mathbf{x}|y)} \underline{P(y)}$$



Conditional probability table



Y = 0

$x_1 = 0$	$\theta_0$
$x_1 = 1$	$\theta_1$

Y = 1

$x_1 = 0$	$\theta_2$
$x_1 = 1$	$\theta_3$



Y = 0	$\theta_4$
Y = 1	$\theta_5$

Learn these during training

# Brute Force Bayes

$$\hat{y} = \operatorname{argmax}_{y=\{0,1\}} \underline{P(\mathbf{x}|y)} \underline{P(y)}$$



Conditional probability table



$x_1 \backslash Y$	0	1
0	$\theta_0$	$\theta_2$
1	$\theta_1$	$\theta_3$



$Y = 0$	$\theta_4$
$Y = 1$	$\theta_5$

Learn these during training

# Training

$x_1$



$y$



$P(\mathbf{x}|y)$

User 1

1

1

User 2

0

0

$\vdots$

User  $n$

0

1

$x_1 \backslash y$	0	1
0	$\theta_0$	$\theta_2$
1	$\theta_1$	$\theta_3$

What is  $P(x_1 | Y = 0)$ ?

What is  $P(x_1 | Y = 1)$ ?

# MLE Estimate

$x_1$



$y$



$P(\mathbf{x}|y)$

User 1

1

1

User 2

0

0

⋮

User  $n$

0

1

$x_1 \backslash y$	0	1
0	0.0	0.4
1	1.0	0.6

MLE: Just count

# MAP Estimate

$x_1$



$y$



$P(\mathbf{x}|y)$

User 1

1

1

User 2

0

0

⋮

User  $n$

0

1

$x_1 \backslash y$	0	1
0	0.01	0.42
1	0.99	0.58

MAP: Just count  
and add imaginary  
trials

# Testing

$$\hat{y} = \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)$$

$x_1 \backslash Y$	0	1
0	0.01	0.42
1	0.99	0.58

$Y=0$	0.21
$Y=1$	0.79

---

Test user: Likes independence day

$$P(x_1 = 1|y = 0)P(y = 0)$$

vs

$$P(x_1 = 1|y = 1)P(y = 1)$$

# Testing

$$\hat{y} = \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)$$

$x_1 \backslash Y$	0	1
0	0.01	0.42
1	0.99	0.58

$Y=0$	0.21
$Y=1$	0.79

---

Test user: Likes independence day

$$P(x_1 = 1|y = 0)P(y = 0)$$

0.208

vs

$$P(x_1 = 1|y = 1)P(y = 1)$$

# Testing

$$\hat{y} = \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)$$

$x_1 \backslash Y$	0	1
0	0.01	0.42
1	0.99	0.58

$Y=0$	0.21
$Y=1$	0.79

---

Test user: Likes independence day




$$P(x_1 = 1|y = 0)P(y = 0) \quad 0.208$$

vs

$$P(x_1 = 1|y = 1)P(y = 1) \quad 0.458$$


That was pretty good!

# Brute Force Bayes $m = 2$

	$x_1$	$x_2$	$y$
			
User 1	1	0	1
User 2	1	0	0
			⋮
User $n$	0	1	1

# Brute Force Bayes $m = 2$

Simply chose the class label that is the most likely given the data

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x}) \\ &= \operatorname{argmax}_{y=\{0,1\}} \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} \\ &= \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)\end{aligned}$$

$$P(x_1, x_2|y)$$

# Brute Force Bayes

$$\hat{y} = \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)$$

		Y = 0		Y = 1	
		X <sub>1</sub>	X <sub>1</sub>	X <sub>1</sub>	X <sub>1</sub>
X <sub>2</sub>	X <sub>1</sub>	0	1	0	1
	0	$\theta_0$	$\theta_1$	$\theta_4$	$\theta_5$
1	$\theta_2$	$\theta_3$	$\theta_6$	$\theta_7$	

X<sub>1</sub>



X<sub>2</sub>

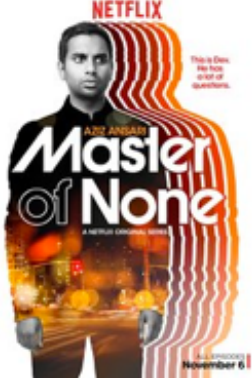


y




Fine

# Brute Force Bayes $m = 3$

	$X_1$	$X_2$	$X_3$	$y$
				
User 1	1	0	1	1
User 2	1	0	1	0
				⋮
User $n$	0	1	1	1

# Brute Force Bayes $m = 3$

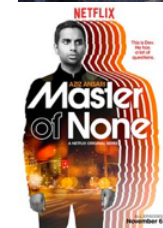
Simply chose the class label that is the most likely given the data

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x}) \\ &= \operatorname{argmax}_{y=\{0,1\}} \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} \\ &= \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)\end{aligned}$$

$$P(x_1, x_2, x_3|y)$$

# Brute Force Bayes

$$\hat{y} = \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)$$

		Y = 0		Y = 1	
		X <sub>1</sub> = 0	X <sub>1</sub> = 1	X <sub>1</sub> = 0	X <sub>1</sub> = 1
X <sub>3</sub> = 0	X <sub>2</sub> = 0	θ <sub>0</sub>	θ <sub>1</sub>	θ <sub>8</sub>	θ <sub>9</sub>
	X <sub>2</sub> = 1	θ <sub>2</sub>	θ <sub>3</sub>	θ <sub>10</sub>	θ <sub>11</sub>
	X <sub>2</sub> = 0	θ <sub>4</sub>	θ <sub>5</sub>	θ <sub>12</sub>	θ <sub>13</sub>
X <sub>3</sub> = 1	X <sub>2</sub> = 0	θ <sub>6</sub>	θ <sub>7</sub>	θ <sub>14</sub>	θ <sub>15</sub>
	X <sub>2</sub> = 1				
	X <sub>2</sub> = 0				



And if  $m=100$ ?

# Brute Force Bayes $m = 100$

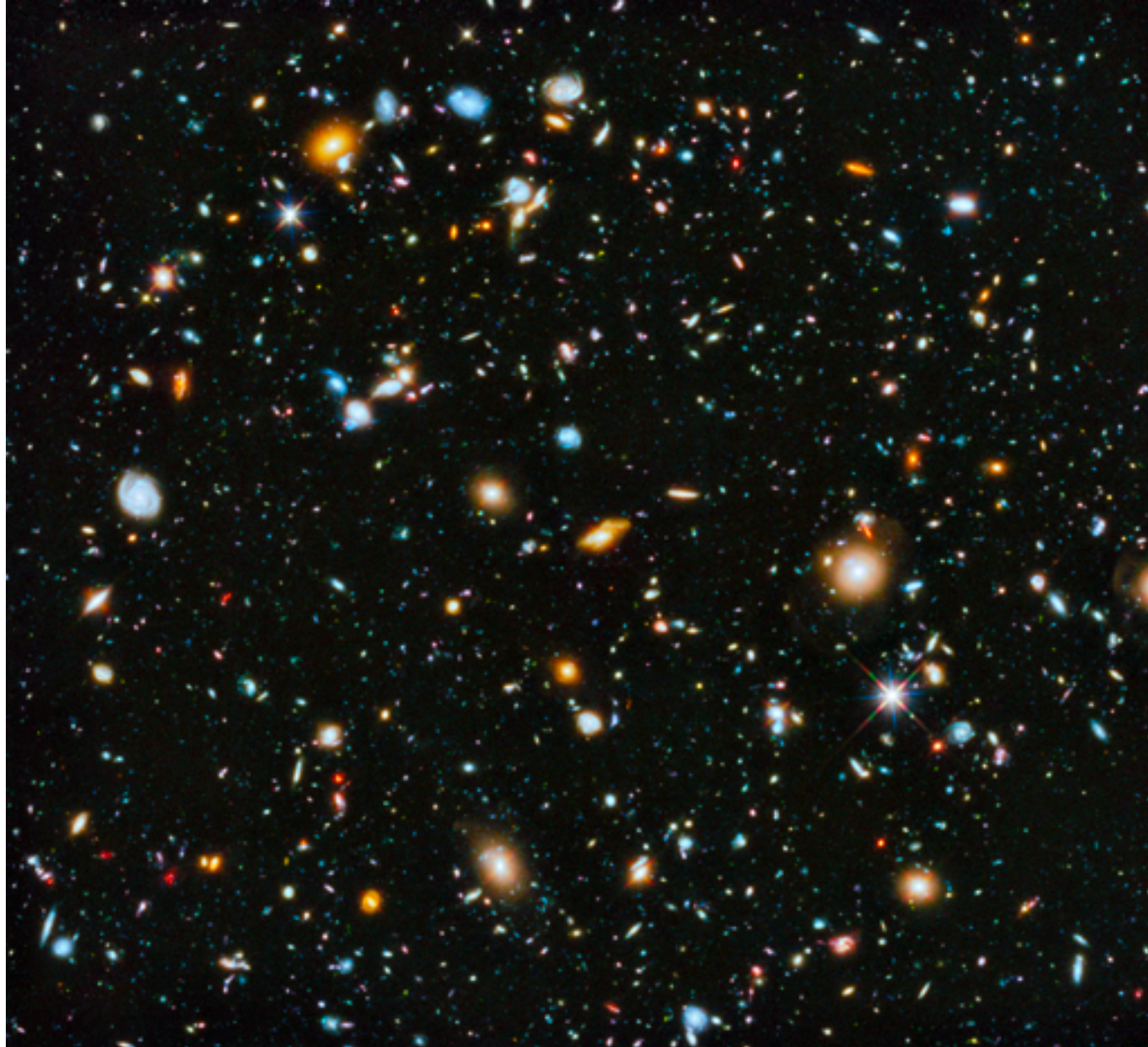
Simply chose the class label that is the most likely given the data

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x}) \\ &= \operatorname{argmax}_{y=\{0,1\}} \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} \\ &= \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)\end{aligned}$$



$$P(x_1, x_2, x_3, \dots, x_{100}|y)$$

# Oops... Number of atoms in the universe



What is the big O for # parameters?  
 $m = \# \text{ features.}$

# Big O of Brute Force Joint

What is the big O for # parameters?  
 $m = \# \text{ features.}$

$$O(2^m)$$

*Assuming each feature  
is binary...*

Not going to cut it!

# What is the problem here?

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x}) \\ &= \operatorname{argmax}_{y=\{0,1\}} \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} \\ &= \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)\end{aligned}$$

---

$$P(\mathbf{x}|y) = P(x_1, x_2, \dots, x_m|y)$$

# Naïve Bayes Assumption

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x}) \\ &= \operatorname{argmax}_{y=\{0,1\}} \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} \\ &= \operatorname{argmax}_{y=\{0,1\}} P(\mathbf{x}|y)P(y)\end{aligned}$$

---

$$\begin{aligned}P(\mathbf{x}|y) &= P(x_1, x_2, \dots, x_m|y) \\ &= \prod_i P(x_i|y)\end{aligned}$$

*The Naïve Bayes  
assumption*



Naïve Bayes Assumption:

$$P(\mathbf{x}|y) = \prod_i P(x_i|y)$$

# Naïve Bayes Classifier

# Naïve Bayes

Our prediction for  $y$

Is a function of  $\mathbf{x}$

That chooses the best value of  $y$  given  $\mathbf{x}$

$$\hat{y} = g(\mathbf{x}) = \operatorname{argmax}_{y \in \{0,1\}} \hat{P}(y|\mathbf{x})$$

$$= \operatorname{argmax}_{y \in \{0,1\}} \hat{P}(\mathbf{x}|y) \hat{P}(y)$$

Bayes rule!

$$= \operatorname{argmax}_y \left( \prod_{i=1}^n \hat{P}(x_i|y) \right) \hat{P}(y)$$

Naïve Bayes Assumption

$$= \operatorname{argmax}_y \log \hat{P}(y) + \sum_{i=1}^m \log \hat{P}(x_i|y)$$

This log version is useful for numerical stability

# Naïve Bayes Example

- Predict  $Y$  based on observing variables  $X_1$  and  $X_2$ 
  - $X_1$  and  $X_2$  are both indicator variables
    - $X_1$  denotes “likes Star Wars”,  $X_2$  denotes “likes Harry Potter”
  - $Y$  is indicator variable: “likes Lord of the Rings”
    - Use training data to estimate params:  $\hat{P}(x_i|y)$   $\hat{P}(y)$

$Y \backslash X_1$	0	1	MLE estimates		$Y \backslash X_2$	0	1	MLE estimates		$Y$	#	MLE est.
0	3	10	0.23	0.77	0	5	8	0.38	0.62	0	13	0.43
1	4	13	0.24	0.76	1	7	10	0.41	0.59	1	17	0.57

- Say someone likes **Star Wars ( $X_1 = 1$ )**, but not **Harry Potter ( $X_2 = 0$ )**
- Will they like “Lord of the Rings”? Need to predict  $Y$ :

$$\hat{y} = \operatorname{argmax}_{y \in \{0,1\}} \hat{P}(\mathbf{x}|y)\hat{P}(y) = \operatorname{argmax}_{y \in \{0,1\}} \hat{P}(x_1|y)\hat{P}(x_2|y)\hat{P}(y)$$

# Naïve Bayes Example

- Predict  $Y$  based on observing variables  $X_1$  and  $X_2$ 
  - $X_1$  and  $X_2$  are both indicator variables
    - $X_1$  denotes “likes Star Wars”,  $X_2$  denotes “likes Harry Potter”
  - $Y$  is indicator variable: “likes Lord of the Rings”
    - Use training data to estimate params:  $\hat{P}(x_i|y)$   $\hat{P}(y)$

$Y \backslash X_1$	0	1	MLE estimates		$Y \backslash X_2$	0	1	MLE estimates		$Y$	#	MLE est.
0	3	10	0.23	0.77	0	5	8	0.38	0.62	0	13	0.43
1	4	13	0.24	0.76	1	7	10	0.41	0.59	1	17	0.57

- Say someone likes **Star Wars ( $X_1 = 1$ )**, but not **Harry Potter ( $X_2 = 0$ )**
- Will they like “Lord of the Rings”? Need to predict  $Y$ .

$$\hat{y} = \operatorname{argmax}_{y \in \{0,1\}} \hat{P}(X_1 = x_1 | Y = y) \hat{P}(X_2 = x_2 | Y = y) \hat{P}(Y = y)$$

# Naïve Bayes Example

- Predict  $Y$  based on observing variables  $X_1$  and  $X_2$ 
  - $X_1$  and  $X_2$  are both indicator variables
    - $X_1$  denotes “likes Star Wars”,  $X_2$  denotes “likes Harry Potter”
  - $Y$  is indicator variable: “likes Lord of the Rings”
    - Use training data to estimate params:  $\hat{P}(x_i|y)$   $\hat{P}(y)$

$Y \backslash X_1$	0	1	MLE estimates		$Y \backslash X_2$	0	1	MLE estimates		$Y$	#	MLE est.
0	3	10	0.23	0.77	0	5	8	0.38	0.62	0	13	0.43
1	4	13	0.24	0.76	1	7	10	0.41	0.59	1	17	0.57

- Say someone likes **Star Wars ( $X_1 = 1$ )**, but not **Harry Potter ( $X_2 = 0$ )**
- Will they like “Lord of the Rings”? Need to predict  $Y$ :

$$\hat{y} = \operatorname{argmax}_{y \in \{0,1\}} \hat{P}(X_1 = 1|Y = y) \hat{P}(X_2 = 0|Y = y) \hat{P}(Y = y)$$

# Naïve Bayes Example

$Y \backslash X_1$	0	1	MLE estimates		$Y \backslash X_2$	0	1	MLE estimates		Y	#	MLE est.
0	3	10	0.23	0.77	0	5	8	0.38	0.62	0	13	0.43
1	4	13	0.24	0.76	1	7	10	0.41	0.59	1	17	0.57

$$\hat{y} = \operatorname{argmax}_{y \in \{0,1\}} \hat{P}(X_1 = 1|Y = y)\hat{P}(X_2 = 0|Y = y)\hat{P}(Y = y)$$

- Let  $Y = 0$ 

$$\hat{P}(X_1 = 1|Y = 0)\hat{P}(X_2 = 0|Y = 0)\hat{P}(Y = 0)$$

$$= (0.77)(0.38)(0.43) = 0.126$$
- Let  $Y = 1$ 

$$\hat{P}(X_1 = 1|Y = 1)\hat{P}(X_2 = 0|Y = 1)\hat{P}(Y = 1)$$

$$= (0.76)(0.41)(0.57) = 0.178$$

Since term is greatest when  $Y = 1$ , we predict  $\hat{Y} = 1$

$$P(Y = 1) = K \cdot 0.178 \quad P(Y = 0) = K \cdot 0.126 \quad K = \frac{1}{0.126 + 0.178}$$

# MAP Naïve Bayes

- Predict Y based on observing variables  $X_1$  and  $X_2$ 
  - $X_1$  and  $X_2$  are both indicator variables
    - $X_1$  denotes “likes Star Wars”,  $X_2$  denotes “likes Harry Potter”
  - Y is indicator variable: “likes Lord of the Rings”
    - Use training data to estimate PMFs:  $\hat{P}(x_i|y)$   $\hat{P}(y)$

$X_1 \backslash Y$	0	1	MAP estimates
0	3	10	
1	4	13	

$X_2 \backslash Y$	0	1	MAP estimates
0	5	8	
1	7	10	

Y	#	MAP est.
0	13	
1	17	

What prior?

# MAP Naïve Bayes

- Predict Y based on observing variables  $X_1$  and  $X_2$ 
  - $X_1$  and  $X_2$  are both indicator variables
    - $X_1$  denotes “likes Star Wars”,  $X_2$  denotes “likes Harry Potter”
  - Y is indicator variable: “likes Lord of the Rings”
    - Use training data to estimate PMFs:  $\hat{P}(x_i|y)$   $\hat{P}(y)$

$X_1 \backslash Y$	0	1	MAP estimates	
0	3	10	0.27	0.73
1	4	13		

$X_2 \backslash Y$	0	1	MAP estimates	
0	5	8		
1	7	10		

Y	#	MAP est.
0	13	
1	17	

Laplace!

$$p_i = \frac{n_i + 1}{n + m}$$

$$p_i = \frac{n_i + 1}{n + 2}$$

# MAP Naïve Bayes

- Predict  $Y$  based on observing variables  $X_1$  and  $X_2$ 
  - $X_1$  and  $X_2$  are both indicator variables
    - $X_1$  denotes “likes Star Wars”,  $X_2$  denotes “likes Harry Potter”
  - $Y$  is indicator variable: “likes Lord of the Rings”
    - Use training data to estimate PMFs:  $\hat{P}(x_i|y)$   $\hat{P}(y)$

$Y \backslash X_1$	0	1	MAP estimates	
0	3	10	0.27	0.73
1	4	13	0.26	0.74

$Y \backslash X_2$	0	1	MAP estimates	
0	5	8	0.4	0.6
1	7	10	0.42	0.58

$Y$	#	MAP est.
0	13	0.44
1	17	0.56

Laplace!

$$p_i = \frac{n_i + 1}{n + m}$$

$$p_i = \frac{n_i + 1}{n + 2}$$



Training Naïve Bayes is just estimating conditional probabilities.

Thus training is just counting.

# What is Bayes Doing in my Mail Server

- This is spam:

<p>From: Abey Chavez [tristranu@deletemail.com] To: sahani@robotics.stanford.edu Cc: Subject: For excellent metabolism</p> <p><b>Canadian Pharmacy</b> #1 Internet Online Dispensary</p> <table border="0"><tr><td><b>Viagra</b> Our price <b>\$1.15</b></td><td><b>Cialis</b> Our price <b>\$1.99</b></td><td><b>Viagra Professional</b> Our price <b>\$3.73</b></td></tr><tr><td><b>Cialis Professional</b> Our price <b>\$4.17</b></td><td><b>Viagra Super Active</b> Our price <b>\$2.82</b></td><td><b>Cialis Super Active</b> Our price <b>\$3.66</b></td></tr><tr><td><b>Levitra</b> Our price <b>\$2.93</b></td><td><b>Viagra Soft Tabs</b> Our price <b>\$1.64</b></td><td><b>Cialis Soft Tabs</b> Our price <b>\$3.51</b></td></tr></table> <p>And more...</p> <p><a href="#">Click here</a></p>	<b>Viagra</b> Our price <b>\$1.15</b>	<b>Cialis</b> Our price <b>\$1.99</b>	<b>Viagra Professional</b> Our price <b>\$3.73</b>	<b>Cialis Professional</b> Our price <b>\$4.17</b>	<b>Viagra Super Active</b> Our price <b>\$2.82</b>	<b>Cialis Super Active</b> Our price <b>\$3.66</b>	<b>Levitra</b> Our price <b>\$2.93</b>	<b>Viagra Soft Tabs</b> Our price <b>\$1.64</b>	<b>Cialis Soft Tabs</b> Our price <b>\$3.51</b>	<h2>Let's get Bayesian on your spam:</h2> <p>Content analysis details: (49.5 hits, 7.0 required)</p> <table border="0"><tr><td>0.9 RCVD_IN_PBL</td><td>RBL: Received via a relay in Spamhaus PBL [93.40.189.29 listed in zen.spamhaus.org]</td></tr><tr><td>1.5 URIBL_WS_SURBL</td><td>Contains an URL listed in the WS SURBL blacklist [URIs: recragas.cn]</td></tr><tr><td>5.0 URIBL_JP_SURBL</td><td>Contains an URL listed in the JP SURBL blacklist [URIs: recragas.cn]</td></tr><tr><td>5.0 URIBL_OB_SURBL</td><td>Contains an URL listed in the OB SURBL blacklist [URIs: recragas.cn]</td></tr><tr><td>5.0 URIBL_SC_SURBL</td><td>Contains an URL listed in the SC SURBL blacklist [URIs: recragas.cn]</td></tr><tr><td>2.0 URIBL_BLACK</td><td>Contains an URL listed in the URIBL blacklist [URIs: recragas.cn]</td></tr></table> <p><b>8.0 BAYES_99</b> <b>BODY: Bayesian spam probability is 99 to 100% [score: 1.0000]</b></p>	0.9 RCVD_IN_PBL	RBL: Received via a relay in Spamhaus PBL [93.40.189.29 listed in zen.spamhaus.org]	1.5 URIBL_WS_SURBL	Contains an URL listed in the WS SURBL blacklist [URIs: recragas.cn]	5.0 URIBL_JP_SURBL	Contains an URL listed in the JP SURBL blacklist [URIs: recragas.cn]	5.0 URIBL_OB_SURBL	Contains an URL listed in the OB SURBL blacklist [URIs: recragas.cn]	5.0 URIBL_SC_SURBL	Contains an URL listed in the SC SURBL blacklist [URIs: recragas.cn]	2.0 URIBL_BLACK	Contains an URL listed in the URIBL blacklist [URIs: recragas.cn]
<b>Viagra</b> Our price <b>\$1.15</b>	<b>Cialis</b> Our price <b>\$1.99</b>	<b>Viagra Professional</b> Our price <b>\$3.73</b>																				
<b>Cialis Professional</b> Our price <b>\$4.17</b>	<b>Viagra Super Active</b> Our price <b>\$2.82</b>	<b>Cialis Super Active</b> Our price <b>\$3.66</b>																				
<b>Levitra</b> Our price <b>\$2.93</b>	<b>Viagra Soft Tabs</b> Our price <b>\$1.64</b>	<b>Cialis Soft Tabs</b> Our price <b>\$3.51</b>																				
0.9 RCVD_IN_PBL	RBL: Received via a relay in Spamhaus PBL [93.40.189.29 listed in zen.spamhaus.org]																					
1.5 URIBL_WS_SURBL	Contains an URL listed in the WS SURBL blacklist [URIs: recragas.cn]																					
5.0 URIBL_JP_SURBL	Contains an URL listed in the JP SURBL blacklist [URIs: recragas.cn]																					
5.0 URIBL_OB_SURBL	Contains an URL listed in the OB SURBL blacklist [URIs: recragas.cn]																					
5.0 URIBL_SC_SURBL	Contains an URL listed in the SC SURBL blacklist [URIs: recragas.cn]																					
2.0 URIBL_BLACK	Contains an URL listed in the URIBL blacklist [URIs: recragas.cn]																					

### A Bayesian Approach to Filtering Junk E-Mail

Mehran Sahami\*   Susan Dumais†   David Heckerman†   Eric Horvitz†

\*Gates Building 1A  
Computer Science Department  
Stanford University  
Stanford, CA 94305-5010  
sahami@cs.stanford.edu

†Microsoft Research  
Redmond, WA 98052-6399  
(sdumais, heckerma, horvitz}@microsoft.com

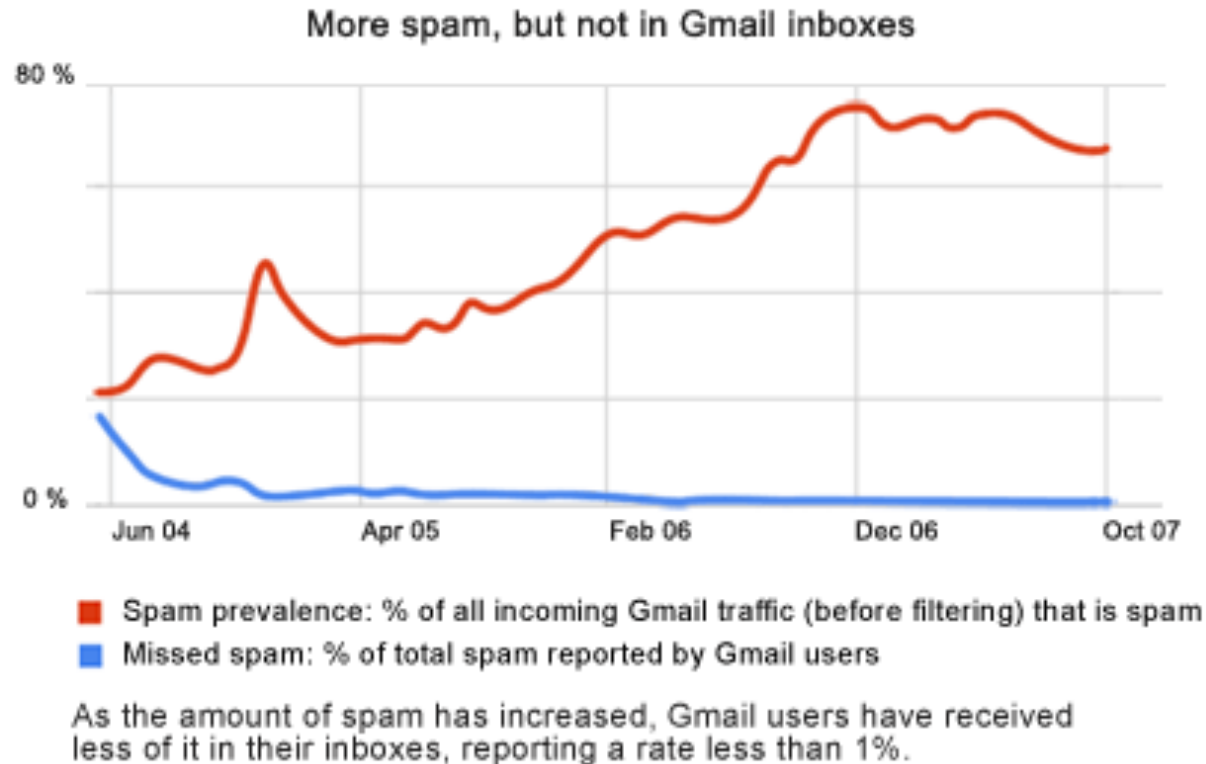
**Abstract**

In addressing the growing problem of junk E-mail on the Internet, we examine methods for the automated

contains offensive material (such as graphic pornography), there is often a higher cost to users of actually viewing this mail than simply the time to sort out the junk. Lastly, junk mail not only wastes user time, but

# Spam, Spam... Go Away!

- The constant battle with spam



*“And machine-learning algorithms developed to merge and rank large sets of Google search results allow us to combine hundreds of factors to classify spam.”*

Source: <http://www.google.com/mail/help/fightspam/spamexplained.html>

# Email Classification

- Want to predict if an email is spam or not
  - Start with the input data
    - Consider a lexicon of  $m$  words (Note: in English  $m \approx 100,000$ )
    - Define  $m$  indicator variables  $\mathbf{X} = \langle X_1, X_2, \dots, X_m \rangle$
    - Each variable  $X_i$  denotes if word  $i$  appeared in a document or not
    - Note:  $m$  is huge, so make “Naive Bayes” assumption
  - Define output classes  $Y$  to be: {spam, non-spam}
  - Given training set of  $N$  previous emails
    - For each email message, we have a training instance:  
 $\mathbf{X} = \langle X_1, X_2, \dots, X_m \rangle$  noting for each word, if it appeared in email
    - Each email message is also marked as spam or not (value of  $Y$ )

# Training the Classifier

- Given  $N$  training pairs:

$$(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})$$

- Learning

- Estimate probabilities  $P(y)$  and  $P(x_i | y)$  for all  $i$

- Many words are likely to not appear at all in given set of email

- Laplace estimate:  $\hat{p}(X_i = 1 | Y = spam)_{Laplace} = \frac{(\# \text{spam emails with word } i) + 1}{\text{total \# spam emails} + 2}$

- Classification

- For a new email, generate  $\mathbf{X} = \langle X_1, X_2, \dots, X_m \rangle$

- Classify as spam or not using:  $\hat{y} = \operatorname{argmax}_{y \in \{0,1\}} \hat{P}(\mathbf{x}|y) \hat{P}(y)$

- Employ Naive Bayes assumption:  $P(\mathbf{x}|y) = \prod_i P(x_i|y)$

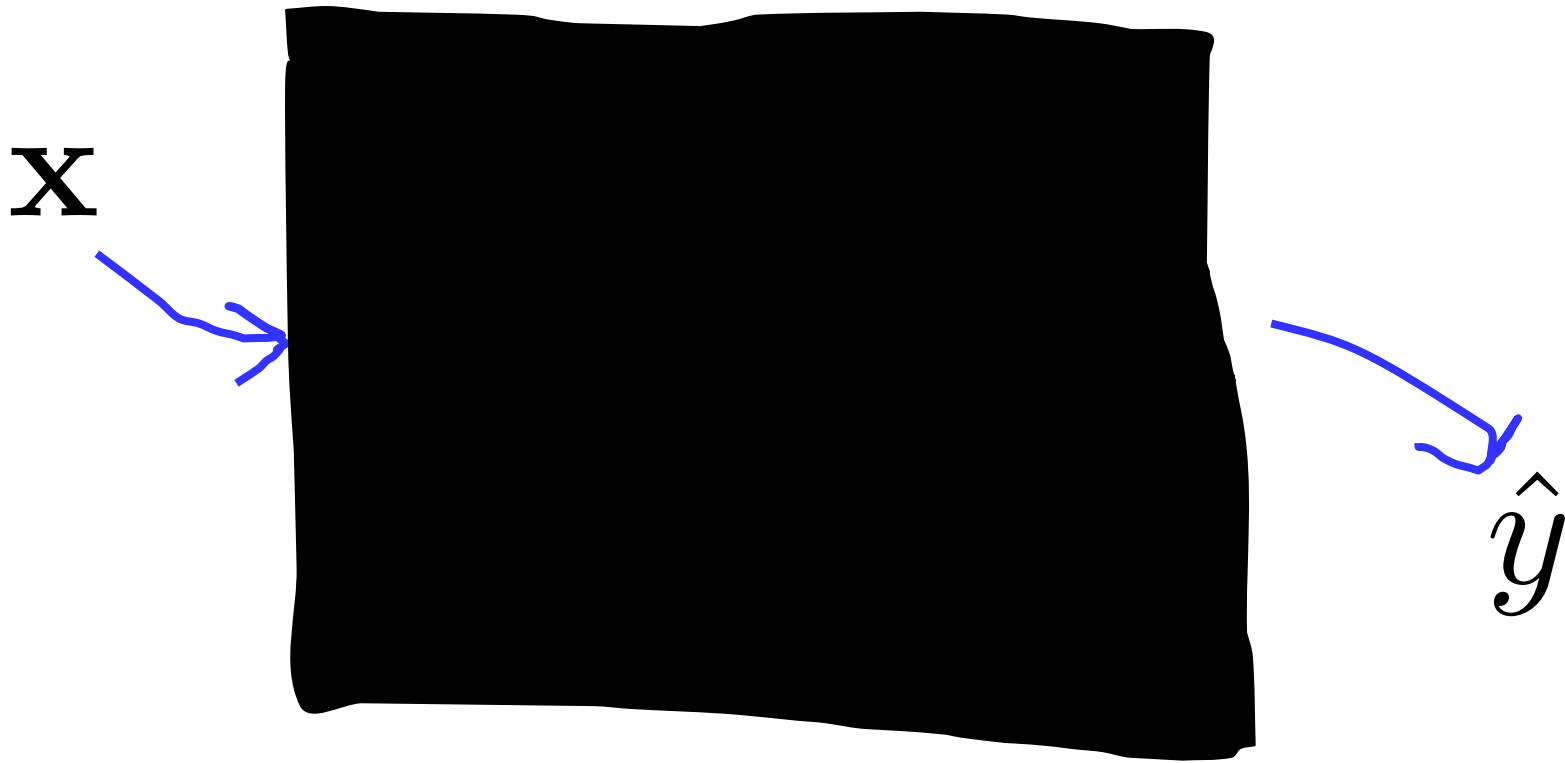
# How Does This Do?

- After training, can test with another set of data
  - “Testing” set also has known values for Y, so we can see how often we were right/wrong in predictions for Y
  - Spam data
    - Email data set: 1789 emails (1578 spam, 211 non-spam)
    - First, 1538 email messages (by time) used for training
    - Next 251 messages used to test learned classifier
  - Criteria:
    - Precision = # *correctly* predicted class Y / # predicted class Y
    - Recall = # *correctly* predicted class Y / # real class Y messages

	Spam		Non-spam	
	Precision	Recall	Precision	Recall
<b>Words only</b>	<b>97.1%</b>	<b>94.3%</b>	<b>87.7%</b>	<b>93.4%</b>
<b>Words + add'l features</b>	<b>100%</b>	<b>98.3%</b>	<b>96.2%</b>	<b>100%</b>

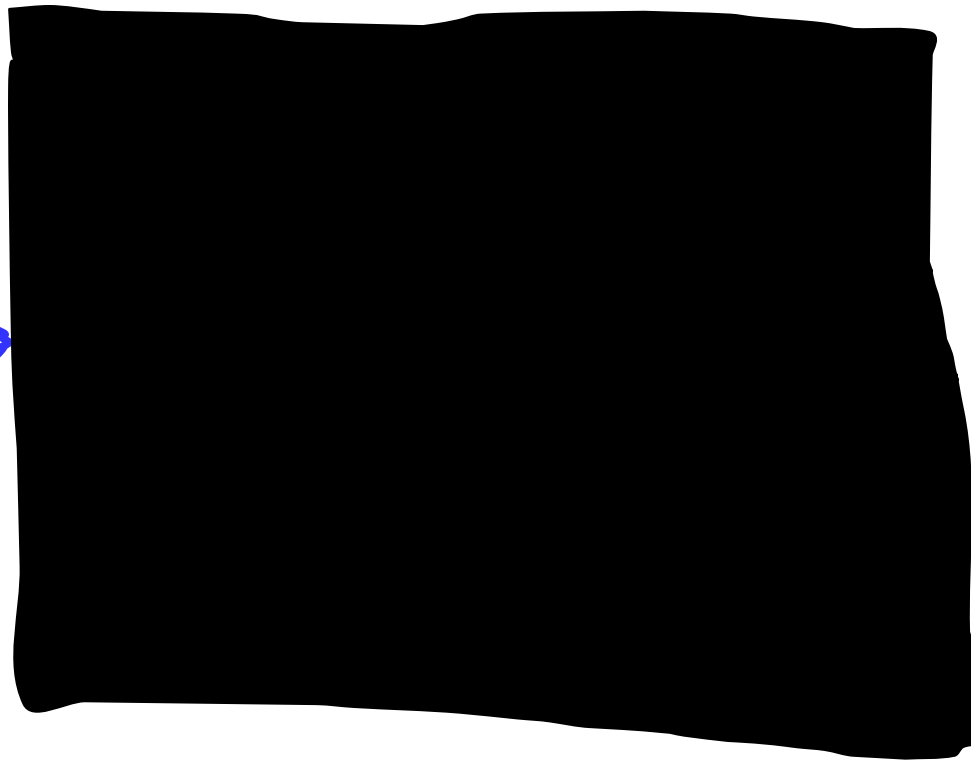
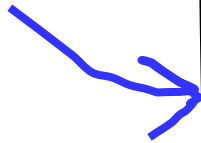
# Naïve Bayes Classification

$$g_{\theta}(\mathbf{x})?$$



$$g_{\theta}(\mathbf{x})?$$

$\mathbf{x}$   
[0, 1, 1, 0]

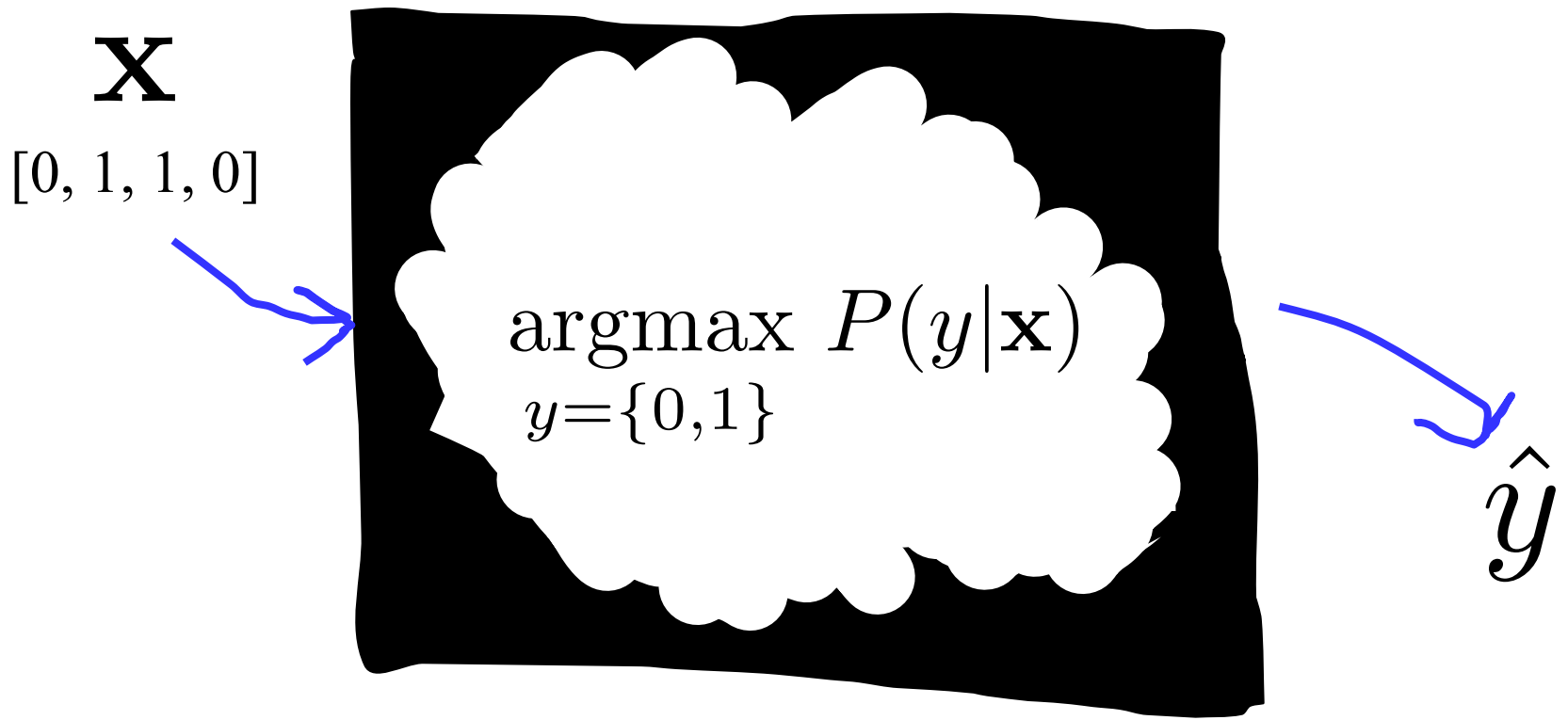


$\hat{y}$



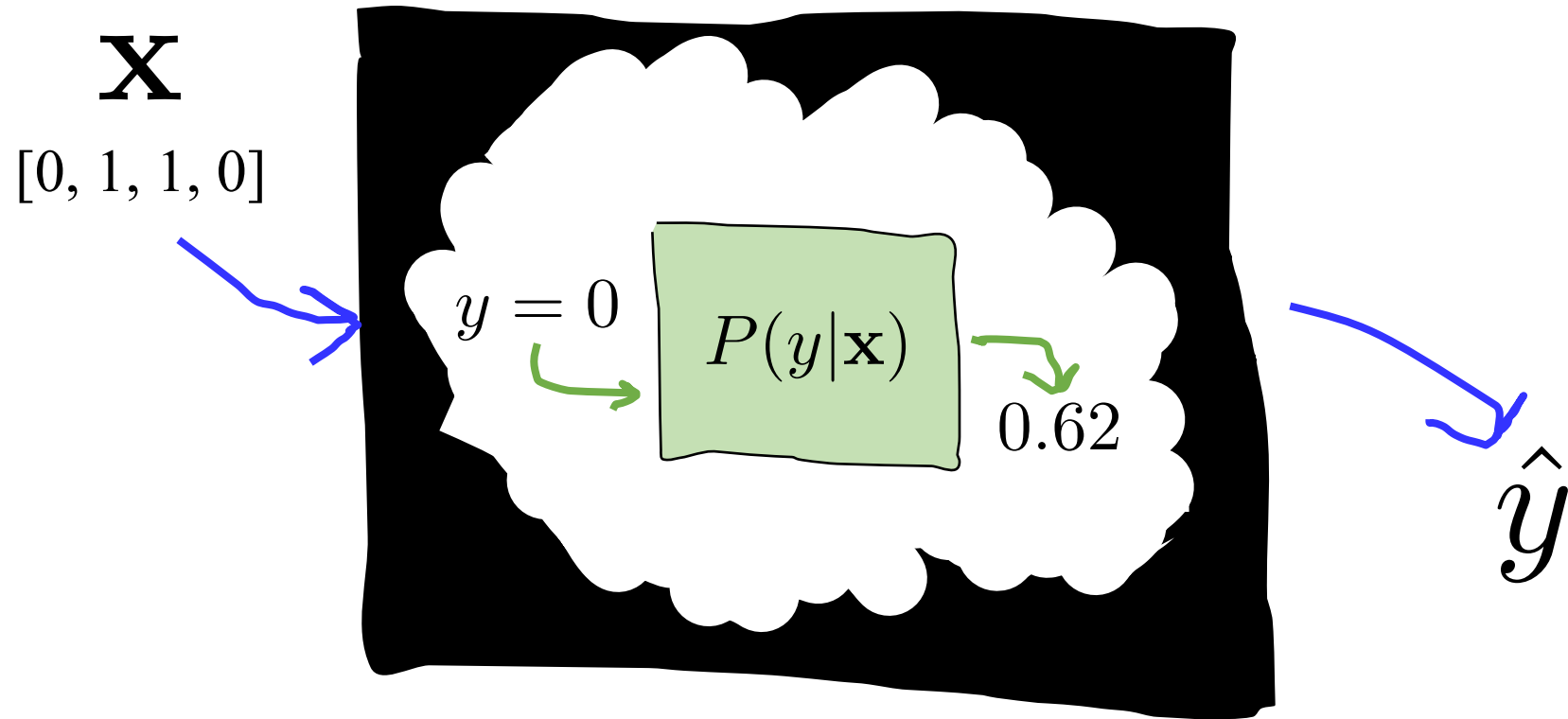
$$g_{\theta}(\mathbf{x})?$$

$\mathbf{x}$   
[0, 1, 1, 0]

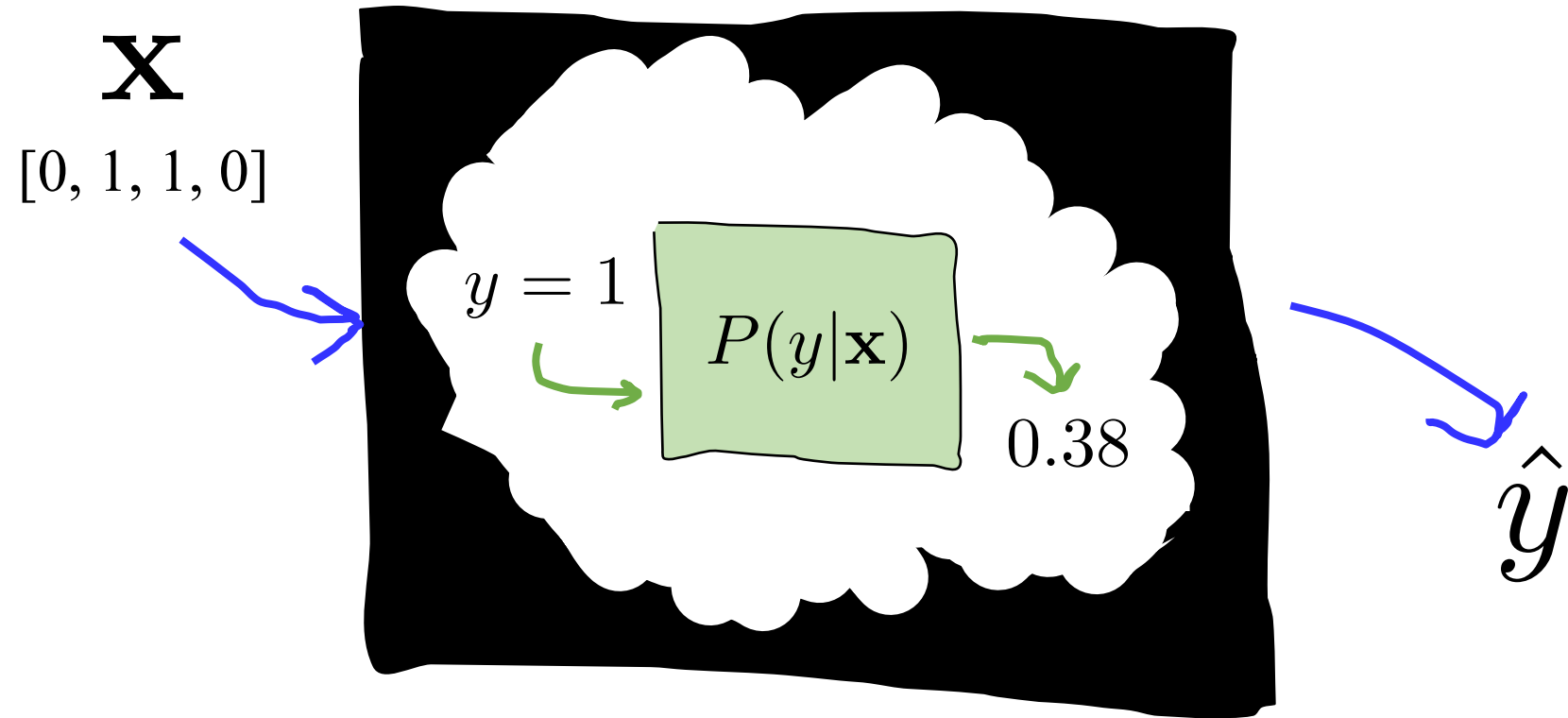

$$\operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x})$$

$\hat{y}$

$$g_{\theta}(\mathbf{x})?$$

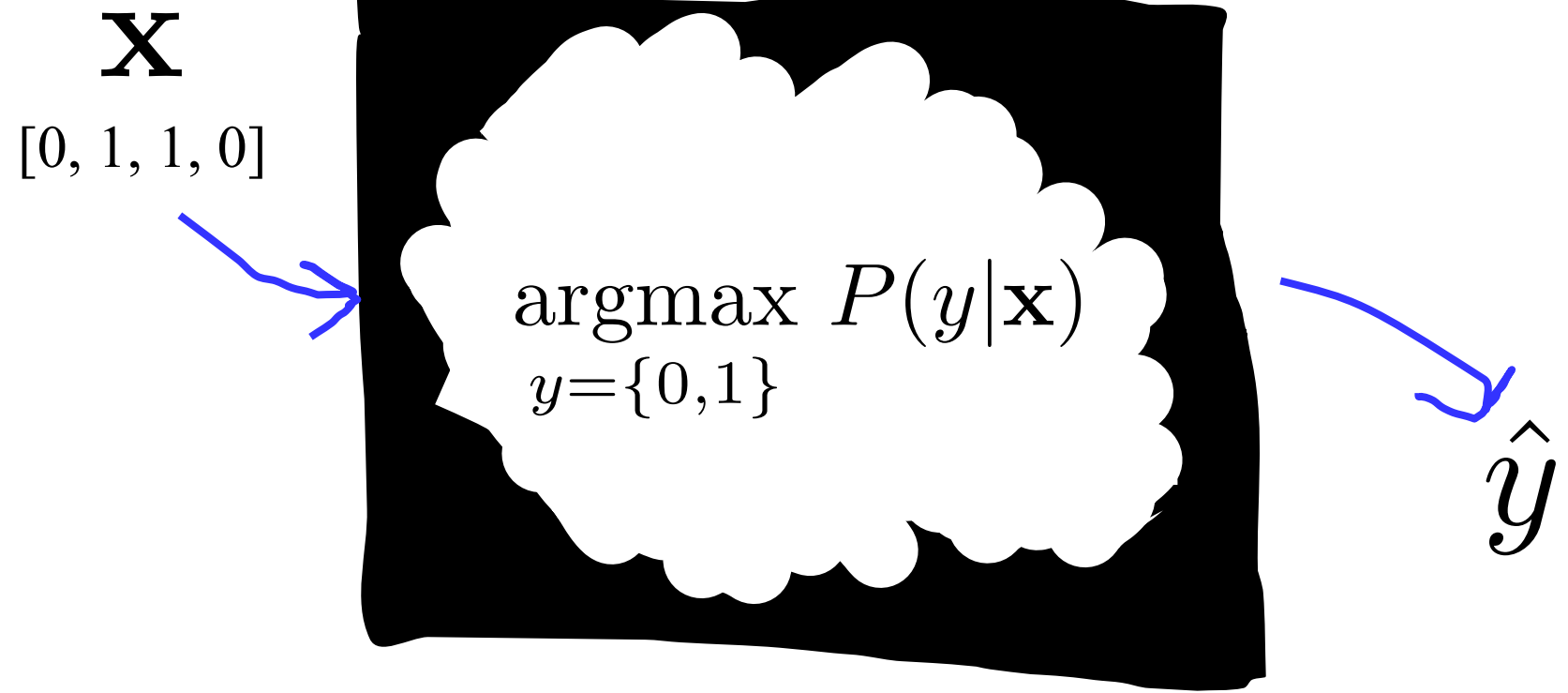


$$g_{\theta}(\mathbf{x})?$$



$$g_{\theta}(\mathbf{x})?$$

$\mathbf{x}$   
[0, 1, 1, 0]


$$\operatorname{argmax}_{y=\{0,1\}} P(y|\mathbf{x})$$

$\hat{y}$

$$g_{\theta}(\mathbf{x})?$$

$\mathbf{x}$   
[0, 1, 1, 0]

$\operatorname{argmax}_{y \in \{0, 1\}} P(y|\mathbf{x})$

$\hat{y} = 0$

# Naïve Bayes

Simply chose the class label that is the most likely given the data. Make Naïve Bayes assumption

$$\begin{aligned}\hat{y} &= g(\mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} P(Y = y | \mathbf{X} = \mathbf{x})\end{aligned}$$

# Naïve Bayes

Simply chose the class label that is the most likely given the data. Make Naïve Bayes assumption

$$\begin{aligned}\hat{y} &= g(\mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} P(Y = y | \mathbf{X} = \mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} \frac{P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y)}{P(\mathbf{X} = \mathbf{x})} \\ &= \arg \max_{y=\{0,1\}} P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y) \\ &= \arg \max_{y=\{0,1\}} P(Y = y) \prod_i P(X_i = x_i | Y = y) \\ &= \arg \max_{y=\{0,1\}} \log P(Y = y) + \sum_i \log P(X_i = x_i | Y = y)\end{aligned}$$

# Naïve Bayes

Simply chose the class label that is the most likely given the data. Make Naïve Bayes assumption

$$\begin{aligned}\hat{y} &= g(\mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} P(Y = y | \mathbf{X} = \mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} \frac{P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y)}{P(\mathbf{X} = \mathbf{x})} \\ &= \arg \max_{y=\{0,1\}} P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y) \\ &= \arg \max_{y=\{0,1\}} P(Y = y) \prod_i P(X_i = x_i | Y = y) \\ &= \arg \max_{y=\{0,1\}} \log P(Y = y) + \sum_i \log P(X_i = x_i | Y = y)\end{aligned}$$

Woot, Bayes!

# Naïve Bayes

Simply chose the class label that is the most likely given the data. Make Naïve Bayes assumption

$$\begin{aligned}\hat{y} &= g(\mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} P(Y = y | \mathbf{X} = \mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} \frac{P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y)}{P(\mathbf{X} = \mathbf{x})} \\ &= \arg \max_{y=\{0,1\}} P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y) \\ &= \arg \max_{y=\{0,1\}} P(Y = y) \prod_i P(X_i = x_i | Y = y) \\ &= \arg \max_{y=\{0,1\}} \log P(Y = y) + \sum_i \log P(X_i = x_i | Y = y)\end{aligned}$$

argmax is unaffected by P(X)

# Naïve Bayes

Simply chose the class label that is the most likely given the data. Make Naïve Bayes assumption

$$\begin{aligned}\hat{y} &= g(\mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} P(Y = y | \mathbf{X} = \mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} \frac{P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y)}{P(\mathbf{X} = \mathbf{x})} \\ &= \arg \max_{y=\{0,1\}} P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y) \\ &= \arg \max_{y=\{0,1\}} P(Y = y) \prod_i P(X_i = x_i | Y = y) \\ &= \arg \max_{y=\{0,1\}} \log P(Y = y) + \sum_i \log P(X_i = x_i | Y = y)\end{aligned}$$

Naïve Bayes Assumption

# Naïve Bayes

Simply chose the class label that is the most likely given the data. Make Naïve Bayes assumption

$$\begin{aligned}\hat{y} &= g(\mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} P(Y = y | \mathbf{X} = \mathbf{x}) \\ &= \arg \max_{y=\{0,1\}} \frac{P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y)}{P(\mathbf{X} = \mathbf{x})} \\ &= \arg \max_{y=\{0,1\}} P(Y = y) P(\mathbf{X} = \mathbf{x} | Y = y) \\ &= \arg \max_{y=\{0,1\}} P(Y = y) \prod_i P(X_i = x_i | Y = y) \\ &= \arg \max_{y=\{0,1\}} \log P(Y = y) + \sum_i \log P(X_i = x_i | Y = y)\end{aligned}$$



Argmax of log

# Computing Probabilities from Data

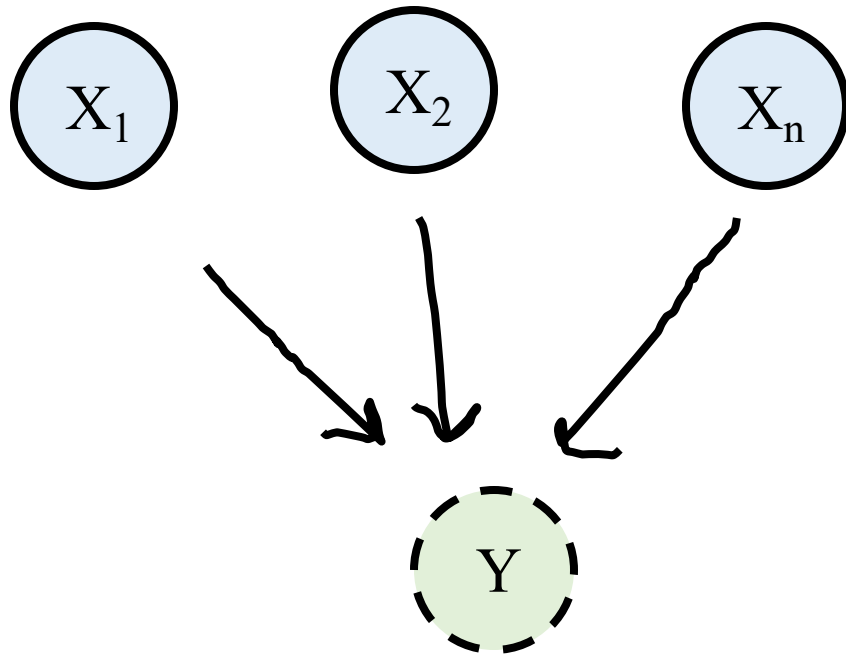
- Various probabilities you will need to compute for Naive Bayesian Classifier (using MLE here):

$$\hat{p}(X_i = 1|Y = 0) = \frac{(\# \text{ training examples where } X_i = 1 \text{ and } Y = 0)}{(\# \text{ training examples where } Y = 0)}$$

$$\hat{p}(Y = 1) = \frac{(\# \text{ training examples where } Y = 1)}{(\# \text{ training examples})}$$

Deeper Understanding

# Brute Force Bayes Bayes Net

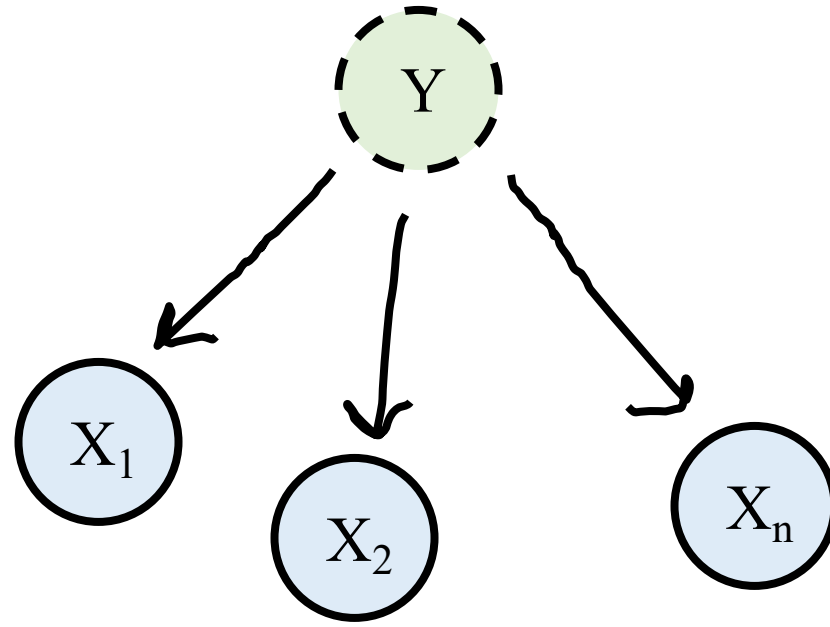


Parameters:

$$P(Y = y | \text{Parents of } Y \text{ take on specified values})$$

$$P(X_i = x_i)$$

# Naïve Bayes Bayes Net



Assumption:

$$P(\mathbf{x}, y) = P(y) \prod_i P(x_i|y)$$

Parameters:

$$P(X_i = x_i | \text{Parents of } X_i \text{ take on specified values})$$

$$P(Y = y)$$