



#### 10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

## Naïve Bayes



# Generative vs. Discriminative

Matt Gormley Lecture 18 Mar. 23, 2020

#### Reminders

- Homework 6: Learning Theory / Generative Models
  - Out: Fri, Mar 20
  - Due: Fri, Mar 27 at 11:59pm
- Midterm Exam 2
  - Thu, Apr 2 evening exam, details announced on Piazza
- Today's In-Class Poll
  - http://poll.mlcourse.org

#### Q&A

**Q:** Why would we use Naïve Bayes? Isn't it too Naïve?

**A:** Naïve Bayes has one **key advantage** over methods like Perceptron, Logistic Regression, Neural Nets:

#### Training is lightning fast!

While other methods require slow iterative training procedures that might require hundreds of epochs, Naïve Bayes computes its parameters in closed form by counting.

## NAÏVE BAYES

Flip weighted coin



 $\chi_3$ 

 $x_M$ 

If HEADS, flip each red coin



 $\chi_2$ 

 $\mathcal{Y}$ 

 $x_1$ 

If TAILS, flip each blue coin



We can **generate** data in this fashion. Though in practice we never would since our data is **given**.

Instead, this provides an explanation of **how** the data was generated (albeit a terrible one).

Each red coin corresponds to  $an x_m$ 

# What's wrong with the Naïve Bayes Assumption?

#### The features might not be independent!!

- Example 1:
  - If a document contains the word "Donald", it's extremely likely to contain the word "Trump"
  - These are not independent!

\* ELECTION 2016 \* MORE ELECTION COVERAGE

Trump Spends Entire Classified National Security Briefing Asking About Egyptian Mummies



NEWS IN BRIEF August 18, 2016 VOL 52 ISSUE 32 - Politics - Politicians - Election 2016 - Donald Trum

#### • Example 2:

If the petal width is very high,
 the petal length is also likely to
 be very high



#### Recipe for Closed-form MLE

1. Assume data was generated i.i.d. from some model (i.e. write the generative story)  $x^{(i)} \sim p(x|\theta)$ 

2. Write log-likelihood

$$\ell(\boldsymbol{\theta}) = \log p(\mathbf{x}^{(1)}|\boldsymbol{\theta}) + \dots + \log p(\mathbf{x}^{(N)}|\boldsymbol{\theta})$$

3. Compute partial derivatives (i.e. gradient)

$$\frac{\partial \ell(\mathbf{\Theta})}{\partial \theta_1} = \dots$$
$$\frac{\partial \ell(\mathbf{\Theta})}{\partial \theta_2} = \dots$$
$$\frac{\partial \ell(\mathbf{\Theta})}{\partial \theta_M} = \dots$$

4. Set derivatives to zero and solve for  $\theta$ 

$$\partial \ell(\theta)/\partial \theta_{\rm m} = 0$$
 for all m  $\in \{1, ..., M\}$   
 $\theta^{\rm MLE} =$  solution to system of M equations and M variables

5. Compute the second derivative and check that  $\ell(\theta)$  is concave down at  $\theta^{\text{MLE}}$ 

## Naïve Bayes: Learning from Data

#### Whiteboard

- Data likelihood
- MLE for Naive Bayes
- Example: MLE for Naïve Bayes with Two Features
- MAP for Naive Bayes

## NAÏVE BAYES: MODEL DETAILS

Data: Binary feature vectors, Binary labels

$$\mathbf{x} \in \{0,1\}^M$$

#### **Generative Story:**

$$y \sim \mathsf{Bernoulli}(\phi)$$

$$x_1 \sim \mathsf{Bernoulli}(\theta_{y,1})$$

$$x_2 \sim \mathsf{Bernoulli}(\theta_{y,2})$$

:

 $x_M \sim \mathsf{Bernoulli}(\theta_{y,M})$ 

$$y \in \{0, 1\}$$

#### Model:

$$p_{\phi,\theta}(x,y) = p_{\phi,\theta}(x_1, \dots, x_M, y)$$

$$= p_{\phi}(y) \prod_{m=1}^{M} p_{\theta}(x_m | y)$$

$$= \left[ (\phi)^y (1 - \phi)^{(1-y)} \right]$$

$$\prod_{m=1}^{M} (\theta_{y,m})^{x_m} (1 - \theta_{y,m})^{(1-x_m)}$$

#### **Maximum Likelihood Estimation**

#### Training: Find the class-conditional MLE parameters

Count 
$$N_{y=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 1)$$

$$N_{y=0} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0)$$

$$N_{y=0,x_m=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0 \land x_m^{(i)} = 1)$$

Maximum Likelihood **Estimators:** 

$$\phi = \frac{N_{y=1}}{N}$$

$$\theta_{0,m} = \frac{N_{y=0,x_m=1}}{N_{y=0}}$$

$$\theta_{1,m} = \frac{N_{y=1,x_m=1}}{N_{y=1}}$$

$$\forall m \in \{1, \dots, M\}$$

#### **Maximum Likelihood Estimation**

#### Training: Find the class-conditional MLE parameters

Count 
$$N_{y=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 1)$$

$$N_{y=0} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0)$$

$$N_{y=0,x_m=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0 \land x_m^{(i)} = 1)$$

Maximum Likelihood **Estimators:** 

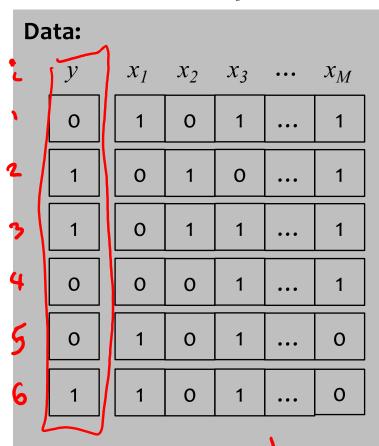
$$\phi = \frac{N_{y=1}}{N}$$

$$\theta_{0,m} = \frac{N_{y=0,x_m=1}}{N_{y=0}}$$

$$\theta_{0,m} = \frac{N_{y=0,x_m=1}}{N_{y=0}}$$

$$\theta_{1,m} = \frac{N_{y=1,x_m=1}}{N_{y=1}}$$

$$\forall m \in \{1, \dots, M\}$$



## Question 1:

What is the MLE of  $\phi$ ?

(A) 0/6 (B) 1/6 (C) 2/6 (D) 3/6

(E) 4/6 (F) 5/6 (G) 6/6 (H) None of the above

#### **Maximum Likelihood Estimation**

#### Training: Find the class-conditional MLE parameters

Count 
$$N_{y=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 1)$$

$$N_{y=0} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0)$$

$$N_{y=0,x_m=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0 \land x_m^{(i)} = 1)$$

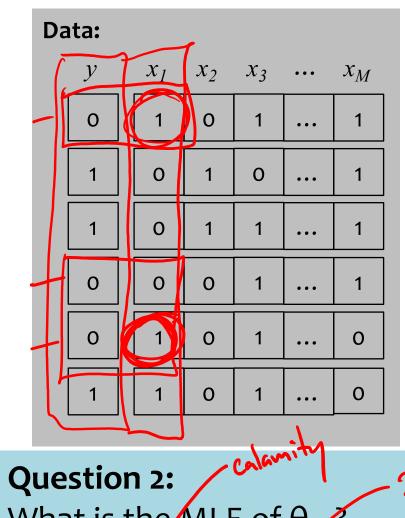
Maximum Likelihood Estimators:

$$\phi = \frac{N_{y=1}}{N}$$

$$\theta_{0,m} \neq \frac{N_{y=0,x_m=1}}{N_{y=0,x_m=1}}$$

$$\theta_{1,m} = \frac{N_{y=1,x_m=1}}{N_{y=1}}$$

$$\forall m \in \{1,\dots,M\}$$



What is the MLE of  $\theta_{\alpha}$ ?

(A) 0/6 (B) 1/6 (C) 2/6 (D) 3/6

E) 4/6 (F) 5/6 (G) 6/6 (H) None of the above

#### **Maximum Likelihood Estimation**

Training: Find the class-conditional MLE parameters

Count 
$$N_{y=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 1)$$

$$N_{y=0} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0)$$

$$N_{y=0,x_m=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0 \land x_m^{(i)} = 1)$$

Maximum Likelihood **Estimators:** 

$$\phi = \frac{N_{y=1}}{N}$$

$$\theta_{0,m} = \frac{N_{y=0,x_m=1}}{N_{y=0}}$$

$$\theta_{1,m} = \frac{N_{y=1,x_m=1}}{N_{y=1}}$$

$$\forall m \in \{1, \dots, M\}$$



MLE for Naïve Bayes is a splendid learning algorithm for when you have say billions of training examples and hundreds of millions of features!

You only need one pass through the data to perform some counting.

#### MLE

What does maximizing likelihood accomplish?

- There is only a finite amount of probability mass (i.e. sum-to-one constraint)
- MLE tries to allocate as much probability mass as possible to the things we have observed...

... at the expense of the things we have not observed

## A Shortcoming of MLE

For Naïve Bayes, suppose we **never** observe the word "unicorn" in a real news article.

In this case, what is the MLE of the following quantity?

$$p(x_{unicorn}|y=real) = \bigcirc$$

Recall: 
$$\theta_{k,0} = \frac{\sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0 \land x_k^{(i)} = 1)}{\sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0)}$$

Now suppose we observe the word "unicorn" at test time. What is the posterior probability that the article was a real article?

article? 
$$p(y \neq real | \mathbf{x}) = \frac{p(\mathbf{x}|y \neq real)p(y \neq real)}{p(\mathbf{x})}$$

## Recipe for Closed-form MAP **Estimation**

- Assume data was generated i.i.d. from some model 1. (i.e. write the generative story)
  - $\theta \sim p(\theta)$  and then for all i:  $x^{(i)} \sim p(x|\theta)$
- Write log-likelinood 2.
  - $\ell_{MAP}(\theta) = \log p(\theta) + \log p(x^{(1)}|\theta) + ... + \log p(x^{(N)}|\theta)$ Compute partial derivatives (i.e. gradient)
- 3.

$$\partial \ell_{MAP}(\mathbf{\theta})/\partial \theta_1 = \dots$$

$$\partial \ell_{\mathsf{MAP}}(\mathbf{\theta})/\partial \theta_2 = \dots$$

$$\partial \ell_{MAP}(\boldsymbol{\Theta})/\partial \boldsymbol{\Theta}_{M} = \dots$$

Set derivatives to zero and solve for  $\boldsymbol{\theta}$ 4.

$$\partial \ell_{MAP}(\boldsymbol{\theta})/\partial \theta_{m} = o \text{ for all } m \in \{1, ..., M\}$$

 $\Theta^{MAP}$  = solution to system of M equations and M variables

Compute the second derivative and check that  $\ell(\theta)$  is concave down 5. at  $\theta^{MAP}$ 

## & ~ Uniform ([0,1])

## Model 1: Bernoulli Naïve Bayes

#### MAP Estimation (Beta Prior)

#### 1. Generative Story:

The parameters are drawn once for the entire dataset.

for 
$$m \in \{1, \dots, M\}$$
: for  $y \in \{0, 1\}$ :

$$iggraphi$$
 Beta $(lpha,eta)$ 

for 
$$i \in \{1, \dots, N\}$$
:

$$y^{(i)} \sim \text{Bernoulli}(\phi)$$

for 
$$m \in \{1, ..., M\}$$
:

$$x_m^{(i)} \sim \mathsf{Bernoulli}(\theta_{y^{(i)},m})$$

$$N_{y=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 1)$$

$$N_{y=0} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0)$$

$$N_{y=0,x_m=1} = \sum_{i=1}^{N} \mathbb{I}(y^{(i)} = 0 \land x_m^{(i)} = 1)$$

#### 2. Likelihood:

$$\ell_{MAP}(\phi, \theta) = \lim_{\theta \to 0} \int_{\Omega(\theta, \theta)} \int_$$

$$= \log[p(\phi, \boldsymbol{\theta}|\alpha, \beta)p(\mathcal{D}|\phi, \boldsymbol{\theta})]$$

$$= \log \left[ \left( p(\phi | \alpha, \beta) \prod_{m=1}^{M} p(\theta_{0,m} | \alpha, \beta) \right) \left( \prod_{i=1}^{N} p(\mathbf{x}^{(i)}, y^{(i)} | \phi, \boldsymbol{\theta}) \right) \right]$$

3. MAP Estimators: 
$$(\phi^{MAP}, \boldsymbol{\theta}^{MAP}) = \operatorname*{argmax}_{\phi, \boldsymbol{\theta}} \ell_{MAP}(\phi, \boldsymbol{\theta})$$

Take derivatives, set to zero and solve...

$$\phi = \frac{N_{y=1}}{N}$$

$$\theta_{0,m} = (\alpha - 1) + N_{y=0,x_m=1}$$

$$(\alpha - 1) + (\beta - 1) + N_{y=0}$$

$$\theta_{1,m} = \frac{(\alpha - 1) + N_{y=1,x_m=1}}{(\alpha - 1) + (\beta - 1) + N_{y=1}}$$

$$\forall m \in \{1, \dots, M\}$$

#### Other NB Models

- Bernoulli Naïve Bayes:
  - for binary features
- 2. Multinomial Naïve Bayes:
  - for integer features
- 3. Gaussian Naïve Bayes:
  - for continuous features
- 4. Multi-class Naïve Bayes:
  - for classification problems with > 2 classes
  - event model could be any of Bernoulli, Gaussian, Multinomial, depending on features

event model

## Model 2: Multinomial Naïve Bayes

# Support: Option 1: Integer vector (word IDs) $\mathbf{x} = [x_1, x_2, \dots, x_M] \text{ where } x_m \in \{1, \dots, K\} \text{ a word id.}$

#### **Generative Story:**

$$for \ i \in \{1, \dots, N\};$$
 
$$y^{(i)} \sim \operatorname{Bernoulli}(\phi) + \mathbf{for} \ j \in \{1, \dots, M_i\};$$
 
$$x_j^{(i)} \sim \operatorname{Multinomial}(\theta_{y^{(i)}}, 1)$$

#### Model:

$$p_{\phi,\theta}(\boldsymbol{x},y) = p_{\phi}(y) \prod_{k=1}^{K} p_{\theta_k}(x_k|y)$$
$$= (\phi)^y (1-\phi)^{(1-y)} \prod_{j=1}^{M_i} \theta_{y,x_j}$$

## Model 3: Gaussian Naïve Bayes

#### **Support:**

$$\mathbf{x} \in \mathbb{R}^{\mathsf{MM}}$$

Model: Product of prior and the event model

$$p(\boldsymbol{x}, y) = p(x_1, \dots, x_{m}, y)$$

$$= p(y) \prod_{k=1}^{m} p(x_k|y)$$
Gaussin

Gaussian Naive Bayes assumes that  $p(x_k|y)$  is given by a Normal distribution.

## Model 4: Multiclass Naïve Bayes

#### Model:

The only change is that we permit y to range over C classes.

$$p(\mathbf{x}, y) = p(x_1, \dots, x_K, y)$$
  
=  $p(y) \prod_{k=1}^{K} p(x_k|y)$ 

Now,  $y \sim \text{Multinomial}(\phi, 1)$  and we have a separate conditional distribution  $p(x_k|y)$  for each of the C classes.

# Generic Naïve Bayes Model South South Louise Louise

**Support:** Depends on the choice of event model  $P(X_k|Y)$ 

Model: Product of prior and the event model,

$$P(\mathbf{X}, Y) = P(Y) \prod_{k=1}^{K} P(X_k | Y)$$

Training: Find the class-conditional MLE parameters

For P(Y), we find the MLE using all the data. For each  $P(X_k|Y)$  we condition on the data with the corresponding

Classification: Find the class that maximizes the posterior

$$\hat{y} = \operatorname*{argmax}_{y} p(y|\mathbf{x})$$

## Generic Naïve Bayes Model

## **Classification:** $\hat{y} = \operatorname{argmax} p(y|\mathbf{x})$ (posterior) $= \operatorname{argmax} \frac{p(\mathbf{x}|y)p(y)}{p(x)}$ (by Bayes' rule) $= \operatorname{argmax} p(\mathbf{x}|y)p(y)$

## VISUALIZING GAUSSIAN NAÏVE BAYES

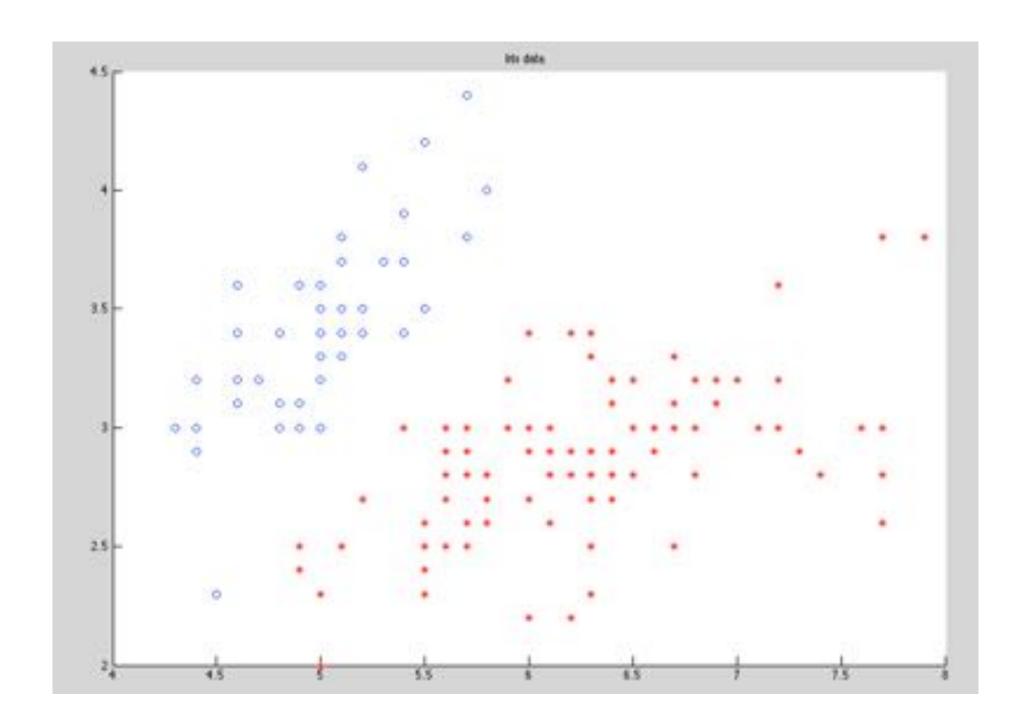


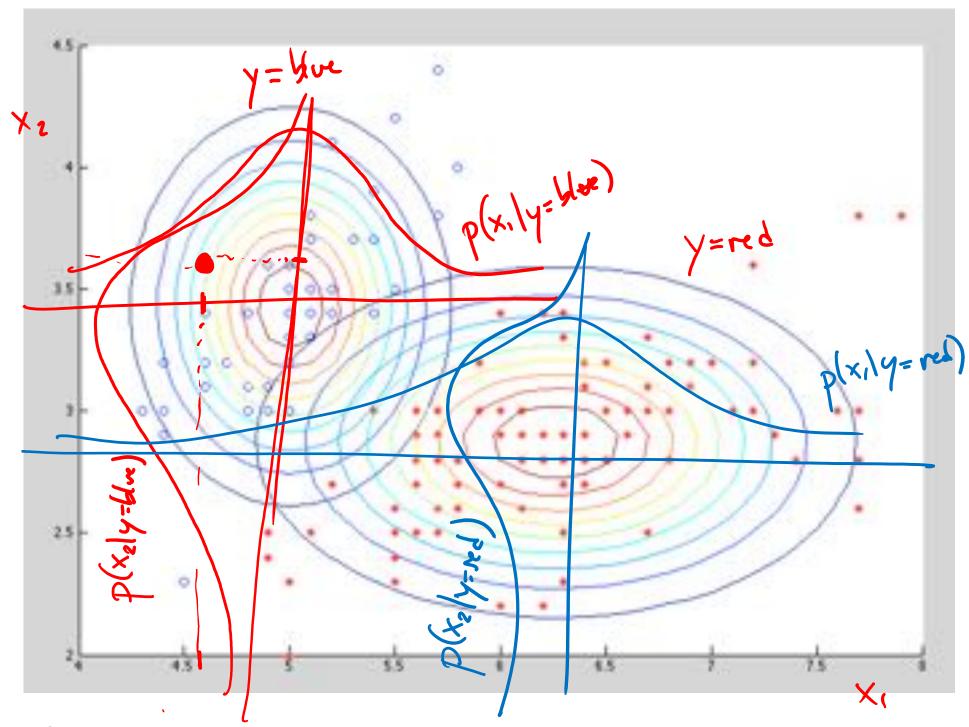


#### Fisher Iris Dataset

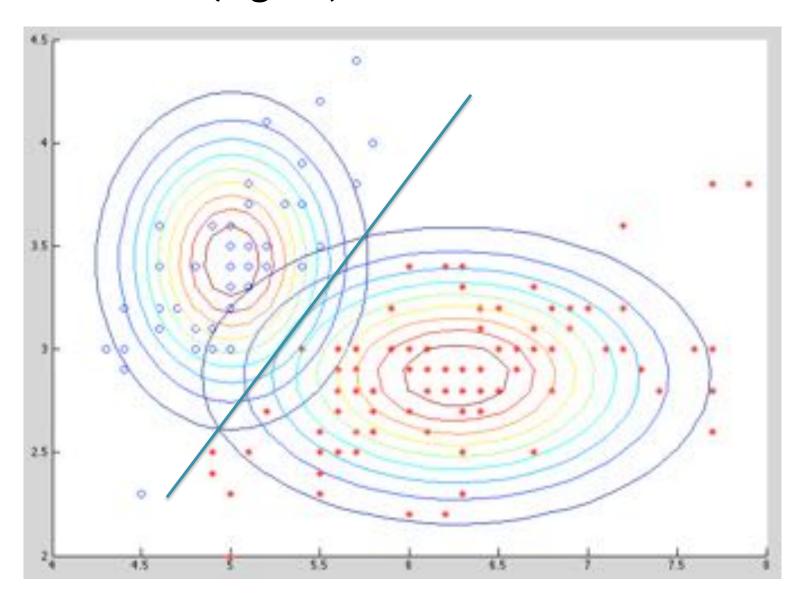
Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
1	6.7	3.0	5.0	1.7

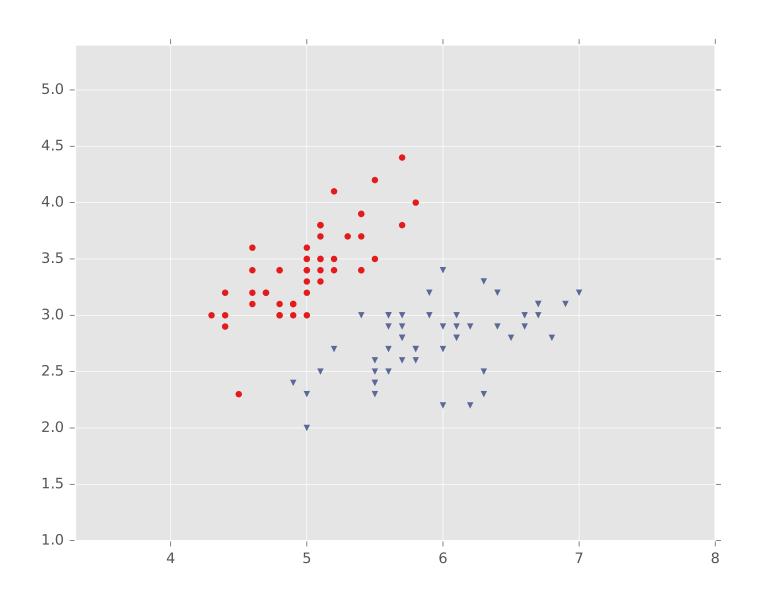




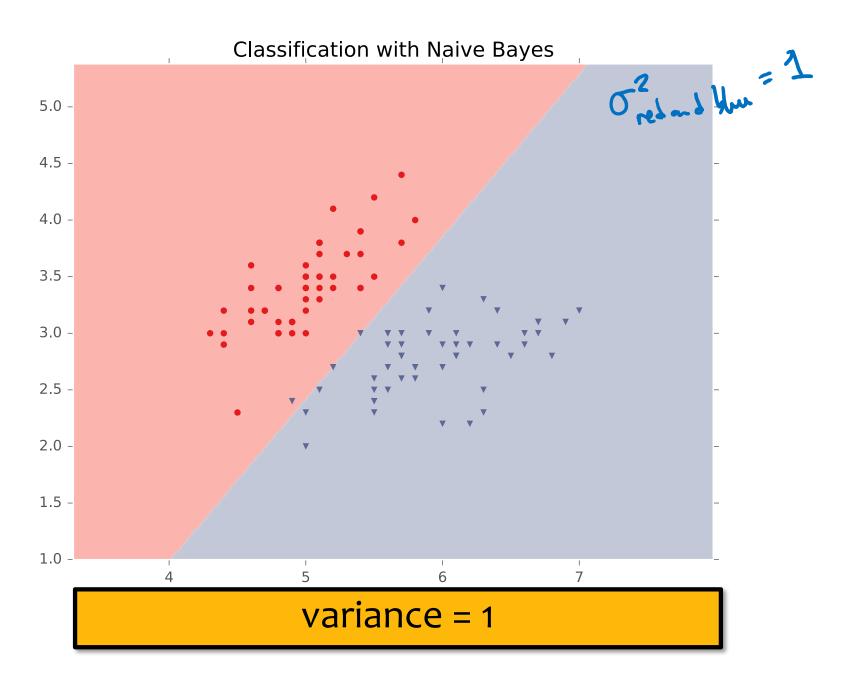
## Naïve Bayes has a **linear** decision boundary if variance (sigma) is constant across classes



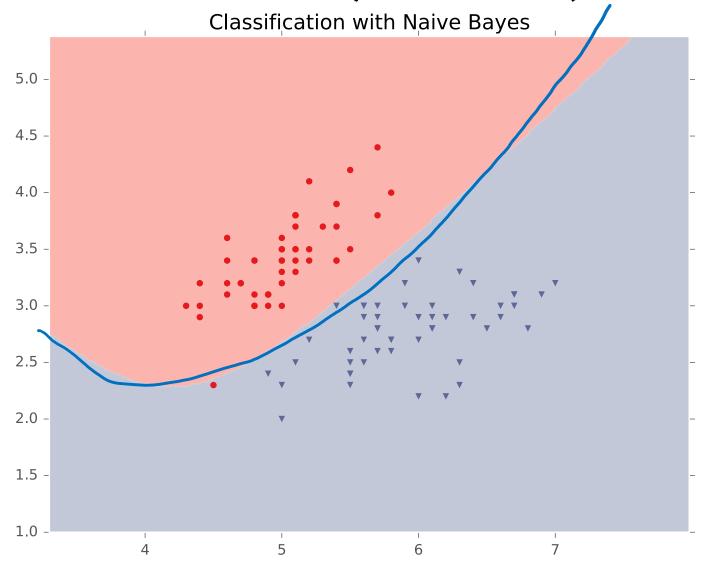
## Iris Data (2 classes)



## Iris Data (2 classes)

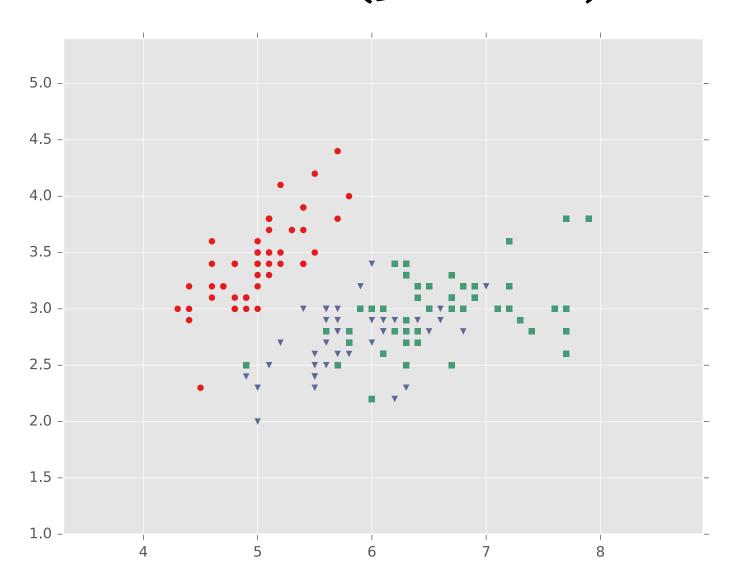


## Iris Data (2 classes)

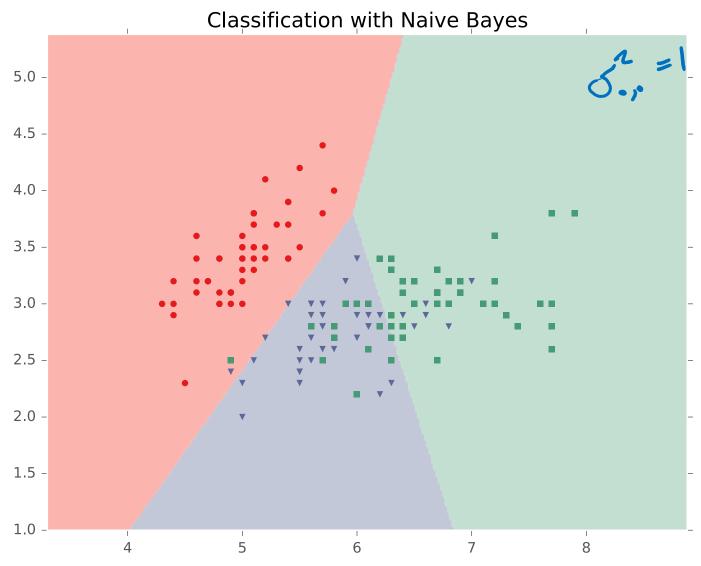


variance learned for each class

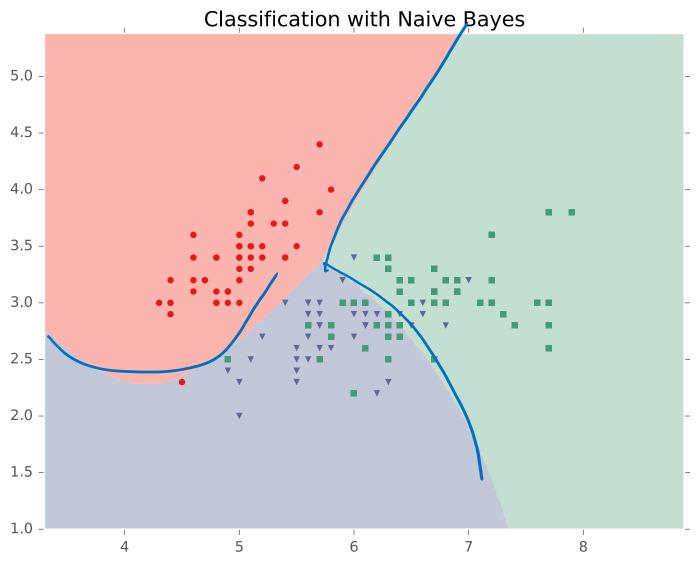
## Iris Data (3 classes)



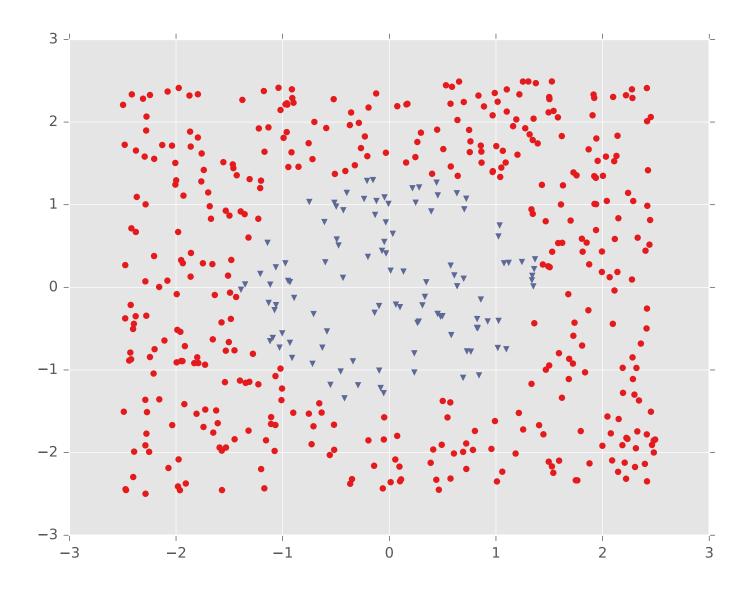
## Iris Data (3 classes)



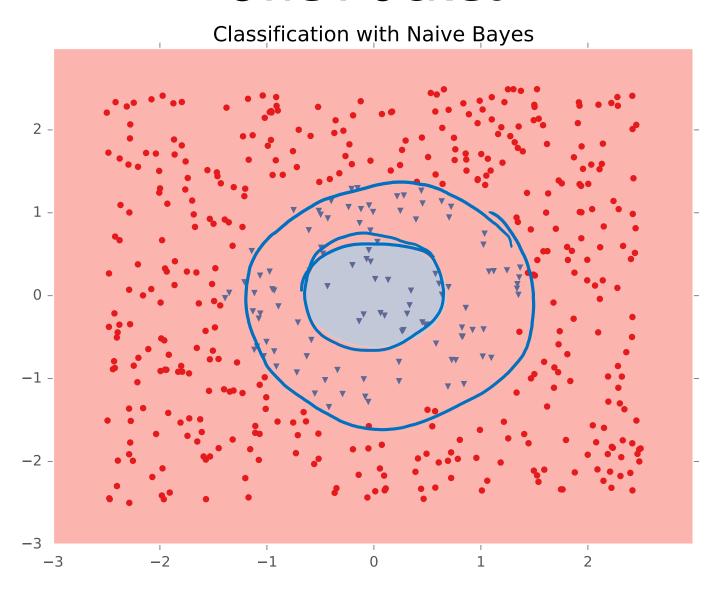
## Iris Data (3 classes)



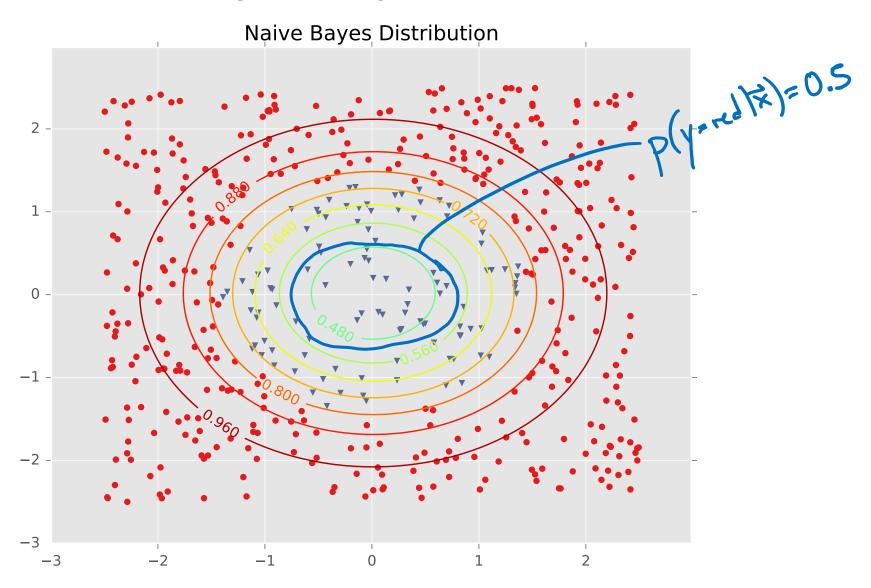
# One Pocket



## One Pocket



## One Pocket



## Summary

- Naïve Bayes provides a framework for generative modeling
- 2. Choose  $p(x_m | y)$  appropriate to the data (e.g. Bernoulli for binary features, Gaussian for continuous features)
- 3. Train by MLE or MAP
- 4. Classify by maximizing the posterior

## Learning Objectives

#### **Naïve Bayes**

#### You should be able to...

- 1. Write the generative story for Naive Bayes
- 2. Create a new Naive Bayes classifier using your favorite probability distribution as the event model
- 3. Apply the principle of maximum likelihood estimation (MLE) to learn the parameters of Bernoulli Naive Bayes
- 4. Motivate the need for MAP estimation through the deficiencies of MLE
- 5. Apply the principle of maximum a posteriori (MAP) estimation to learn the parameters of Bernoulli Naive Bayes
- 6. Select a suitable prior for a model parameter
- 7. Describe the tradeoffs of generative vs. discriminative models
- 8. Implement Bernoulli Naives Bayes
- 9. Employ the method of Lagrange multipliers to find the MLE parameters of Multinomial Naive Bayes
- 10. Describe how the variance affects whether a Gaussian Naive Bayes model will have a linear or nonlinear decision boundary

# DISCRIMINATIVE AND GENERATIVE CLASSIFIERS

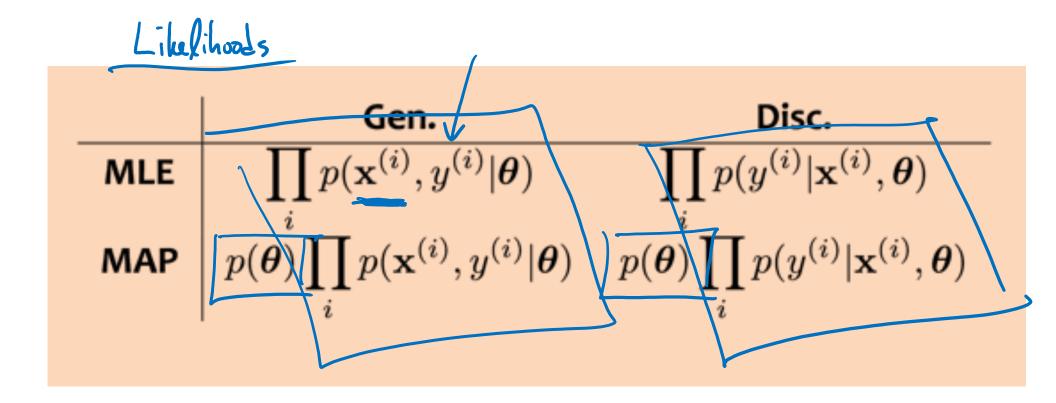
#### Generative Classifiers:

- Example: Naïve Bayes
- Define a joint model of the observations  ${\bf x}$  and the labels  ${\bf y}$ :  $p({\bf x},y)$
- Learning maximizes (joint) likelihood
- Use Bayes' Rule to classify based on the posterior:

$$p(y|\mathbf{x}) = p(\mathbf{x}|y)p(y)/p(\mathbf{x})$$

### Discriminative Classifiers:

- Example: Logistic Regression
- Directly model the conditional  $p(y|\mathbf{x})$
- Learning maximizes conditional likelihood



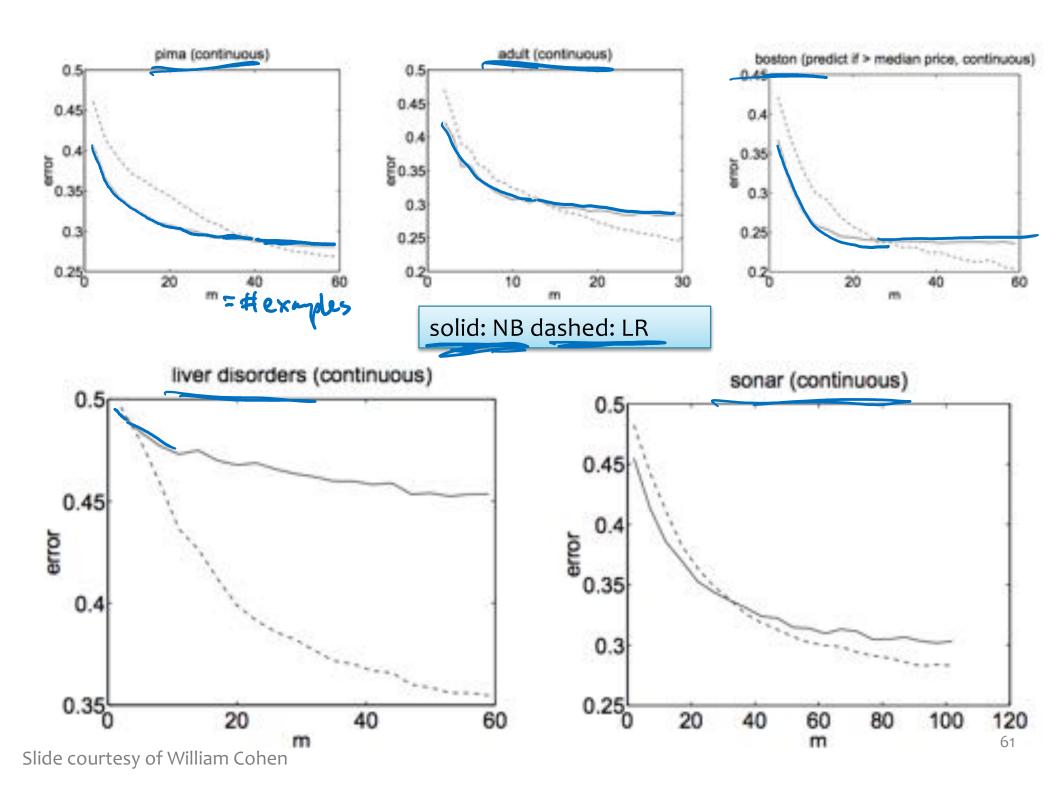
Finite Sample Analysis (Ng & Jordan, 2002)

[Assume that we are learning from a finite training dataset]

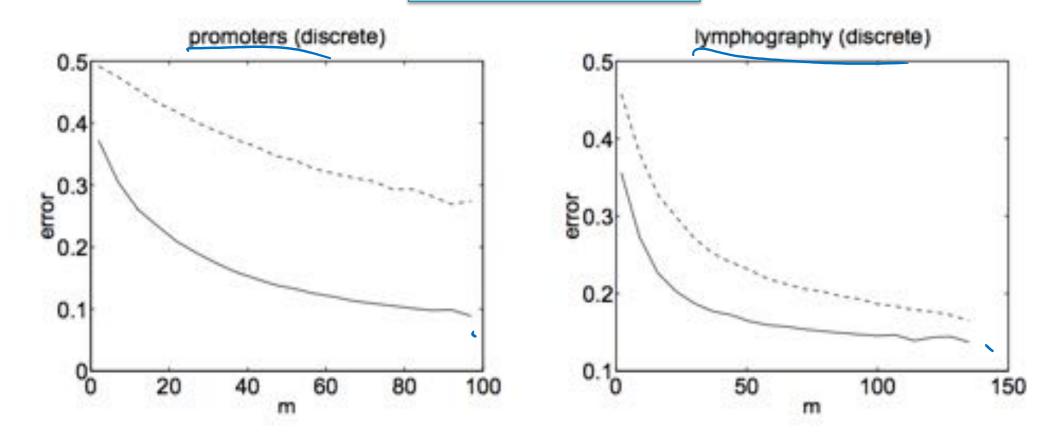
If model assumptions are correct: Naive Bayes is a more efficient learner (requires fewer samples) than Logistic Regression

If model assumptions are incorrect: Logistic Regression has lower asymtotic error, and does better than Naïve Bayes

Gyuldin efficiency



solid: NB dashed: LR



Naïve Bayes makes stronger assumptions about the data but needs fewer examples to estimate the parameters

"On Discriminative vs Generative Classifiers: ...." Andrew Ng and Michael Jordan, NIPS 2001.

## Learning (Parameter Estimation)

### **Naïve Bayes:**

Parameters are decoupled -> Closed form solution for MLE

### **Logistic Regression:**

Parameters are coupled  $\rightarrow$  No closed form solution – must use iterative optimization techniques instead

# Naïve Bayes vs. Logistic Reg.

## Learning (MAP Estimation of Parameters)

### **Bernoulli Naïve Bayes:**

Parameters are probabilities > Beta prior (usually) pushes probabilities away from zero / one extremes

### **Logistic Regression:**

Parameters are not probabilities 

Gaussian prior encourages parameters to be close to zero

(effectively pushes the probabilities away from zero / one extremes)

## Naïve Bayes vs. Logistic Reg.

#### **Features**

#### **Naïve Bayes:**

Features x are assumed to be conditionally independent given y. (i.e. Naïve Bayes Assumption)

### **Logistic Regression:**

No assumptions are made about the form of the features x. They can be dependent and correlated in any fashion.

# MOTIVATION: STRUCTURED PREDICTION

## Structured Prediction

 Most of the models we've seen so far were for classification

– Given observations: 
$$\mathbf{x} = (x_1, x_2, ..., x_K)$$

- Predict a (binary) label
- Many real-world problems require structured prediction

– Given observations: 
$$\mathbf{x} = (x_1, x_2, ..., x_K)$$

- Predict a structure:  $y \ni (y_1, y_2, ..., y_J)$
- Some classification problems benefit from latent structure

## Structured Prediction Examples

## Examples of structured prediction

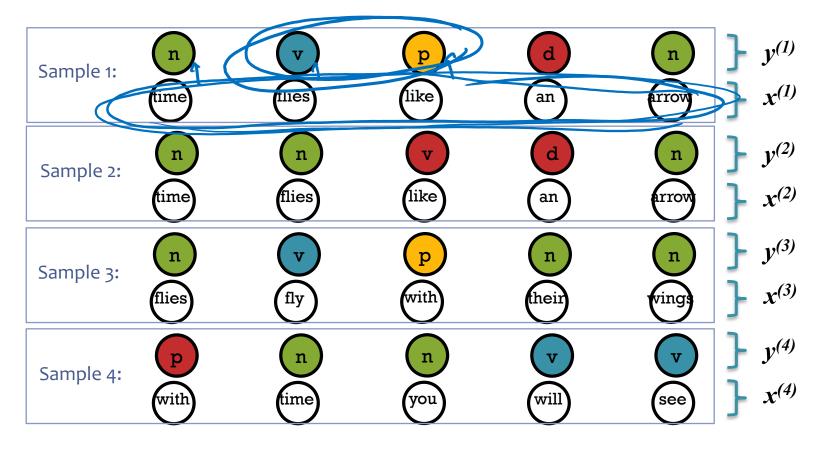
- Part-of-speech (POS) tagging
- Handwriting recognition
- Speech recognition
- Word alignment
- Congressional voting

### Examples of latent structure

Object recognition

# Dataset for Supervised Part-of-Speech (POS) Tagging

Data:  $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$ 



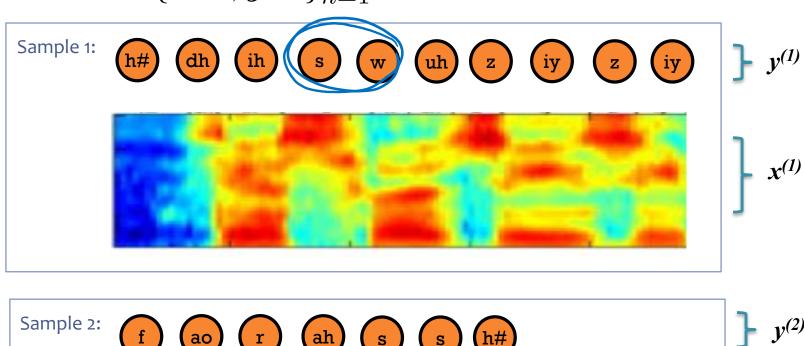
# Dataset for Supervised Handwriting Recognition

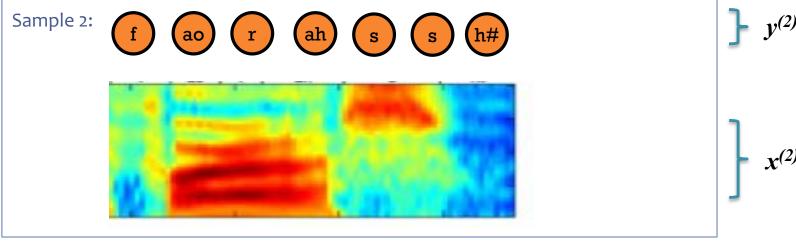
Data:  $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$ 



# Dataset for Supervised Phoneme (Speech) Recognition

Data:  $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$ 



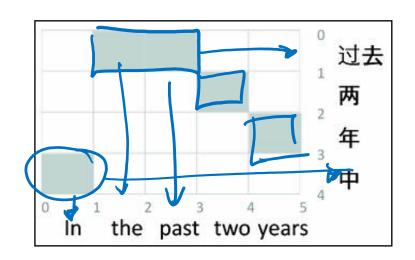


### Application:

## Word Alignment / Phrase Extraction

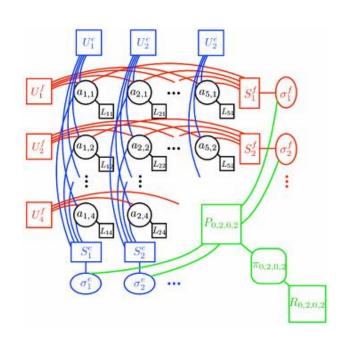
### Variables (boolean):

For each (Chinese phrase, English phrase) pair, are they linked?



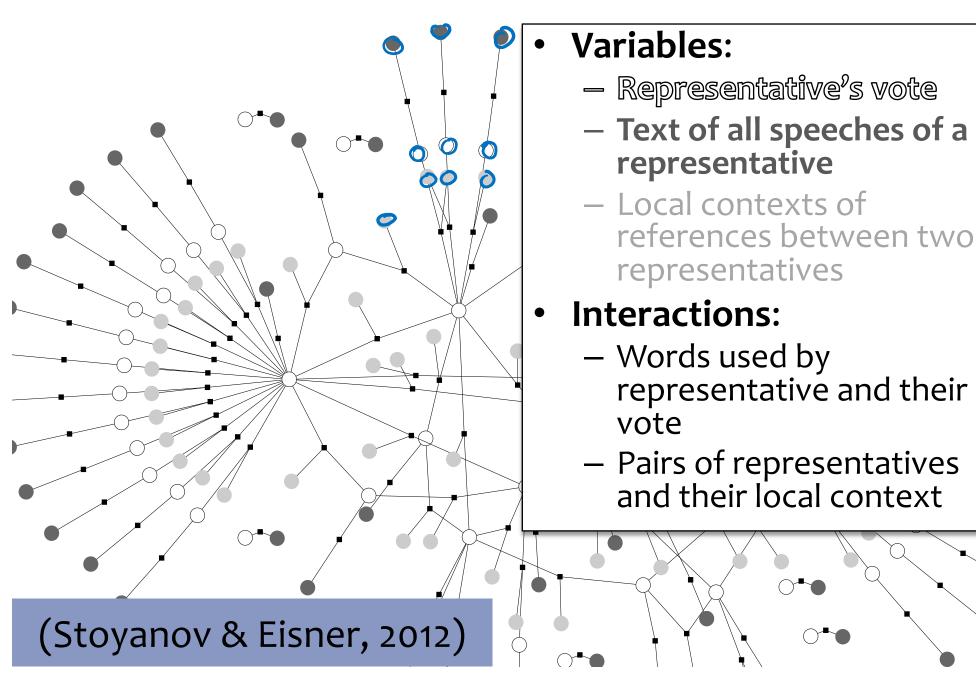
#### Interactions:

- Word fertilities
- Few "jumps" (discontinuities)
- Syntactic reorderings
- "ITG contraint" on alignment
- Phrases are disjoint (?)



### Application:

# Congressional Voting



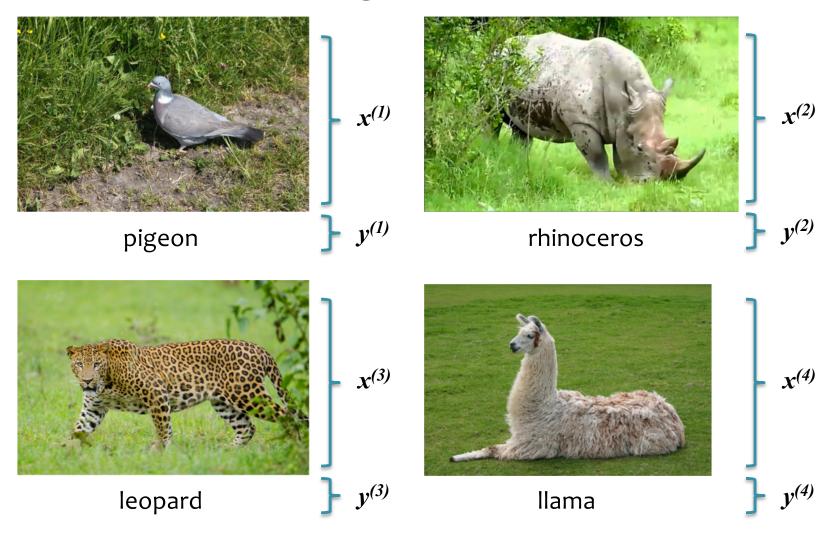
## Structured Prediction Examples

## Examples of structured prediction

- Part-of-speech (POS) tagging
- Handwriting recognition
- Speech recognition
- Word alignment
- Congressional voting

## Examples of latent structure

Object recognition

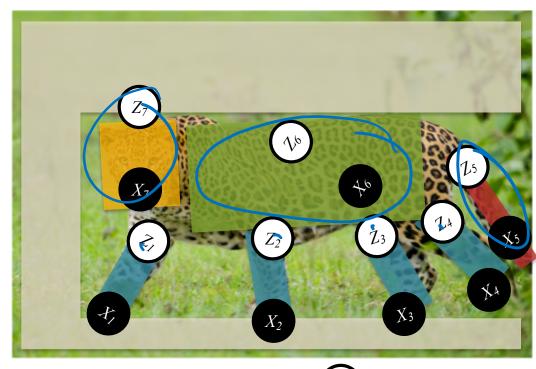


- Preprocess data into "patches"
- Posit a latent labeling z describing the object's parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time

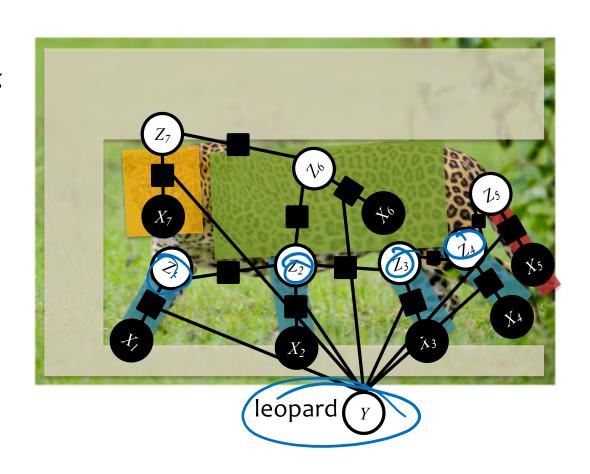


leopard

- Preprocess data into "patches"
- Posit a latent labeling z
   describing the object's
   parts (e.g. head, leg,
   tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time



- Preprocess data into "patches"
- Posit a latent labeling z
   describing the object's
   parts (e.g. head, leg,
   tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time



## Structured Prediction

# Preview of challenges to come...

Consider the task of finding the most probable
 assignment to the output

Classification 
$$\hat{y} = \operatorname*{argmax} p(y|\mathbf{x})$$
 where  $y \in \{+1, -1\}$ 

Structured Prediction 
$$\hat{\mathbf{y}} = \operatorname*{argmax} p(\mathbf{y}|\mathbf{x})$$
  $\mathbf{y} \in \mathcal{Y}$  where  $\mathbf{y} \in \mathcal{Y}$  and  $|\mathcal{Y}|$  is very large

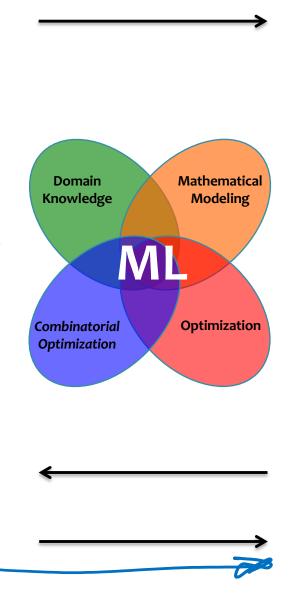
## Machine Learning

The data inspires the structures we want to predict



{best structure, marginals, partition function} for a new observation

(Inference is usually called as a subroutine in learning)

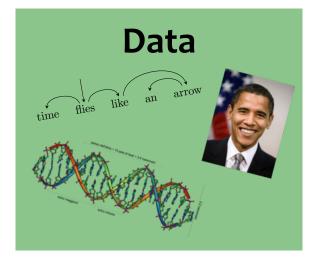


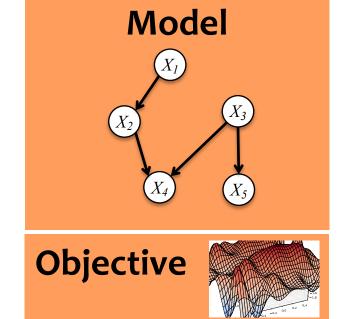
Our **model**defines a score
for each structure

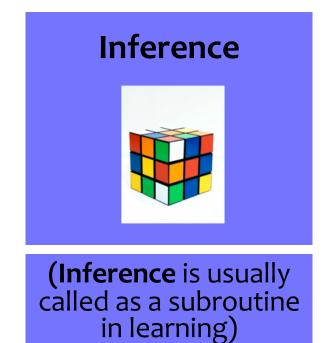
it also tells us what to optimize

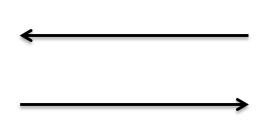
**Learning** tunes the parameters of the model

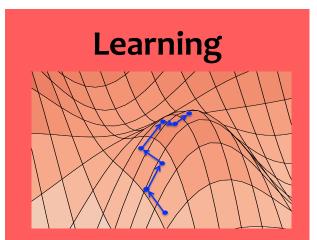
# Machine Learning











## **BACKGROUND**

# Background: Chain Rule of Probability

For random variables A and B:

$$P(A,B) = P(A|B)P(B)$$

For random variables  $X_1, X_2, X_3, X_4$ :

$$P(X_1, X_2, X_3, X_4) = P(X_1|X_2, X_3, X_4)$$

$$P(X_2|X_3, X_4)$$

$$P(X_3|X_4)$$

$$P(X_4)$$

# Background: Conditional Independence

Random variables A and B are conditionally independent given C if:

$$P(A,B|C) = P(A|C)P(B|C)$$
 (1)

or equivalently:

$$P(A|B,C) = P(A|C) \tag{2}$$

We write this as:

$$A \perp \!\!\! \perp B | C$$

Later we will also write: I < A,  $\{C\}$ , B >

# HIDDEN MARKOV MODEL (HMM)

### From Mixture Model to HMM

