

10-301/601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

Exam 2 Review



Hidden Markov Models

Matt Gormley Lecture 18 Nov. 4, 2022

Reminders

- Homework 6: Learning Theory / Generative Models
 - Out: Thu, Oct. 27
 - Due: Fri, Nov. 4 at 11:59pm
 - IMPORTANT: only 2 grace/late days permitted
- Practice Problems: Exam 2
 - Out: Fri, Nov. 4
- Exam 2
 - Thu, Nov. 10, 6:30pm 8:30pm

EXAM 2 LOGISTICS

Exam 2

- Time / Location
 - Time: Thu, Nov. 10, 6:30pm 8:30pm
 - Location & Seats: You have all been split across multiple rooms.
 Everyone has an assigned seat in one of these room. Please watch Piazza carefully for announcements.

Logistics

- Covered material: Lecture 8 Lecture 17
- Format of questions:
 - Multiple choice
 - True / False (with justification)
 - Derivations
 - Short answers
 - Interpreting figures
 - Implementing algorithms on paper
- No electronic devices
- You are allowed to bring one 8½ x 11 sheet of notes (front and back, handwritten with pen/pencil or tablet)

Topics for Exam 1

- Foundations
 - Probability, Linear
 Algebra, Geometry,
 Calculus
 - Optimization
- Important Concepts
 - Overfitting
 - Experimental Design

- Classification
 - Decision Tree
 - KNN
 - Perceptron
- Regression
 - Linear Regression

Topics for Exam 2

- Classification
 - Binary LogisticRegression
- Important Concepts
 - Stochastic Gradient
 Descent
 - Regularization
 - Feature Engineering
- Feature Learning
 - Neural Networks
 - Basic NN Architectures
 - Backpropagation

- Learning Theory
 - PAC Learning
- Generative Models
 - MLE / MAP
 - Naïve Bayes
 - Generative vs.
 Discriminative

- Regression
 - Linear Regression

SAMPLE QUESTIONS

3.2 Logistic regression

Given a training set $\{(x_i, y_i), i = 1, ..., n\}$ where $x_i \in \mathbb{R}^d$ is a feature vector and $y_i \in \{0, 1\}$ is a binary label, we want to find the parameters \hat{w} that maximize the likelihood for the training set, assuming a parametric model of the form

$$p(y = 1|x; w) = \frac{1}{1 + \exp(-w^T x)}.$$

The conditional log likelihood of the training set is

$$\ell(w) = \sum_{i=1}^{n} y_i \log p(y_i, | x_i; w) + (1 - y_i) \log(1 - p(y_i, | x_i; w)),$$

and the gradient is

$$\nabla \ell(w) = \sum_{i=1}^{n} (y_i - p(y_i|x_i; w))x_i.$$

- (b) [5 pts.] What is the form of the classifier output by logistic regression?
- (c) [2 pts.] **Extra Credit:** Consider the case with binary features, i.e, $x \in \{0,1\}^d \subset \mathbb{R}^d$, where feature x_1 is rare and happens to appear in the training set with only label 1. What is \hat{w}_1 ? Is the gradient ever zero for any finite w? Why is it important to include a regularization term to control the norm of \hat{w} ?

on? (integrated in d-dims)

$$\widehat{W} \in \mathbb{R}^d, \quad \widehat{W} \times = 0$$

$$\widehat{W} = \{0,1\}^d \subset \mathbb{R}^d, \quad \text{with only label 1.}$$

2.1 Train and test errors



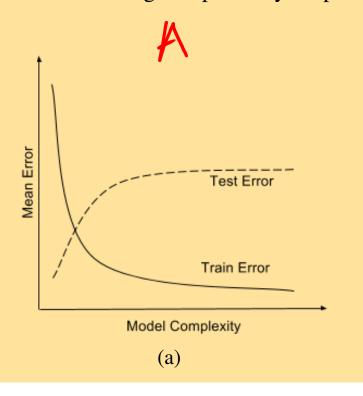
In this problem, we will see how you can debug a classifier by looking at its train and test errors. Consider a classifier trained till convergence on some training data $\mathcal{D}^{\text{train}}$, and tested on a separate test set $\mathcal{D}^{\text{test}}$. You look at the test error, and find that it is very high. You then compute the training error and find that it is close to 0.

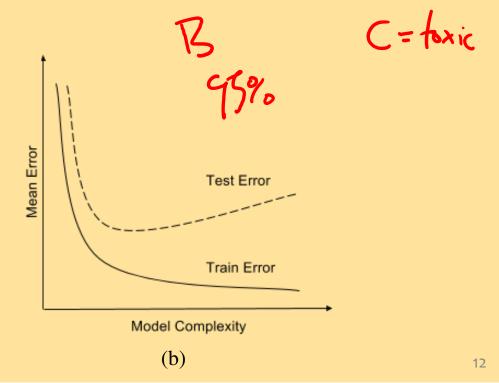
- 1. [4 pts] Which of the following is expected to help? Select all that apply.
 - (a) Increase the training data size. \$2%
 - (b) Decrease the training data size. 7%
 - (c) Increase model complexity (For example, if your classifier is an SVM, use a more complex kernel. Or if it is a decision tree, increase the depth).
 - (d) Decrease model complexity. 90%
 - (e) Train on a combination of $\mathcal{D}^{\text{train}}$ and $\mathcal{D}^{\text{test}}$ and test on $\mathcal{D}^{\text{test}}$
 - (f) Conclude that Machine Learning does not work.

2.1 Train and test errors

In this problem, we will see how you can debug a classifier by looking at its train and test errors. Consider a classifier trained till convergence on some training data $\mathcal{D}^{\text{train}}$, and tested on a separate test set $\mathcal{D}^{\text{test}}$. You look at the test error, and find that it is very high. You then compute the training error and find that it is close to 0.

[1 pts] Say you plot the train and test errors as a function of the model complexity. Which of the following two plots is your plot expected to look like?





5 Learning Theory [20 pts.]

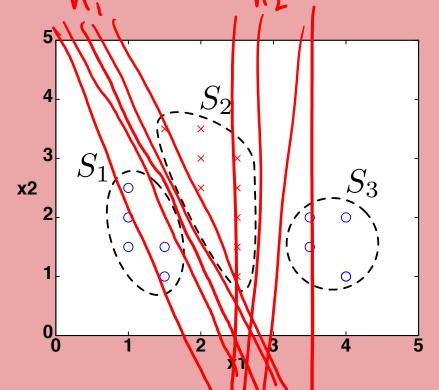
(a) [3 pts.] **T** or **F**: It is possible to label 4 points in \mathbb{R}^2 in all possible 2^4 ways via linear separators in \mathbb{R}^2 .

(d) [3 pts.] **T** or **F**: The VC dimension of a hypothesis space with infinite size is also infinite.

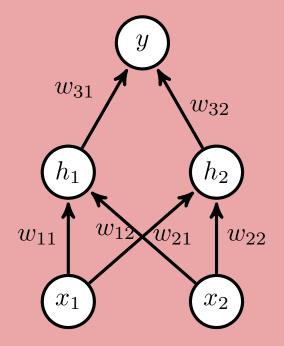




Neural Networks Q_3 A = bxi B = 1es C = No Can the neural network in Figure (b) correctly classify the dataset given in Figure (a)?



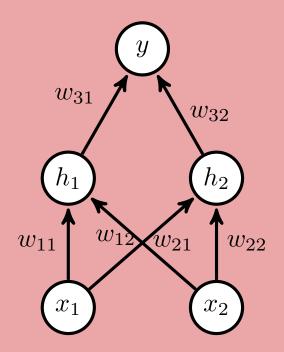
(a) The dataset with groups S_1 , S_2 , and S_3 .



(b) The neural network architecture

Neural Networks

Apply the backpropagation algorithm to obtain the partial derivative of the mean-squared error of y with the true value y^* with respect to the weight w_{22} assuming a sigmoid nonlinear activation function for the hidden layer.



(b) The neural network architecture

1.2 Maximum Likelihood Estimation (MLE)

Assume we have a random sample that is Bernoulli distributed $X_1, \ldots, X_n \sim \text{Bernoulli}(\theta)$. We are going to derive the MLE for θ . Recall that a Bernoulli random variable X takes values in $\{0,1\}$ and has probability mass function given by

$$P(X;\theta) = \theta^X (1-\theta)^{1-X}.$$

(a) [2 pts.] Derive the likelihood, $L(\theta; X_1, \ldots, X_n)$.

(c) **Extra Credit:** [2 pts.] Derive the following formula for the MLE: $\hat{\theta} = \frac{1}{n} \left(\sum_{i=1}^{n} X_i \right)$.

1.3 MAP vs MLE

Answer each question with **T** or **F** and **provide a one sentence explanation of your answer:**

(a) [2 pts.] **T** or **F**: In the limit, as *n* (the number of samples) increases, the MAP and MLE estimates become the same.

1.1 Naive Bayes

You are given a data set of 10,000 students with their sex, height, and hair color. You are trying to build a classifier to predict the sex of a student, so you randomly split the data into a training set and a testing set. Here are the specifications of the data set:

- $sex \in \{male, female\}$
- height $\in [0,300]$ centimeters
- hair \in {brown, black, blond, red, green}
- 3240 men in the data set
- 6760 women in the data set

Under the assumptions necessary for Naive Bayes (not the distributional assumptions you might naturally or intuitively make about the dataset) answer each question with **T** or **F** and **provide a one sentence explanation of your answer**:

(a) [2 pts.] **T** or **F**: As height is a continuous valued variable, Naive Bayes is not appropriate since it cannot handle continuous valued variables.

(c) [2 pts.] **T** or **F**: P(height|sex,hair) = P(height|sex).

DISCRIMINATIVE AND GENERATIVE CLASSIFIERS

Generative vs. Discriminative

- Naïve Bayes is a generative model:
 - Models the joint distribution P(X,Y) = P(X|Y)P(Y)
 - Learn by optimizing the (joint) likelihood
 - By modelling P(X|Y) and P(Y), we can generate new data points:
 - 1. Sample a label $y \sim P(Y)$
 - 2. Sample features $x_d \sim P(X_d | Y = y)$
- Logistic regression is a discriminative model
 - Directly models the conditional distribution P(Y|X)
 - Learn by optimizing the conditional likelihood
 - By just modelling P(Y|X), we can only discriminate (or distinguish) between classes.

Generative vs. Discriminative Finite Sample Analysis (Ng & Jordan, 2001)

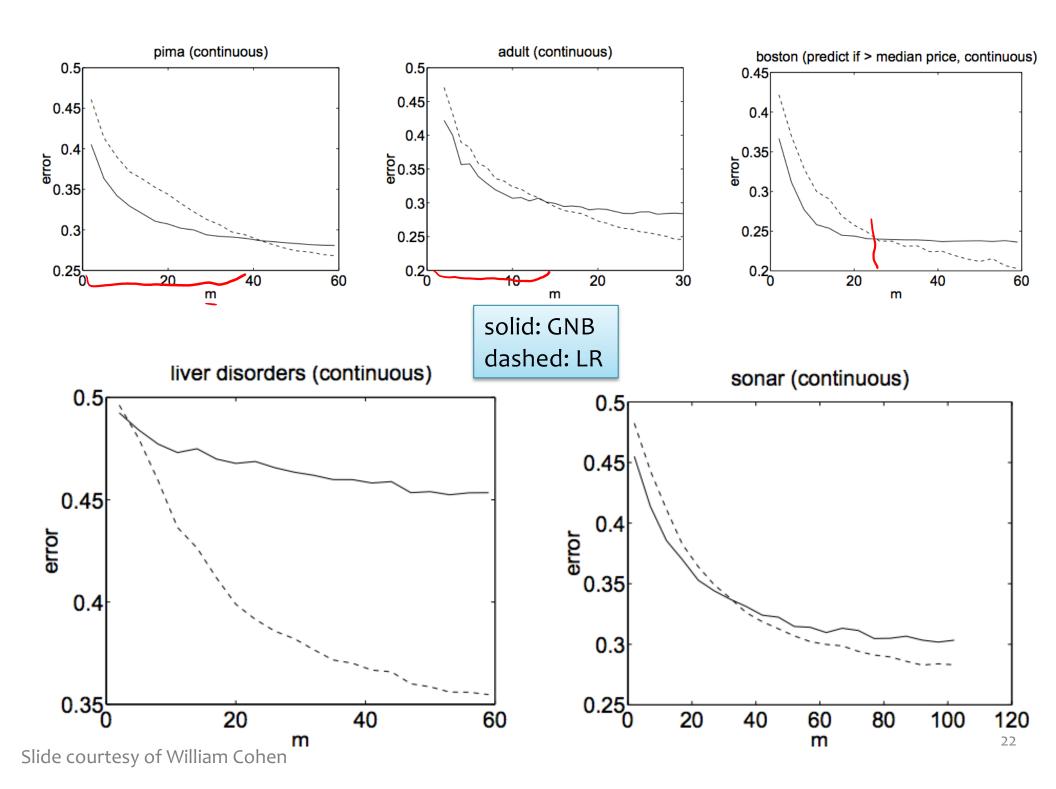
[Assume that we are learning from a finite training dataset]

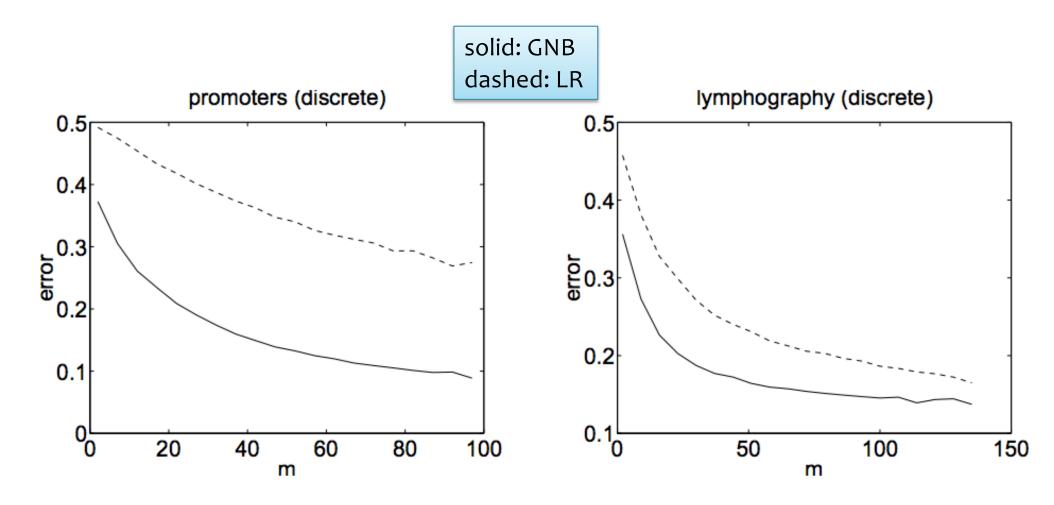
Naïve Bayes and logistic regression form a generativediscriminative model pair:

If model assumptions are correct: as the amount of training data increases, Gaussian Naïve Bayes and logistic regression approach the same (linear) decision boundary!

Furthermore, Gaussian Naïve Bayes is a more efficient learner (requires fewer samples) than Logistic Regression

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





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.

Naïve Bayes vs. Logistic Regression

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.

Naïve Bayes vs. Logistic Regression

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 Regression

Learning (MAP Estimation of Parameters)

Bernoulli Naïve Bayes:

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

Logistic Regression:

Parameters are not probabilities \rightarrow Gaussian prior encourages parameters to be close to zero

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

Summary

- Naïve Bayes provides a framework for generative modeling
- 2. Choose the feature distributions $P(X_d|Y)$ based on the data (e.g., Bernoulli for binary features, Gaussian for continuous features)
- 3. Train using MLE or MAP estimation
- 4. Make predictions by maximizing the posterior P(Y|x')

Learning Objectives

Naïve Bayes

You should be able to...

- 1. Write the generative story for Naïve Bayes
- 2. Create a new Naïve 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 Naïve 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 Naïve Bayes
- 6. Describe the tradeoffs of generative vs. discriminative models
- 7. Implement Bernoulli Naïve Bayes
- 8. Describe how the variance affects whether a Gaussian Naïve Bayes model will have a linear or nonlinear decision boundary

THE BIG PICTURE

ML Big Picture

Learning Paradigms:

What data is available and when? What form of prediction?

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning —
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

Theoretical Foundations:

What principles guide learning?

- probabilistic
- ☐ information theoretic
- evolutionary search
- ☐ ML as optimization

Problem Formulation:

What is the structure of our output prediction?

boolean Binary Classification

categorical Multiclass Classification

ordinal Ordinal Classification

real Regression

ordering Ranking

multiple discrete Structured Prediction

multiple continuous (e.g. dynamical systems)

both discrete & (e.g. mixed graphical models)

cont.

Big Ideas in ML:

Which are the ideas driving development of the field?

- inductive bias
- generalization / overfitting

splication Area

- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

Facets of Building ML Systems:

How to build systems that are robust, efficient, adaptive, effective?

- 1. Data prep
- 2. Model selection
- 3. Training (optimization / search)
- 4. Hyperparameter tuning on validation data
- 5. (Blind) Assessment on test

Classification and Regression: The Big Picture

Recipe for Machine Learning

- 1. Given data $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$
- 2. (a) Choose a decision function $h_{\theta}(\mathbf{x}) = \cdots$ (parameterized by θ)
 - (b) Choose an objective function $J_{\mathcal{D}}(\boldsymbol{\theta}) = \cdots$ (relies on data)
- 3. Learn by choosing parameters that optimize the objective $J_{\mathcal{D}}(\boldsymbol{\theta})$

$$\hat{\boldsymbol{\theta}} pprox \operatorname*{argmin}_{\boldsymbol{\theta}} J_{\mathcal{D}}(\boldsymbol{\theta})$$

4. Predict on new test example $\mathbf{x}_{\mathsf{new}}$ using $h_{\boldsymbol{\theta}}(\cdot)$

$$\hat{y} = h_{\boldsymbol{\theta}}(\mathbf{x}_{\mathsf{new}})$$

Optimization Method

- Gradient Descent: $m{ heta} o m{ heta} \gamma
 abla_{m{ heta}} J(m{ heta})$
- SGD: $\theta \to \theta \gamma \nabla_{\theta} J^{(i)}(\theta)$ for $i \sim \mathsf{Uniform}(1,\ldots,N)$

where
$$J(oldsymbol{ heta}) = rac{1}{N} \sum_{i=1}^N J^{(i)}(oldsymbol{ heta})$$

- mini-batch SGD
- closed form
 - 1. compute partial derivatives
 - 2. set equal to zero and solve

Decision Functions

- Perceptron: $h_{\boldsymbol{\theta}}(\mathbf{x}) = \operatorname{sign}(\boldsymbol{\theta}^T \mathbf{x})$
- Linear Regression: $h_{\theta}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$
- Discriminative Models: $h_{\theta}(\mathbf{x}) = \operatorname*{argmax}_{y} p_{\theta}(y \mid \mathbf{x})$
 - \circ Logistic Regression: $p_{\theta}(y = 1 \mid \mathbf{x}) = \sigma(\boldsymbol{\theta}^T \mathbf{x})$
 - $\circ \text{ Neural Net (classification):} \\ p_{\boldsymbol{\theta}}(y=1\mid \mathbf{x}) = \sigma((\mathbf{W}^{(2)})^T \sigma((\mathbf{W}^{(1)})^T \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$
- Generative Models: $h_{\theta}(\mathbf{x}) = \operatorname*{argmax}_{y} p_{\theta}(\mathbf{x}, y)$
 - $\circ \;$ Naive Bayes: $p_{m{ heta}}(\mathbf{x},y) = p_{m{ heta}}(y) \prod_{m=1}^M p_{m{ heta}}(x_m \mid y)$

Objective Function

- ullet MLE: $J(oldsymbol{ heta}) = -\sum_{i=1}^N \log p(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$
- MCLE: $J(oldsymbol{ heta}) = -\sum_{i=1}^N \log p(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)})$
- L2 Regularized: $J'(\theta) = J(\theta) + \lambda ||\theta||_2^2$ (same as Gaussian prior $p(\theta)$ over parameters)
- L1 Regularized: $J'(\theta) = J(\theta) + \lambda ||\theta||_1$ (same as Laplace prior $p(\theta)$ over parameters)

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: y
- Many real-world problems require structured prediction
 - Given observations: $\mathbf{x} = (x_1, x_2, ..., x_K)$
 - Predict a structure: $y = (y_1, y_2, ..., y_J)$
- 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: $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$

Sample 1:	n	flies	p like	an	$\begin{array}{c c} & & \\ & &$
Sample 2:	n	n	like	d	$\begin{array}{c c} & & \\ & &$
Sample 3:	n	fly	with	n	$\begin{cases} \mathbf{n} \\ \mathbf{vings} \end{cases} \mathbf{y}^{(3)}$
Sample 4:	with	time	you	will	$\begin{cases} \mathbf{v} \\ \mathbf{see} \end{cases} \mathbf{y}^{(4)}$

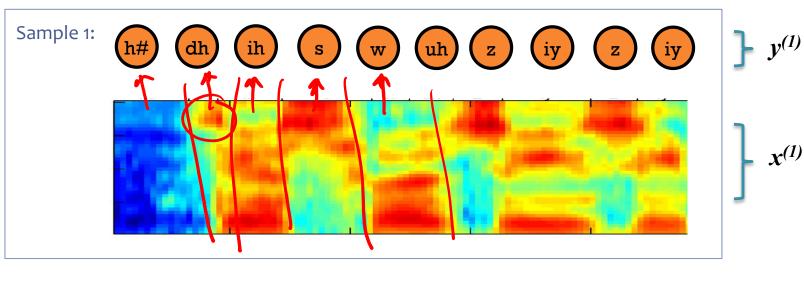
Dataset for Supervised Handwriting Recognition

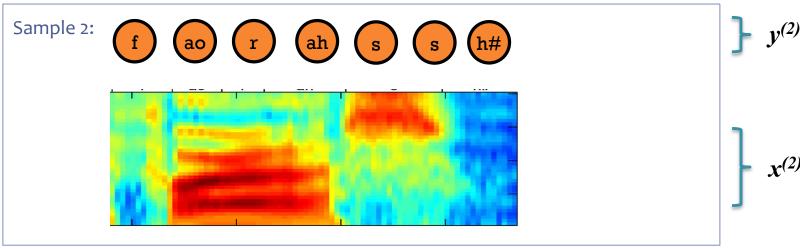
Data: $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$



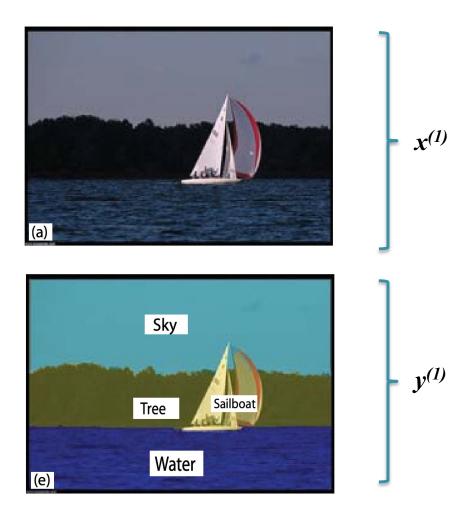
Dataset for Supervised Phoneme (Speech) Recognition

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



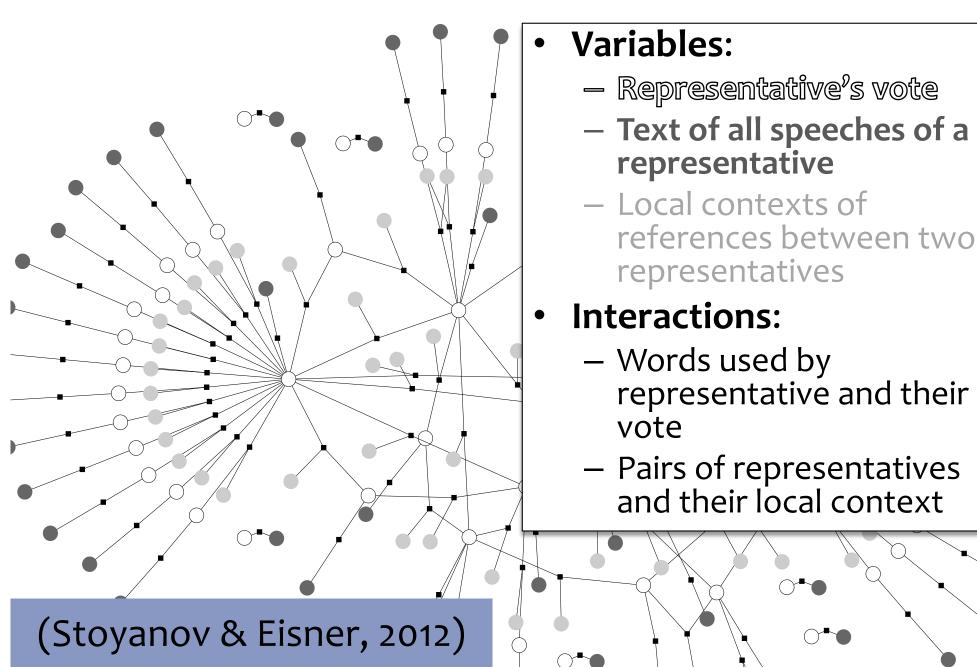


(very small) Dataset for Scene Understanding



Application:

Congressional Voting



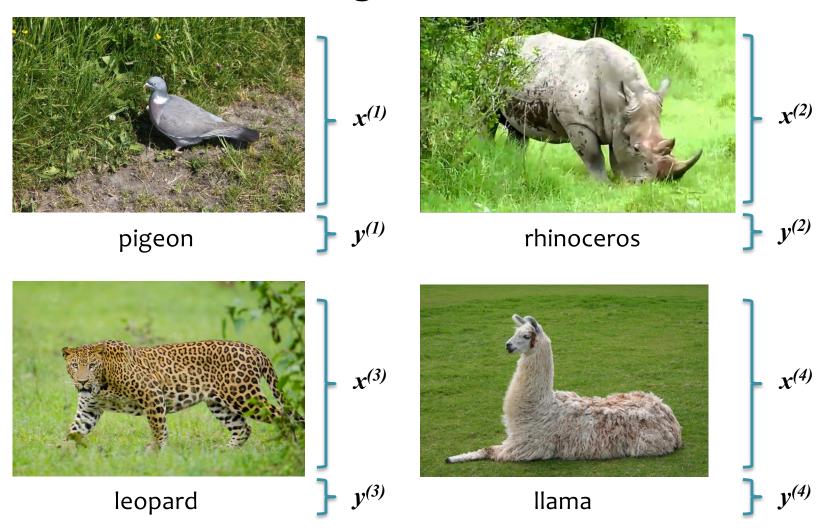
Structured Prediction Examples

Examples of structured prediction

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

Examples of latent structure

Object recognition

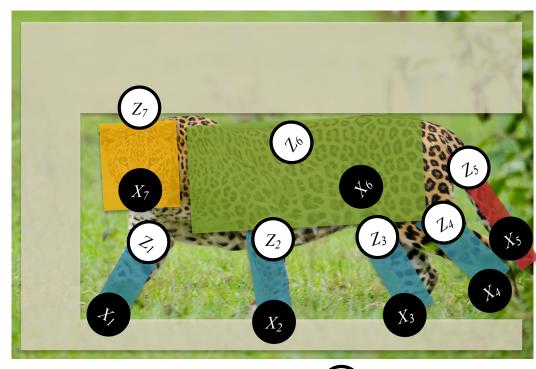


- Preprocess data into "patches"
- Posit a latent labeling z describing the object's parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time

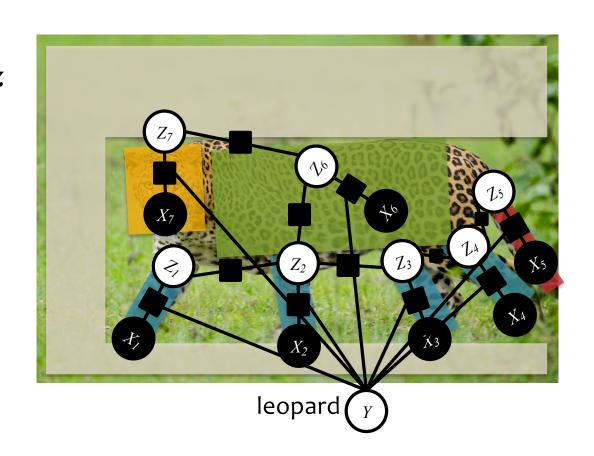


leopard

- 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

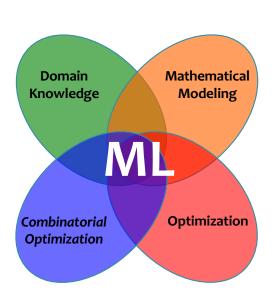
Preview of challenges to come...

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

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

Machine Learning

The data inspires
the structures
we want to
predict



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

It also tells us what to optimize

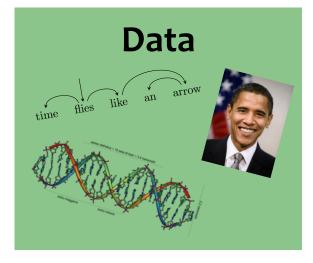
Inference finds

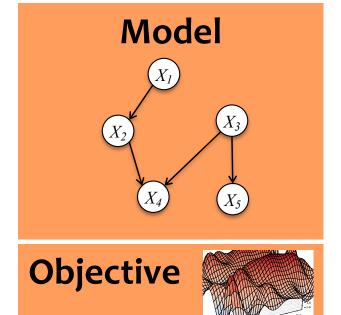
{best structure, marginals, partition function} for a new observation

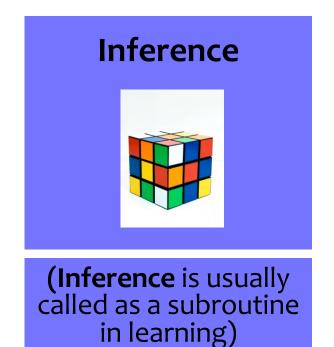
(Inference is usually called as a subroutine in learning)

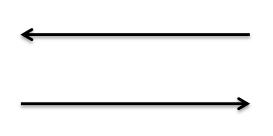
Learning tunes the parameters of the model

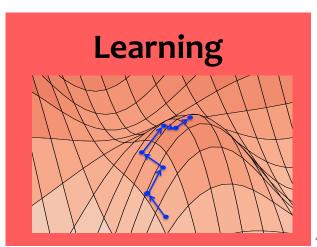
Machine Learning











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:

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

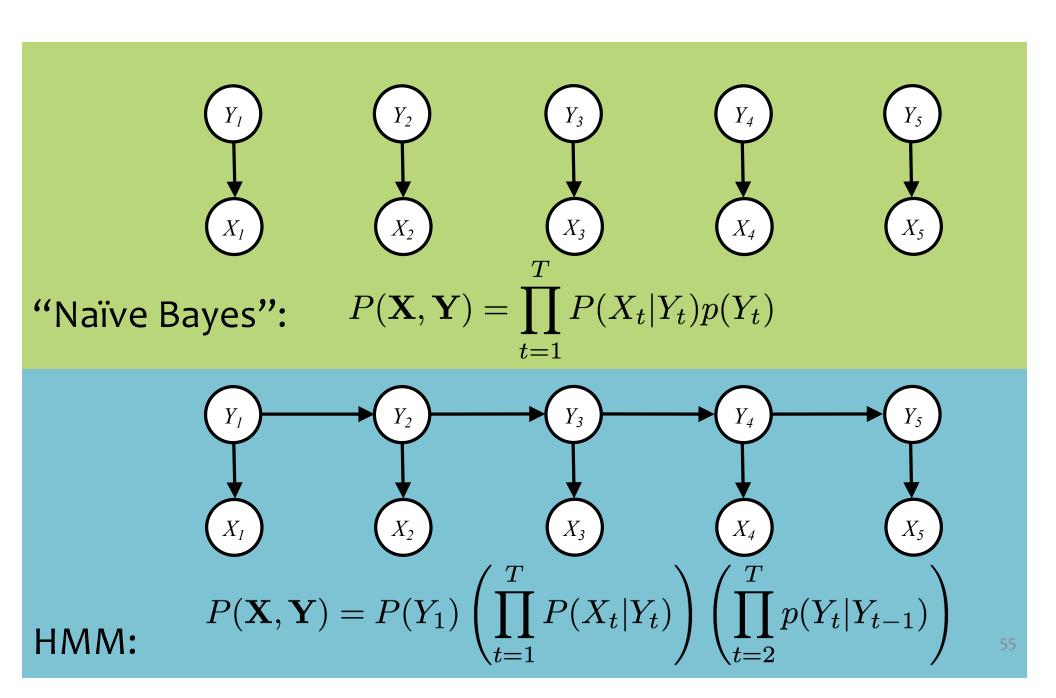
We write this as:

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

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

HIDDEN MARKOV MODEL (HMM)

From Mixture Model to HMM



Markov Models

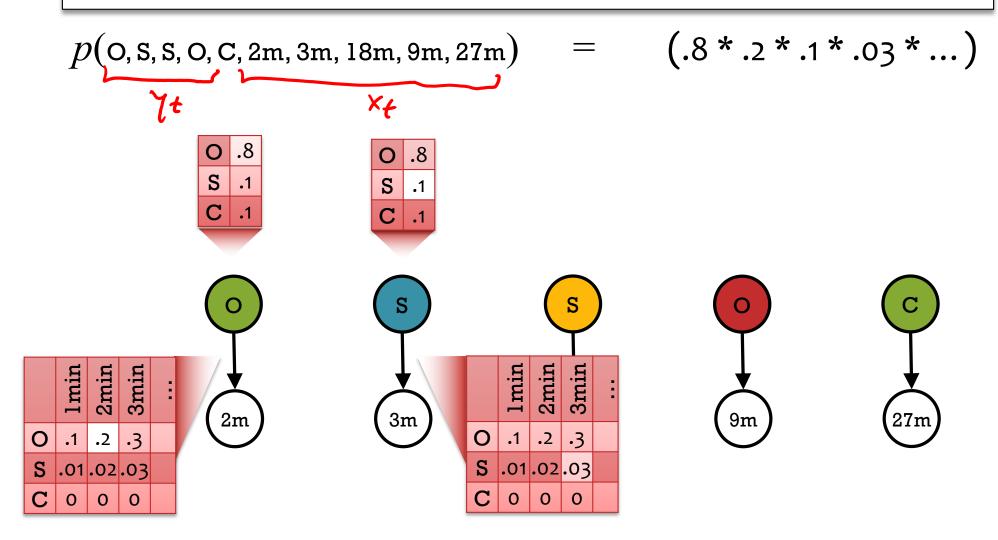
Whiteboard

- Example: Tunnel Closures[courtesy of Roni Rosenfeld]
- First-order Markov assumption
- Conditional independence assumptions



Mixture Model for Time Series Data

We could treat each (tunnel state, travel time) pair as independent. This corresponds to a Naïve Bayes model with a single feature (travel time).



Hidden Markov Model

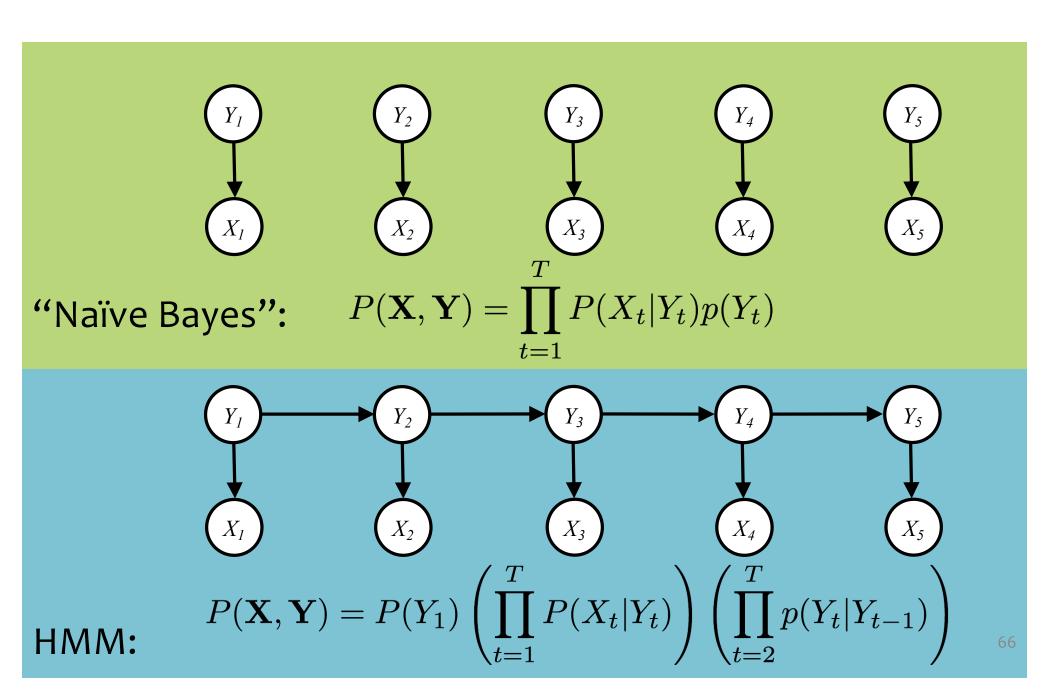
A Hidden Markov Model (HMM) provides a joint distribution over the the tunnel states / travel times with an assumption of dependence between adjacent tunnel states.

.01.02.03

.2

.01.02.03

From Mixture Model to HMM



SUPERVISED LEARNING FOR HMMS

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)

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_1} = \dots$$
$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_2} = \dots$$
$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_M} = \dots$$

4. Set derivatives to zero and solve for θ

$$\partial \ell(\theta)/\partial \theta_{\rm m} = {\rm o \ for \ all \ m} \in \{1, ..., M\}$$

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

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

MLE of Categorical Distribution

1. Suppose we have a **dataset** obtained by repeatedly rolling a M-sided (weighted) die N times. That is, we have data

$$\mathcal{D} = \{x^{(i)}\}_{i=1}^{N}$$

where $x^{(i)} \in \{1, \dots, M\}$ and $x^{(i)} \sim \mathsf{Categorical}(\phi)$.

2. A random variable is **Categorical** written $X \sim \mathsf{Categorical}(\phi)$ iff

$$P(X = x) = p(x; \phi) = \phi_x$$

where $x \in \{1, \dots, M\}$ and $\sum_{m=1}^{M} \phi_m = 1$. The **log-likelihood** of the data becomes:

$$\ell(oldsymbol{\phi}) = \sum_{i=1}^N \log \phi_{x^{(i)}}$$
 s.t. $\sum_{m=1}^M \phi_m = 1$

3. Solving this constrained optimization problem yields the maximum likelihood estimator (MLE):

$$\phi_m^{MLE} = \frac{N_{x=m}}{N} = \frac{\sum_{i=1}^{N} \mathbb{I}(x^{(i)} = m)}{N}$$

