



#### 10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

## **HMMs**



## **Bayesian Networks**

Matt Gormley Lecture 21 Apr. 01, 2020

### Reminders

- Practice Problems for Exam 2
  - Out: Fri, Mar 20
- Midterm Exam 2
  - Thu, Apr 2 evening exam, details announced on Piazza
- Homework 7: HMMs
  - Out: Thu, Apr 02
  - Due: Fri, Apr 10 at 11:59pm
- Today's In-Class Poll
  - http://poll.mlcourse.org

# THE FORWARD-BACKWARD ALGORITHM

## Forward-Backward Algorithm

Define: 
$$\alpha_{\xi}(k) \triangleq p(x_1, ..., x_{\xi}, y_{\xi} = k)$$
 $\beta_{\xi}(k) \triangleq p(x_{\xi+1}, ..., x_{\xi}, y_{\xi} = k)$ 
 $\beta_{\xi}(k) \triangleq p(x_{\xi+1}, ..., x_{\xi}, y_{\xi} = k)$ 
 $\beta_{\xi}(k) \triangleq p(x_{\xi+1}, ..., x_{\xi}, y_{\xi} = k)$ 
 $\beta_{\xi}(k) = p(x_{\xi+1}, ..., x_{\xi}, y_{\xi} = k)$ 
 $\beta_{\xi}(k) = \beta_{\xi}(k) = \beta_{\xi}(k)$ 
 $\beta_{\xi}(k) = \beta_{\xi}(k)$ 
 $\beta_{\xi}(k)$ 

## Inference for HMMs

#### Whiteboard

- Forward-backward algorithm (edge weights version)
- Viterbi algorithm(edge weights version)

## Forward-Backward Algorithm

Define: 
$$\alpha_{t}(k) \triangleq p(x_{1},...,x_{t},y_{t}=k)$$
 $\beta_{t}(k) \triangleq p(x_{1},...,x_{t}|y_{t}=k)$ 
 $\beta_{t}(k) \triangleq p(x_{t+1},...,x_{t}|y_{t}=k)$ 
 $\gamma_{t+1} = END$ 

Define:  $\gamma_{t+1} = \gamma_{t+1} = \gamma_{t+1} = \gamma_{t} = \gamma_{t+1} = \gamma_{t+1}$ 

## Derivation of Forward Algorithm

Definition: 
$$X_{\xi}(k) \triangleq p(x_1, ..., x_{\xi}, y_{\xi} = k)$$

Derivation:
$$(X_{T}(ENO) = p(x_1, ..., x_{\tau}, y_{T} = END))$$

$$= p(x_1, ..., x_{\tau}, y_{\tau}) p(y_{\tau})$$

$$= p(x_1 | y_{\tau}) p(x_1, ..., x_{\tau-1}, y_{\tau}) p(y_{\tau})$$

$$= p(x_{\tau} | y_{\tau}) p(x_1, ..., x_{\tau-1}, y_{\tau})$$

$$= p(x_{\tau} | y_{\tau}) p(x_1, ..., x_{\tau-1}, y_{\tau})$$

$$= p(x_{\tau} | y_{\tau}) p(x_1, ..., x_{\tau-1}, y_{\tau-1}, y_{\tau})$$

$$= p(x_{\tau} | y_{\tau}) p(x_1, ..., x_{\tau-1}, y_{\tau-1}, y_{\tau})$$

$$= p(x_{\tau} | y_{\tau}) p(x_1, ..., x_{\tau-1}, y_{\tau-1}, y_{\tau-1}, y_{\tau-1}) p(y_{\tau-1})$$

$$= p(x_{\tau} | y_{\tau}) p(x_1, ..., x_{\tau-1}, y_{\tau-1}, y_{\tau-1}, y_{\tau-1}) p(y_{\tau-1})$$

$$= p(x_{\tau} | y_{\tau}) p(x_1, ..., x_{\tau-1}, y_{\tau-1}, y_{\tau-$$

## Viterbi Algorithm

Define: 
$$\omega_{\xi}(k) \triangleq \max_{y_1, \dots, y_{\xi-1}, y_{\xi}} p(x_1, \dots, x_{\xi}, y_1, \dots, y_{\xi-1}, y_{\xi} = k)$$

"bulk points"

 $b_{\xi}(k) \triangleq \alpha_{\xi} \max_{y_1, \dots, y_{\xi-1}} p(x_1, \dots, x_{\xi}, y_1, \dots, y_{\xi-1}, y_{\xi} = k)$ 

Assume  $y_0 = START$ 

(2) For  $t = 1, \dots, T$ :

For  $k = 1, \dots, K$ :

 $\omega_{\xi}(k) = \max_{j \in \{1, \dots, K\}} p(x_{\xi} | y_{\xi} = k) \omega_{\xi_{\xi-1}}(j) p(y_{\xi} = k | y_{\xi-1} = j)$ 
 $b_{\xi}(k) = \max_{j \in \{1, \dots, K\}} p(x_{\xi} | y_{\xi} = k) \omega_{\xi_{\xi-1}}(j) p(y_{\xi} = k | y_{\xi-1} = j)$ 

(3) Compute Most Probable Assignment

 $\hat{y}_T = b_{T+1}(END) = \sum_{j \in \{1, \dots, k\}} p(x_{\xi} | y_{\xi} = k) \omega_{\xi_{\xi}}(j) p(y_{\xi} = k | y_{\xi-1} = j)$ 

For  $t = T-1, \dots, 1$ 
 $\hat{y}_t = b_{t+1}(\hat{y}_{t+1})$ 

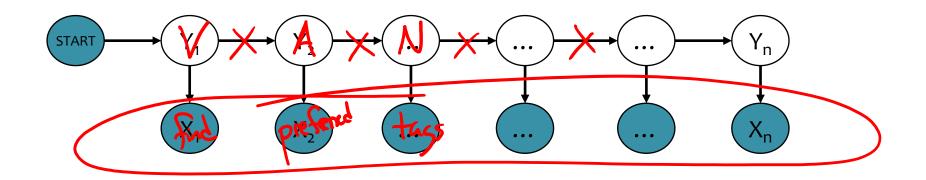
"bulk points"

## Inference in HMMs

What is the **computational complexity** of inference for HMMs?

- The naïve (brute force) computations for Evaluation, Decoding, and Marginals take exponential time, O(K<sup>T</sup>)
- The forward-backward algorithm and Viterbi algorithm run in polynomial time, O(T\*K²)
  - Thanks to dynamic programming!

# Shortcomings of Hidden Markov Models



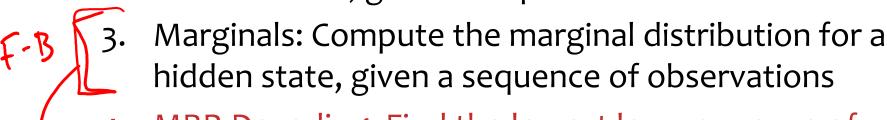
- HMM models capture dependences between each state and only its corresponding observation
  - NLP example: In a sentence segmentation task, each segmental state may depend not just on a single word (and the adjacent segmental stages), but also on the (nonlocal) features of the whole line such as line length, indentation, amount of white space, etc.
- Mismatch between learning objective function and prediction objective function
  - HMM learns a joint distribution of states and observations P(Y, X), but in a prediction task, we need the conditional probability P(Y|X).

## **MBR DECODING**

## Inference for HMMs

Four

- Three Inference Problems for an HMM
  - 1. Evaluation: Compute the probability of a given sequence of observations
  - Viterbi Decoding: Find the most-likely sequence of hidden states, given a sequence of observations



4. MBR Decoding: Find the lowest loss sequence of hidden states, given a sequence of observations (Viterbi decoding is a special case)

## Minimum Bayes Risk Decoding

- Suppose we given a loss function l(y', y) and are asked for a single tagging
- How should we choose just one from our probability distribution p(y|x)?
- A minimum Bayes risk (MBR) decoder h(x) returns the variable assignment with minimum **expected** loss under the model's distribution

$$h_{oldsymbol{ heta}}(oldsymbol{x}) = \underset{\hat{oldsymbol{y}}}{\operatorname{argmin}} \ \mathbb{E}_{oldsymbol{y} \sim p_{oldsymbol{ heta}(\cdot | oldsymbol{x})}} [\ell(\hat{oldsymbol{y}}, oldsymbol{y})]$$

$$= \underset{\hat{oldsymbol{y}}}{\operatorname{argmin}} \ \sum_{oldsymbol{y}} p_{oldsymbol{ heta}}(oldsymbol{y} | oldsymbol{x})\ell(\hat{oldsymbol{y}}, oldsymbol{y})$$

Minimum Bayes Risk Decoding

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \underset{\hat{\boldsymbol{y}}}{\operatorname{argmin}} \ \mathbb{E}_{\boldsymbol{y} \sim p_{\boldsymbol{\theta}}(\cdot | \boldsymbol{x})}[\ell(\hat{\boldsymbol{y}}, \boldsymbol{y})]$$

Consider some example loss functions:

The 0-1 loss function returns 1 only if the two assignments are identical and 0 otherwise:

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}) = 1 - \mathbb{I}(\hat{\boldsymbol{y}}, \boldsymbol{y})$$

The MBR decoder is:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \underset{\hat{\boldsymbol{y}}}{\operatorname{argmin}} \sum_{\boldsymbol{y} \in \mathcal{V}_{\boldsymbol{x}}} p_{\boldsymbol{\theta}}(\boldsymbol{y} \mid \boldsymbol{x}) (1 - \mathbb{I}(\hat{\boldsymbol{y}}, \boldsymbol{y}))$$

$$= \underset{\hat{\boldsymbol{y}}}{\operatorname{argmax}} p_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}} \mid \boldsymbol{x})$$

$$= \underset{\hat{\boldsymbol{y}}}{\operatorname{argmax}} p_{\boldsymbol{\theta}}(\hat{\boldsymbol{y}} \mid \boldsymbol{x})$$

which is exactly the Viterbi decoding problem!

## Minimum Bayes Risk Decoding

$$h_{m{ heta}}(m{x}) = \operatorname*{argmin}_{\hat{m{y}}} \mathbb{E}_{m{y} \sim p_{m{ heta}}(\cdot | m{x})} [\ell(\hat{m{y}}, m{y})]_{m{A}}$$

Consider some example loss functions:

The **Hamming loss** corresponds to accuracy and returns the number of incorrect variable assignments:

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}) = \sum_{i=1}^{V} (1 - \mathbb{I}(\hat{y}_i, y_i))$$

The MBR decoder is:

$$\hat{y}_i = h_{\boldsymbol{\theta}}(\boldsymbol{x})_i = \underset{\hat{y}_i}{\operatorname{argmax}} p_{\boldsymbol{\theta}}(\hat{y}_i \mid \boldsymbol{x})$$

This decomposes across variables and requires the variable marginals.

1 from E-B

## Learning Objectives

#### **Hidden Markov Models**

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

- Show that structured prediction problems yield high-computation inference problems
- 2. Define the first order Markov assumption
- 3. Draw a Finite State Machine depicting a first order Markov assumption
- 4. Derive the MLE parameters of an HMM
- 5. Define the three key problems for an HMM: evaluation, decoding, and marginal computation
- 6. Derive a dynamic programming algorithm for computing the marginal probabilities of an HMM
- 7. Interpret the forward-backward algorithm as a message passing algorithm
- 8. Implement supervised learning for an HMM
- 9. Implement the forward-backward algorithm for an HMM
- 10. Implement the Viterbi algorithm for an HMM
- 11. Implement a minimum Bayes risk decoder with Hamming loss for an HMM

## Bayes Nets Outline

#### Motivation

Structured Prediction

#### Background

- Conditional Independence
- Chain Rule of Probability

#### Directed Graphical Models

- Writing Joint Distributions
- Definition: Bayesian Network
- Qualitative Specification
- Quantitative Specification
- Familiar Models as Bayes Nets

#### Conditional Independence in Bayes Nets

- Three case studies
- D-separation
- Markov blanket

#### Learning

- Fully Observed Bayes Net
- (Partially Observed Bayes Net)

#### Inference

- Background: Marginal Probability
- Sampling directly from the joint distribution
- Gibbs Sampling

Bayesian Networks

## DIRECTED GRAPHICAL MODELS

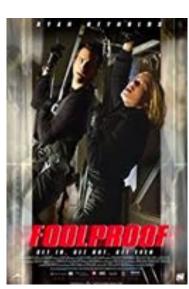
## Example: Ryan Reynolds' Voicemail



## Example: Ryan Reynolds Voicemail











## Example: Ryan Reynolds' Voicemail

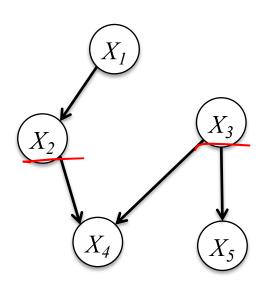


## Directed Graphical Models (Bayes Nets)

#### Whiteboard

- Example: Ryan Reynolds' Voicemail
- Writing Joint Distributions
  - Idea #1: Giant Table
  - Idea #2: Rewrite using chain rule
  - Idea #3: Assume full independence
  - Idea #4: Drop variables from RHS of conditionals
- Definition: Bayesian Network

## Bayesian Network



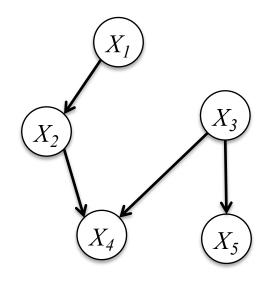
$$p(X_1, X_2, X_3, X_4, X_5) =$$

$$p(X_5|X_3)p(X_4|X_2, X_3)$$

$$p(X_3)p(X_2|X_1)p(X_1)$$

## Bayesian Network

#### **Definition:**



$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i \mid parents(X_i))$$

- A Bayesian Network is a directed graphical model
- It consists of a graph G and the conditional probabilities P
- These two parts full specify the distribution:
  - Qualitative Specification: G
  - Quantitative Specification: P



## Qualitative Specification

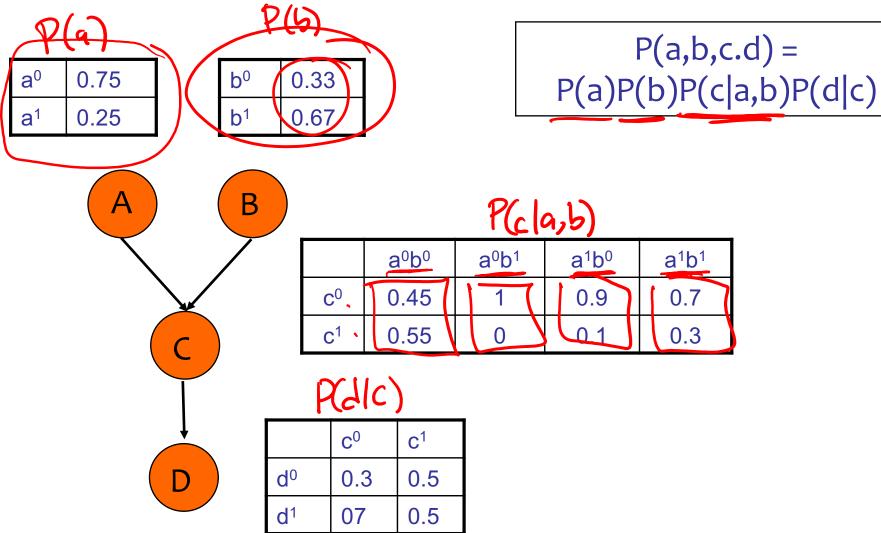
- Where does the qualitative specification come from?
  - Prior knowledge of causal relationships
  - Prior knowledge of modular relationships
  - Assessment from experts
  - Learning from data (i.e. structure learning)
  - We simply prefer a certain architecture (e.g. a layered graph)

**—** ...

## Quantitative Specification

**Example: Conditional probability tables (CPTs)** 

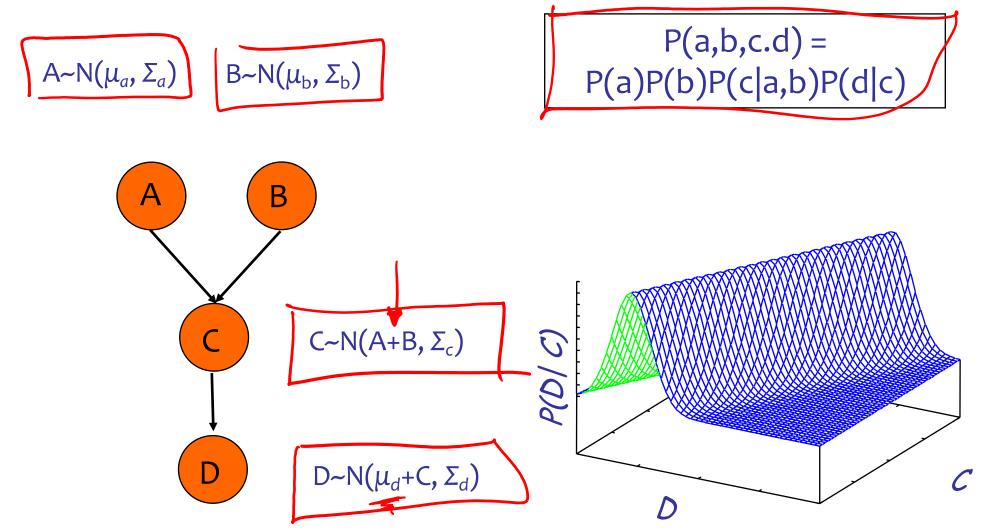
for discrete random variables



## Quantitative Specification

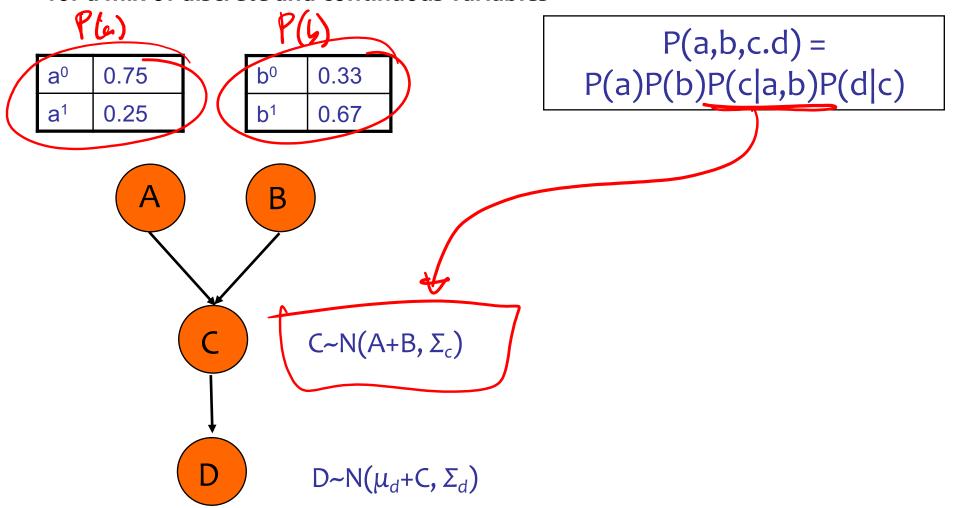
Example: Conditional probability density functions (CPDs)

for continuous random variables



## Quantitative Specification

**Example: Combination of CPTs and CPDs** for a mix of discrete and continuous variables

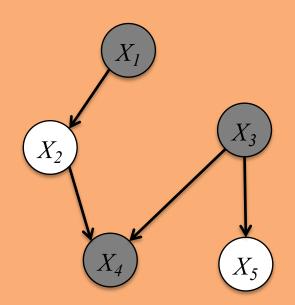


## **Observed Variables**

 In a graphical model, shaded nodes are "observed", i.e. their values are given

## **Example:**

$$P(X_2, X_5 \mid X_1 = 0, X_3 = 1, X_4 = 1)$$



## Familiar Models as Bayesian Networks

#### **Question:**

Match the model name to the corresponding Bayesian Network

- 1. Logistic Regression  $\mathfrak{D}$
- 2. Linear Regression 🤊
- 3. Bernoulli Naïve Bayes 🗛
- 4. Gaussian Naïve Bayes A
- 5. 1D Gaussian E

# (3) p(y,x,,..,xm) = p(y) P(x,1y) --- P(x,n1y) (1) P(y| x,1,...,xm) = p(y) P(x,1y) --- P(x,n1y) (2) x (3) p(x | M, o<sup>2</sup>)

#### **Answer:**

