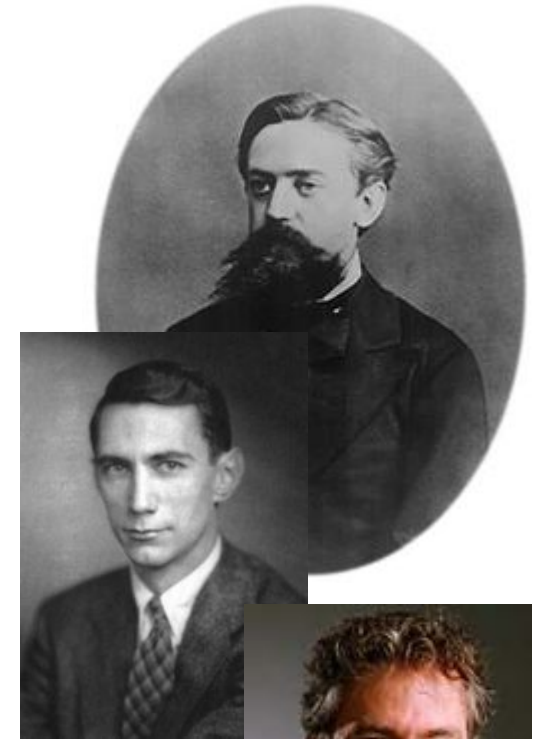# Hidden Markov Models

Matt Gormley
Lecture 20
Mar. 30, 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**
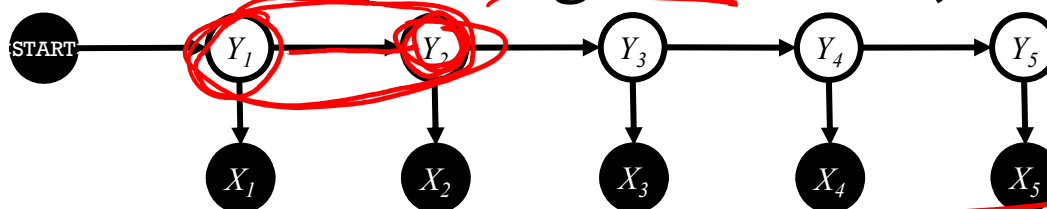  - **http://poll.mlcourse.org**

# HMMs: History

- Markov chains: Andrey Markov (1906)
  - Random walks and Brownian motion
- Used in Shannon's work on information theory (1948)
- Baum-Welsh learning algorithm: late 60's, early 70's.
  - Used mainly for speech in 60s-70s.
- Late 80's and 90's: David Haussler (major player in learning theory in 80's) began to use HMMs for modeling biological sequences
- Mid-late 1990's: Dayne Freitag/Andrew McCallum
  - Freitag thesis with Tom Mitchell on IE from Web using logic programs, grammar induction, etc.
  - McCallum: multinomial Naïve Bayes for text
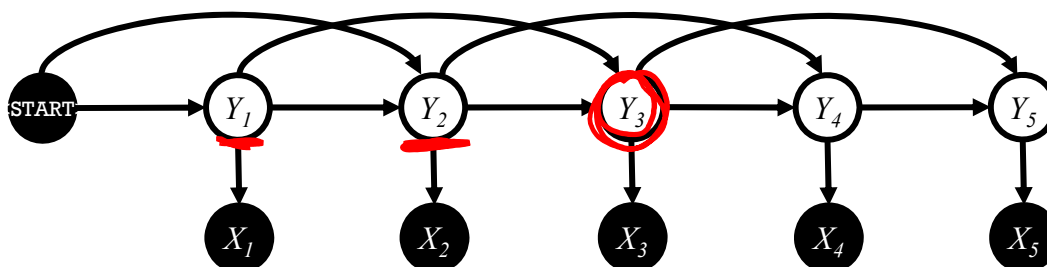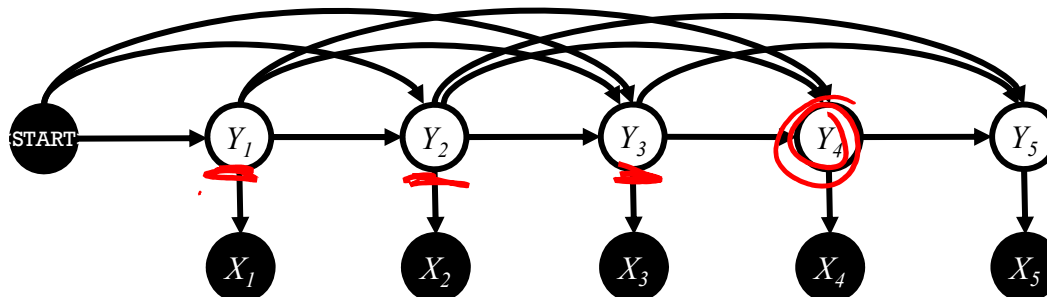  - With McCallum, IE using HMMs on CORA
- …

# Higher-order HMMs

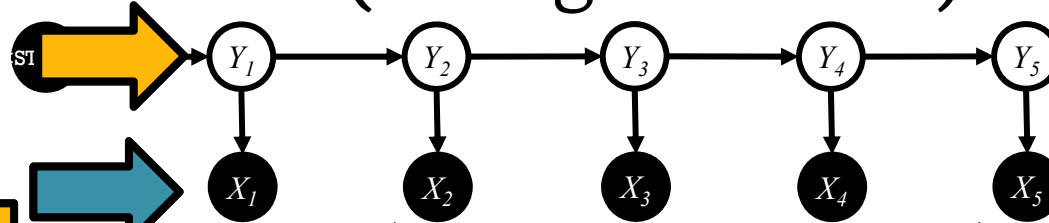- 1st-order HMM (i.e. bigram HMM)



- 2nd-order HMM (i.e. trigram HMM)



- 3rd-order HMM

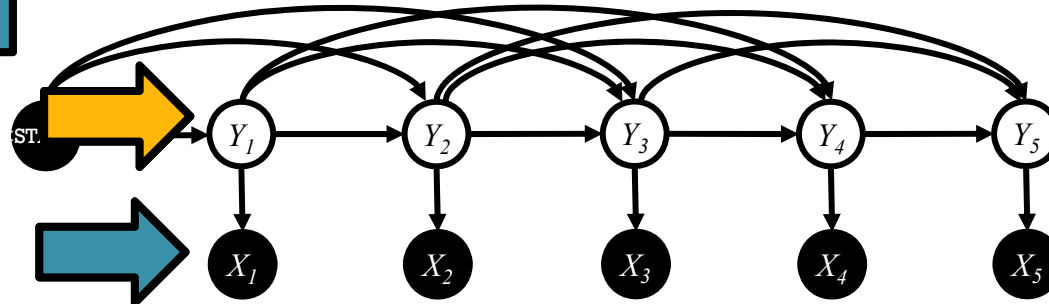# Higher-order HMMs

- 1$^{st}$-order HMM (i.e. bigram HMM)



$n^{th}$-order HMM (i.e. trigram HMM)



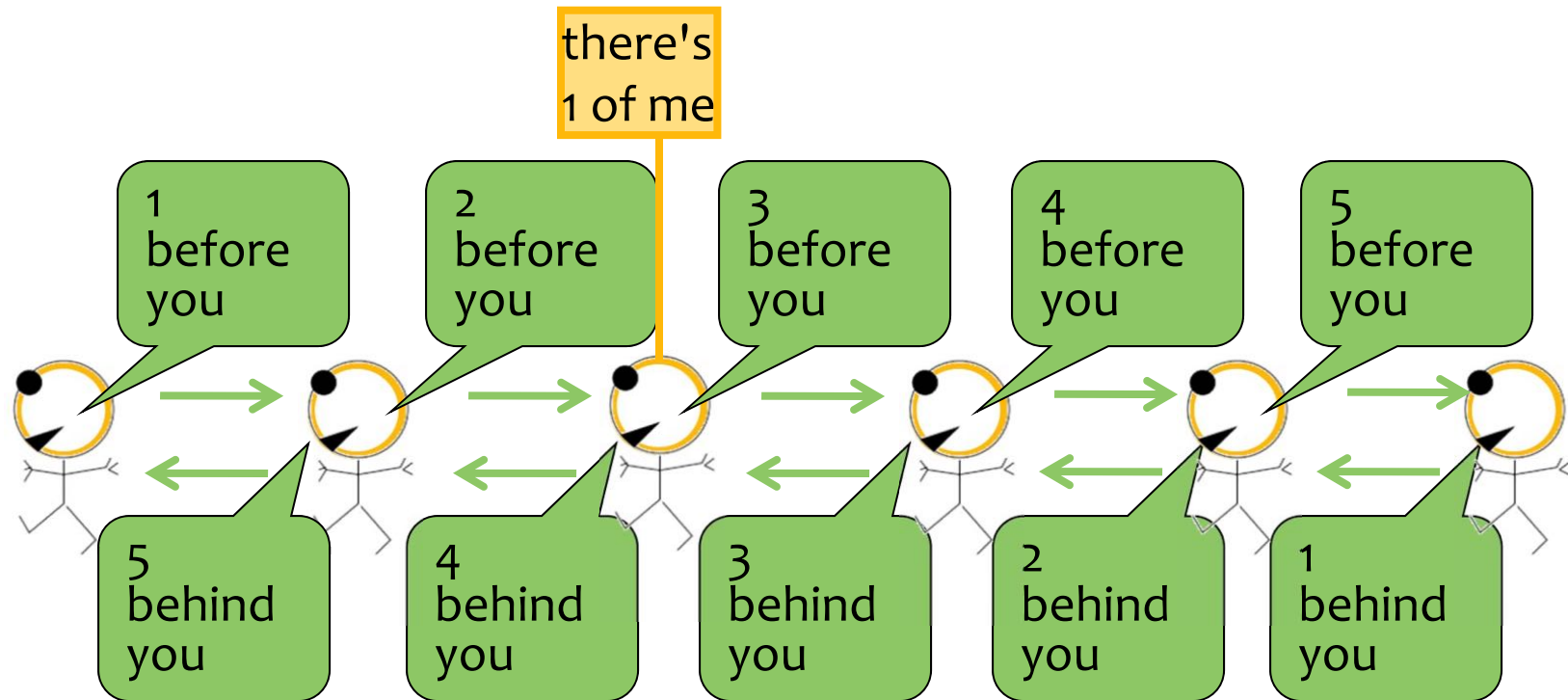Hidden States, **y**

Observa-tions, **x**

$n^{th}$-order HMM

# BACKGROUND: MESSAGE PASSING

# Great Ideas in ML: Message Passing

*Count the soldiers*

adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

*Count the soldiers*



there's
1 of me

Belief:
Must be
2 + 1 + 3 = 6 of
us

2
before
you

only see
my incoming
messages

3
behind
you

# Great Ideas in ML: Message Passing



*Count the soldiers*

adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



3 here

7 here

1 of me

11 here
(= 7+3+1)

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*

adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*

adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*

adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



3 here

7 here

3 here

Belief:
Must be
14 of us

wouldn't work correctly
with a 'loopy' (cyclic) graph

adapted from MacKay (2003) textbook

# THE FORWARD-BACKWARD ALGORITHM

# Inference

*B*    *C*    *A = calamity*

## Question:

*True or False*: The **joint probability of the observations and the hidden states** in an HMM is given by:

70%    30%

$$P(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y}) = C_{y_1} \left[ \prod_{t=1}^{T} A_{y_t, x_t} \right] \left[ \prod_{t=1}^{T-1} B_{y_{t+1}, y_t} \right]$$

*oops!*

## Recall:

Emission matrix, $\mathbf{A}$, where $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$

Transition matrix, $\mathbf{B}$, where $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

Initial probs, $\mathbf{C}$, where $P(Y_1 = k) = C_k, \forall k$

# Inference

A = calamity

**Question:**

*True or False*: The **probability of the observations** in an HMM is given by:

5%   50%

$$P(\mathbf{X} = \mathbf{x}) = \prod_{t=1}^{T} A_{x_t, x_{t-1}}$$

$$P(X = \vec{x}) = \sum_{y \in \mathcal{Y}} P(X = \vec{x}, Y = \vec{y})$$

$$= \sum_{Y_1} \sum_{Y_2} \cdots \sum_{Y_t} p(\cdot, \cdot)$$

**Recall:**

Emission matrix, $\mathbf{A}$, where $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$

Transition matrix, $\mathbf{B}$, where $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

Initial probs, $\mathbf{C}$, where $P(Y_1 = k) = C_k, \forall k$

# Inference

**Question:**

*True or False*: Suppose each hidden state takes K values. The **marginal probability** of a hidden state $\mathbf{y}_t$ given the observations $\mathbf{x}$ is given by:

$$P(Y_t = y_t | \mathbf{X} = \mathbf{x}) = \sum_{j=1}^{K} B_{j,y_t}$$

**Recall:**

Emission matrix, $\mathbf{A}$, where $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$

Transition matrix, $\mathbf{B}$, where $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

Initial probs, $\mathbf{C}$, where $P(Y_1 = k) = C_k, \forall k$

# Inference for HMMs

*Whiteboard*

– Three Inference Problems for an HMM

1. Evaluation: Compute the probability of a given sequence of observations

2. Viterbi Decoding: Find the most-likely sequence of hidden states, given a sequence of observations

3. Marginals: Compute the marginal distribution for a hidden state, given a sequence of observations

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

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

Sample 1:

| n | v | p | d | n | $\boldsymbol{y^{(1)}}$ |
| time | flies | like | an | arrow | $\boldsymbol{x^{(1)}}$ |

Sample 2:

| n | n | v | d | n | $\boldsymbol{y^{(2)}}$ |
| time | flies | like | an | arrow | $\boldsymbol{x^{(2)}}$ |

Sample 3:

| n | v | p | n | n | $\boldsymbol{y^{(3)}}$ |
| flies | fly | with | their | wings | $\boldsymbol{x^{(3)}}$ |

Sample 4:

| p | n | n | v | v | $\boldsymbol{y^{(4)}}$ |
| with | time | you | will | see | $\boldsymbol{x^{(4)}}$ |

49

# Inference for HMMs

*Whiteboard*

- Forward-backward search space

# Hidden Markov Model

A Hidden Markov Model (HMM) provides a joint distribution over the the sentence/tags with an assumption of dependence between adjacent tags.

$$p(\text{n, v, p, d, n, time, flies, like, an, arrow}) \quad = \quad (.3 * .8 * .2 * .5 * \dots)$$

|   | v | n | p | d |
|---|---|---|---|---|
| v | .1 | .4 | .2 | .3 |
| n | .8 | .1 | .1 | 0 |
| p | .2 | .3 | .2 | .3 |
| d | .2 | .8 | 0 | 0 |

|   | v | n | p | d |
|---|---|---|---|---|
| v | .1 | .4 | .2 | .3 |
| n | .8 | .1 | .1 | 0 |
| p | .2 | .3 | .2 | .3 |
| d | .2 | .8 | 0 | 0 |



|   | time | flies | like | … |
|---|---|---|---|---|
| v | .2 | .5 | .2 | |
| n | .3 | .4 | .2 | |
| p | .1 | .1 | .3 | |
| d | .1 | .2 | .1 | |

|   | time | flies | like | … |
|---|---|---|---|---|
| v | .2 | .5 | .2 | |
| n | .3 | .4 | .2 | |
| p | .1 | .1 | .3 | |
| d | .1 | .2 | .1 | |

# Forward-Backward Algorithm

# Forward-Backward Algorithm

# Forward-Backward Algorithm



- Let's show the possible *values* for each variable

# Forward-Backward Algorithm



- Let's show the possible *values* for each variable

# Forward-Backward Algorithm



- Let's show the possible *values* for each variable
- One possible assignment

# Forward-Backward Algorithm



- Let's show the possible *values* for each variable
- One possible assignment
- And what the 7 transition / emission factors think of it ...

# Forward-Backward Algorithm



|   | v | n | a |
|---|---|---|---|
| v | 1 | 6 | 4 |
| n | 8 | 4 | 0.1 |
| a | 0.1 | 8 | 0 |

|   | find | pref. | tags | ... |
|---|------|-------|------|-----|
| v | 3 | 5 | 3 | |
| n | 4 | 5 | 2 | |
| a | 0.1 | 0.2 | 0.1 | |

- Let's show the possible *values* for each variable
- One possible assignment
- And what the 7 transition / emission factors think of it ...

59

# Viterbi Algorithm: Most Probable Assignment



- So p(**v a n**) = (1/Z) * product of 7 numbers
- Numbers associated with edges and nodes of path
- Most probable assignment = **path with highest product**

# Viterbi Algorithm: Most Probable Assignment



- So p(**v a n**|$\vec{x}$) = (1/Z) * product weight of one path

# Forward-Backward Algorithm: Finds Marginals



- So p(**v a n**) = (1/Z) * product weight of one path
- Marginal probability p($Y_2$ = a)
      = (1/Z) * total weight of *all* paths through

# Forward-Backward Algorithm: Finds Marginals



- So p(**v a n**) = (1/Z) * product weight of one path
- Marginal probability p($Y_2$ = n) = (1/Z) * total weight of *all* paths through **n**

63

# Forward-Backward Algorithm: Finds Marginals



- So p(**v a n**) = (1/Z) * product weight of one path
- Marginal probability p($Y_2$ = v)
  = (1/Z) * total weight of *all* paths through

# Forward-Backward Algorithm: Finds Marginals



- So p(**v a n**) = (1/Z) * product weight of one path
- Marginal probability p($Y_2$ = n)
  = (1/Z) * total weight of *all* paths through **n**

65

# Forward-Backward Algorithm: Finds Marginals



$\alpha_2(\mathbf{n})$ = total weight of these path *prefixes* $= a + b + c$

(found by dynamic programming: matrix-vector products)

66

# Forward-Backward Algorithm: Finds Marginals



$\beta_2(\mathbf{n})$ = total weight of these path *suffixes* = x + y + z

(found by dynamic programming: matrix-vector products)

# Forward-Backward Algorithm: Finds Marginals



$\alpha_2(\mathbf{n})$ = total weight of these path *prefixes* (a + b + c)

$\beta_2(\mathbf{n})$ = total weight of these path *suffixes* (x + y + z)

more efficient

$\alpha_2(n)\,\beta_n(n) = (a+b+c)(x+y+z)$

Product gives ax+ay+az+bx+by+bz+cx+cy+cz = total weight of paths

68

# Forward-Backward Algorithm: Finds Marginals

Oops! The weight of a path through a state also includes a weight at that state.
So $\alpha(n) \cdot \beta(n)$ isn't enough.

The extra weight is the opinion of the emission probability at this variable.

$Y_2$

v

n

a

$\alpha_2(n)$

$\beta_2(n)$

"belief that $Y_2 = n$"

$A(\text{pref.}, n)$

$X_2$
preferred

total weight of *all* paths through  n

$= \quad \alpha_2(n) \quad A(\text{pref.}, n) \quad \beta_2(n)$

69

# Forward-Backward Algorithm: Finds Marginals



"belief that $Y_2 = \mathbf{v}$"

"belief that $Y_2 = \mathbf{n}$"

$\alpha_2(\mathbf{v})$

$\beta_2(\mathbf{v})$

$A(\text{pref.}, \mathbf{v})$

total weight of *all* paths through △ v

$= \alpha_2(\mathbf{v}) \quad A(\text{pref.}, \mathbf{v}) \quad \beta_2(\mathbf{v})$

# Forward-Backward Algorithm: Finds Marginals



$P(y_2 \mid \vec{x})$

| | |
|---|---|
| v | 0.1 |
| n | 0 |
| a | 0.4 |

divide by Z=0.5 to get marginal probs

| | |
|---|---|
| v | 0.2 |
| n | 0 |
| a | 0.8 |

$Y_2$

v

n

$\alpha_2(a)$     $\beta_2(a)$

a

A(pref., a)

$X_2$ preferred

"belief that $Y_2 = \mathbf{v}$"

"belief that $Y_2 = \mathbf{n}$"

"belief that $Y_2 = \mathbf{a}$"

sum = $Z$ (total weight of *all* paths)

$P(\vec{x}) = \sum_{\vec{y}} P(\vec{x}, \vec{y})$

total weight of *all* paths through ▲ a

= $\alpha_2(a)$   A(pref., a)   $\beta_2(a)$

71

# Forward-Backward Algorithm



$Y_1$    $Y_2$    $Y_3$

$X_1$    $X_2$    $X_3$

find    preferred    tags

*Could be verb or noun*    *Could be adjective or verb*    *Could be noun or verb*

# Inference for HMMs

*Whiteboard*

- Derivation of Forward algorithm
- Forward-backward algorithm
- Viterbi algorithm

# Forward-Backward Algorithm

Define:
$$\alpha_t(k) \triangleq p(x_1, \ldots, x_t, y_t = k)$$
$$\beta_t(k) \triangleq p(x_{t+1}, \ldots, x_T \mid y_t = k)$$

Assume:
$$y_0 = \text{START}$$
$$y_{T+1} = \text{END}$$

① Initialize
$$\alpha_0(\text{START}) = 1 \qquad \alpha_0(k) = 0 \quad \forall k \neq \text{START}$$
$$\beta_{T+1}(\text{END}) = 1 \qquad \beta_{T+1}(k) = 0 \quad \forall k \neq \text{END}$$

**forward algo**

② For $t = 1, \ldots, T$:

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

*the alphas include the emission probabilities so we don't multiply them in separately*

$$\alpha_t(k) = p(x_t \mid y_t = k) \sum_{j=1}^{K} \alpha_{t-1}(j)\, p(y_t = k \mid y_{t-1} = j)$$

**backward algo**

③ For $t = T, \ldots, 1$:

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

$$\beta_t(k) = \sum_{j=1}^{K} p(x_{t+1} \mid y_{t+1} = j)\, \beta_{t+1}(j)\, p(y_{t+1} = j \mid y_t = k)$$

**Eval.** → ④ Compute $p(\vec{x}) = \alpha_{T+1}(\text{END})$    [Evaluation]

**Marg.** → ⑤ Compute $p(y_t = k \mid \vec{x}) = \dfrac{\alpha_t(k)\, \beta_t(k)}{p(\vec{x})}$    [Marginals]

# Derivation of Forward Algorithm

Definition: $\alpha_t(k) \triangleq p(x_1, \ldots, x_t, y_t = k)$

Derivation:

$\alpha_T(\text{END}) = p(x_1, \ldots, x_T, y_T = \text{END})$

$= p(x_1, \ldots, x_T \mid y_T) \, p(y_T)$    ← by def of joint

$= p(x_T \mid y_T) \, p(x_1, \ldots, x_{T-1} \mid y_T) \, p(y_T)$    ← by cond. indep. of HMM

$= p(x_T \mid y_T) \, p(x_1, \ldots, x_{T-1}, y_T)$    ← by def. of joint

$= p(x_T \mid y_T) \sum_{y_{T-1}} p(x_1, \ldots, x_{T-1}, y_{T-1}, y_T)$    ← by def. of marginal

$= p(x_T \mid y_T) \sum_{y_{T-1}} p(x_1, \ldots, x_{T-1}, y_T \mid y_{T-1}) \, p(y_{T-1})$    ← by def. of joint

$= p(x_T \mid y_T) \sum_{y_{T-1}} p(x_1, \ldots, x_{T-1} \mid y_{T-1}) \, p(y_T \mid y_{T-1}) \, p(y_{T-1})$    ← by cond. indep. of HMM

$= p(x_T \mid y_T) \sum_{y_{T-1}} p(x_1, \ldots, x_{T-1}, y_{T-1}) \, p(y_T \mid y_{T-1})$    ← by def. of joint

$= p(x_T \mid y_T) \sum_{y_{T-1}} \alpha_{T-1}(y_{T-1}) \, p(y_T \mid y_{T-1})$    ← by def. of $\alpha_t(k)$

Herein using "$y_T$" as shorthand for "$y_T = \text{END}$"

# Viterbi Algorithm

<u>Define</u>:  $\omega_t(k) \triangleq \max\limits_{y_1,\ldots,y_{t-1}} p(x_1,\ldots,x_t, y_1,\ldots,y_{t-1}, y_t = k)$

"back pointers" ⟶  $b_t(k) \triangleq \arg\max\limits_{y_1,\ldots,y_{t-1}} p(x_1,\ldots,x_t, y_1,\ldots,y_{t-1}, y_t = k)$

<u>Assume</u>  $y_0 = \text{START}$

① Initialize  $\omega_0(\text{START}) = 1$   $\omega_0(k) = 0 \ \forall k \neq \text{START}$

② For $t = 1,\ldots,T$:

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

$$\omega_t(k) = \max\limits_{j \in \{1,\ldots,K\}} p(x_t | y_t = k)\, \omega_{k-1}(j)\, p(y_t = k | y_{t-1} = j)$$

$$b_t(k) = \arg\max\limits_{j \in \{1,\ldots,K\}} p(x_t | y_t = k)\, \omega_{k-1}(j)\, p(y_t = k | y_{t-1} = j)$$

③ Compute Most Probable Assignment          [Decoding]

$\hat{y}_T = b_{T+1}(\text{END})$

For $t = T-1,\ldots,1$

$\hat{y}_t = b_{t+1}(\hat{y}_{t+1})$          ] follow the "back pointers"

76

# 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**, $O(K^T)$

- The **forward-backward** algorithm and **Viterbi** algorithm run in **polynomial time**, $O(T*K^2)$
  - Thanks to dynamic programming!

# Shortcomings of Hidden Markov Models



- HMM models capture dependences between each state and <span style="color:red">only</span> 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 (non-local) 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 $P(\mathbf{Y}, \mathbf{X})$, but in a prediction task, we need the conditional probability $P(\mathbf{Y}|\mathbf{X})$

# MBR DECODING

# Inference for HMMs

*Four*

– ~~Three~~ Inference Problems for an HMM

1.  Evaluation: Compute the probability of a given sequence of observations

2.  Viterbi Decoding: Find the most-likely sequence of hidden states, given a sequence of observations

3.  Marginals: Compute the marginal distribution for a hidden state, 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

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

$$= \operatorname*{argmin}_{\hat{\boldsymbol{y}}} \ \sum_{\boldsymbol{y}} p_{\boldsymbol{\theta}}(\boldsymbol{y} \mid \boldsymbol{x}) \ell(\hat{\boldsymbol{y}}, \boldsymbol{y})$$

# Minimum Bayes Risk Decoding

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

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:

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

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

which is exactly the Viterbi decoding problem!

# Minimum Bayes Risk Decoding

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

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:

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

This decomposes across variables and requires the variable marginals.

# Learning Objectives

**Hidden Markov Models**

*You should be able to...*

1. 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