



10-301/601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

Hidden Markov Models (Part II)

Matt Gormley Lecture 19 Nov. 7, 2022

Reminders

- Practice Problems: Exam 2
 - Out: Fri, Nov. 4
- Exam 2
 - Thu, Nov. 10, 6:30pm 8:30pm
- Homework 7: Hidden Markov Models
 - Out: Fri, Nov. 11
 - Due: Mon, Nov. 21 at 11:59pm

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

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

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

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 $x \in \{1, ..., M\}$ 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):

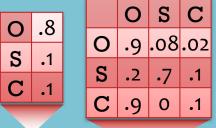
$$\phi_m^{MLE} = \frac{N_{x=m}}{N} = \frac{\sum_{i=1}^{N} \mathbb{I}(x^{(i)} = m)}{N}$$

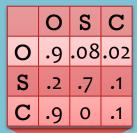


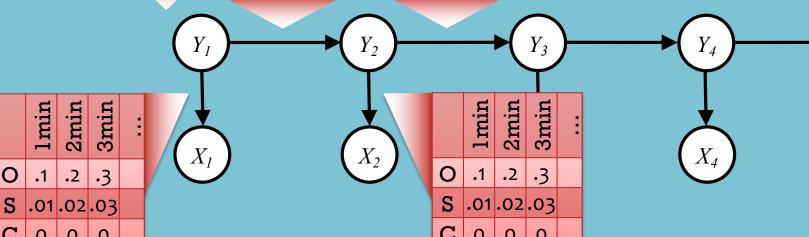
Hidden Markov Model (v1)

HMM Parameters:

Emission matrix, **A**, where $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$ Transition matrix, **B**, where $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$ Initial probs, **C**, where $P(Y_1 = k) = C_k, \forall k$







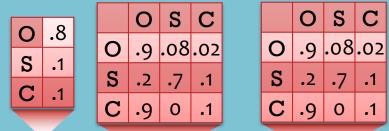
$$P(\mathbf{X}, \mathbf{Y}) = P(Y_1) \left(\prod_{t=1}^{T} P(X_t | Y_t) \right) \left(\prod_{t=2}^{T} p(Y_t | Y_{t-1}) \right)$$

 X_5

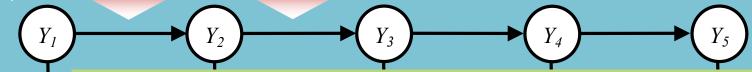
Hidden Markov Model (v1)

HMM Parameters:

Emission matrix, A, where $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$ Transition matrix, **B**, where $P(Y_t = k | Y_{t-1} = j) = B_{i,k}, \forall t, k$ Initial probs, C, where $P(Y_1 = k) = C_k, \forall k$







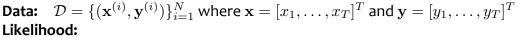
2min .01.02.03 Joint Distribution (probability mass function):

$$(X_t) p(\mathbf{x}, \mathbf{y}) = p(y_1, C) \left(\prod_{t=1}^T p(x_t \mid y_t, A) \right) \left(\prod_{t=2}^T p(y_t \mid y_{t-1}, B) \right)$$

$$= C_{y_1} \left(\prod_{t=1}^{T} A_{y_t, x_t} \right) \left(\prod_{t=2}^{T} B_{y_{t-1}, y_t} \right)$$

Supervised Learning for HMM (v1)

Learning an HMMdecomposes into solving two (independent) Mixture Models



$$\begin{split} \ell(\mathbf{A}, \mathbf{B}, \mathbf{C}) &= \sum_{i=1}^{N} \log p(\mathbf{x}^{(i)}, \mathbf{y}^{(i)} \mid \mathbf{A}, \mathbf{B}, \mathbf{C}) \\ &= \sum_{i=1}^{N} \left[\underbrace{\log p(y_1^{(i)} \mid \mathbf{C})}_{\text{initial}} + \underbrace{\left(\sum_{t=2}^{T} \log p(y_t^{(i)} \mid y_{t-1}^{(i)}, \mathbf{B})\right)}_{\text{transition}} + \underbrace{\left(\sum_{t=1}^{T} \log p(x_t^{(i)} \mid y_t^{(i)}, \mathbf{A})\right)}_{\text{emission}} \right] \end{split}$$

MLE:

$$\hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\mathbf{C}} = \underset{\mathbf{A}, \mathbf{B}, \mathbf{C}}{\operatorname{argmax}} \ell(\mathbf{A}, \mathbf{B}, \mathbf{C})$$

$$\Rightarrow \hat{\mathbf{C}} = \underset{\mathbf{C}}{\operatorname{argmax}} \sum_{i=1}^{N} \log p(y_1^{(i)} \mid \mathbf{C})$$

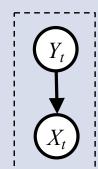
$$\hat{\mathbf{B}} = \underset{\mathbf{B}}{\operatorname{argmax}} \sum_{i=1}^{N} \sum_{t=2}^{T} \log p(y_t^{(i)} \mid y_{t-1}^{(i)}, \mathbf{B})$$

$$\hat{\mathbf{A}} = \underset{\mathbf{A}}{\operatorname{argmax}} \sum_{i=1}^{N} \sum_{t=1}^{T} \log p(x_t^{(i)} \mid y_t^{(i)}, \mathbf{A})$$

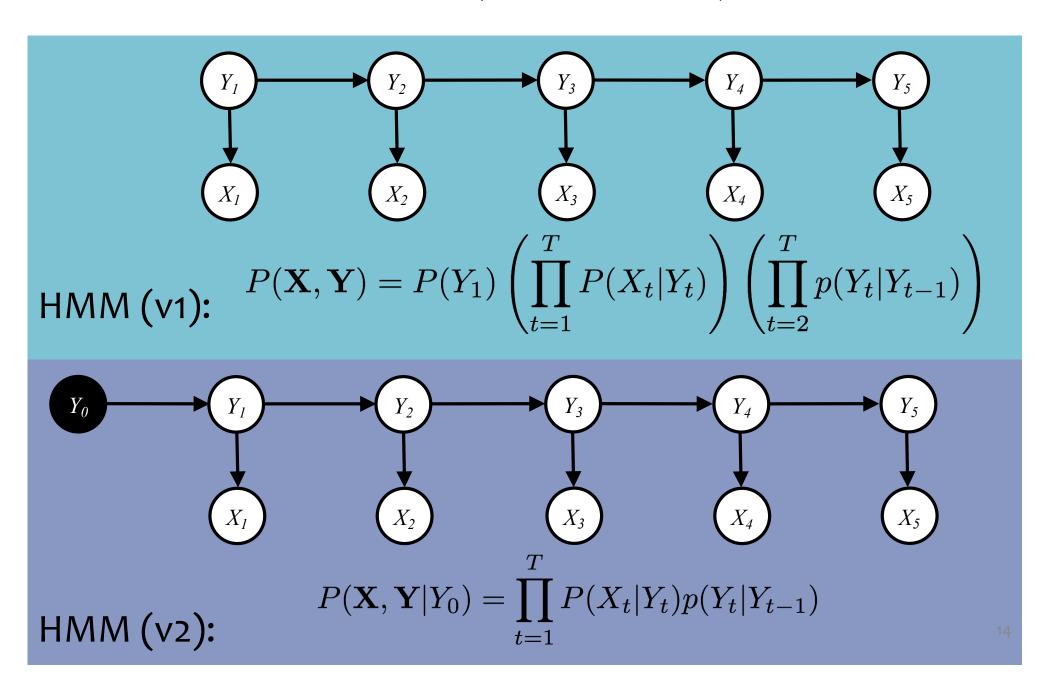
We can solve the above in closed form, which yields...

$$\hat{C}_k = rac{\#(y_1^{(i)} = k)}{N}, \, orall k$$
 $\hat{B}_{j,k} = rac{\#(y_t^{(i)} = k ext{ and } y_{t-1}^{(i)} = j)}{\#(y_{t-1}^{(i)} = j)}, \, orall j, k$
 $\hat{A}_{j,k} = rac{\#(x_t^{(i)} = k ext{ and } y_t^{(i)} = j)}{\#(y_t^{(i)} = j)}, \, orall j, k$

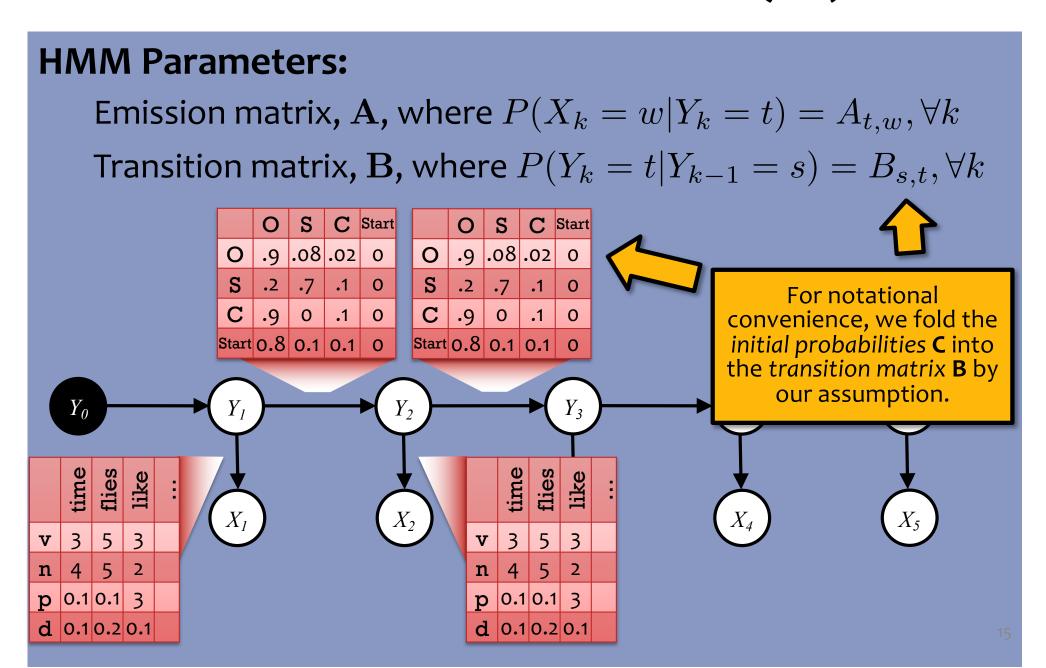




HMM (two ways)



Hidden Markov Model (v2)



Hidden Markov Model (v2)

HMM Parameters:

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

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

Assumption: $y_0 = START$

Generative Story:

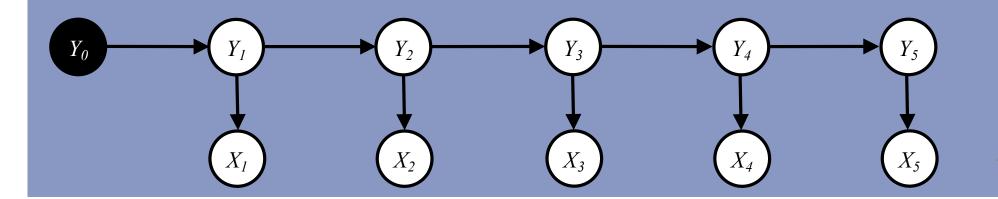
 $Y_t \sim \text{Multinomial}(\mathbf{B}_{Y_{t-1}}) \ \forall t$

 $X_t \sim \text{Multinomial}(\mathbf{A}_{Y_t}) \ \forall t$





For notational convenience, we fold the initial probabilities **C** into the transition matrix **B** by our assumption.



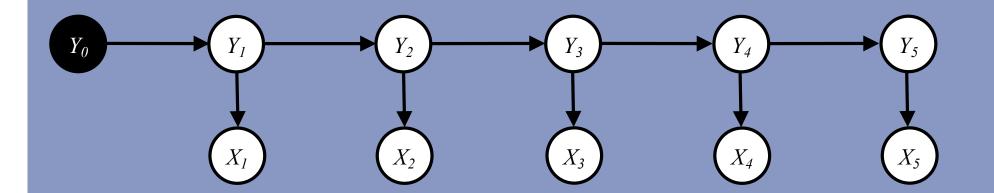
Hidden Markov Model (v2)

Joint Distribution (probability mass function):

$$y_0 = \mathsf{START}$$

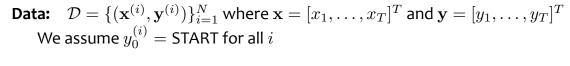
$$p(\mathbf{x}, \mathbf{y}|y_0) = \prod_{t=1}^{1} p(x_t|y_t)p(y_t|y_{t-1})$$

$$= \prod_{t=1}^{I} A_{y_t, x_t} B_{y_{t-1}, y_t}$$



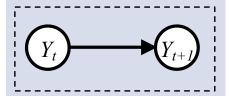
Supervised Learning for HMM (v2)

Learning an HMM decomposes into solving two (independent) Mixture Models



Likelihood:

$$\begin{split} \ell(\mathbf{A}, \mathbf{B}) &= \sum_{i=1}^{N} \log p(\mathbf{x}^{(i)}, \mathbf{y}^{(i)} \mid \mathbf{A}, \mathbf{B}) \\ &= \sum_{i=1}^{N} \left[\sum_{t=1}^{T} \underbrace{\log p(y_t^{(i)} \mid y_{t-1}^{(i)}, \mathbf{B})}_{\text{transition}} + \underbrace{\log p(x_t^{(i)} \mid y_t^{(i)}, \mathbf{A})}_{\text{emission}} \right] \end{split}$$



 Y_t X_t

MLE:

$$\hat{\mathbf{A}}, \hat{\mathbf{B}} = \underset{\mathbf{A}}{\operatorname{argmax}} \ell(\mathbf{A}, \mathbf{B})$$

$$\Rightarrow \hat{\mathbf{B}} = \underset{\mathbf{B}}{\operatorname{argmax}} \sum_{i=1}^{N} \sum_{t=1}^{T} \log p(y_t^{(i)} \mid y_{t-1}^{(i)}, \mathbf{B})$$

$$\hat{\mathbf{A}} = \underset{\mathbf{A}}{\operatorname{argmax}} \sum_{i=1}^{N} \sum_{t=1}^{T} \log p(x_t^{(i)} \mid y_t^{(i)}, \mathbf{A})$$

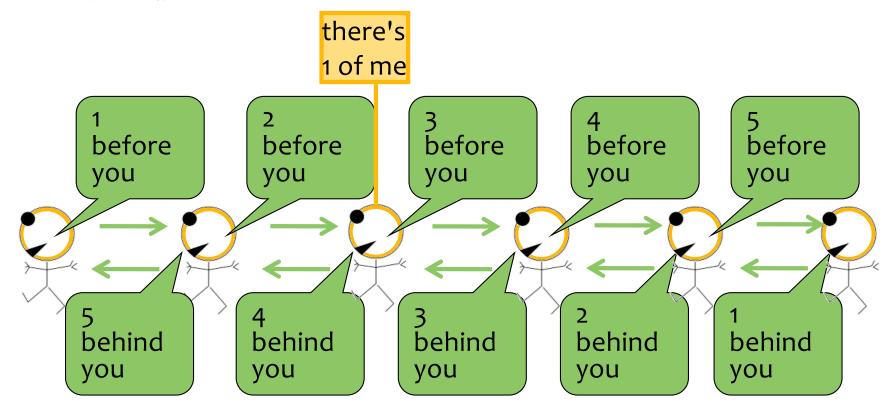
We can solve the above in closed form, which yields...

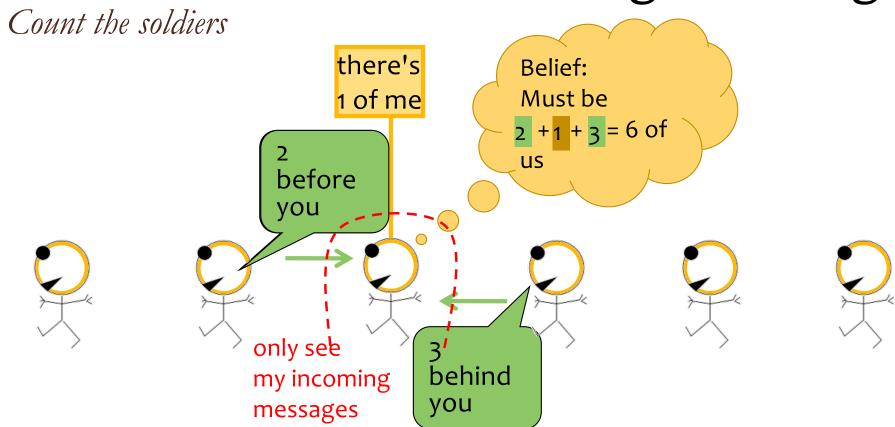
$$\hat{B}_{j,k} = \frac{\#(y_t^{(i)} = k \text{ and } y_{t-1}^{(i)} = j)}{\#(y_{t-1}^{(i)} = j)}, \forall j, k$$

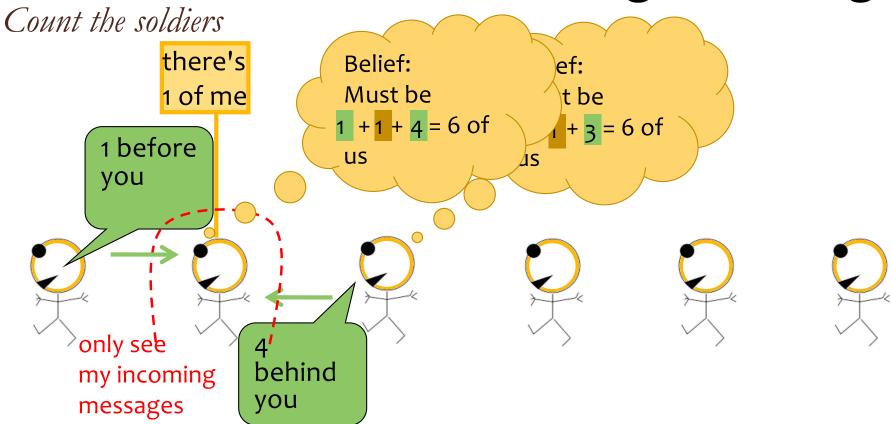
$$\hat{A}_{j,k} = \frac{\#(x_t^{(i)} = k \text{ and } y_t^{(i)} = j)}{\#(y_t^{(i)} = j)}, \forall j, k$$

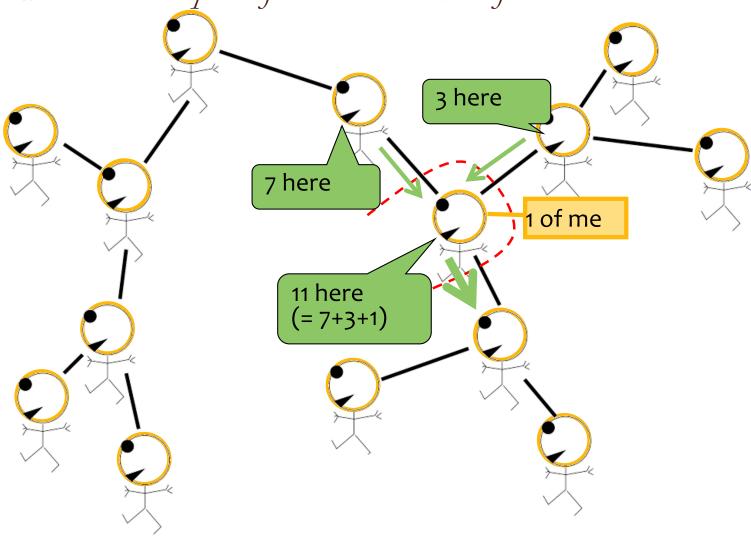
BACKGROUND: MESSAGE PASSING

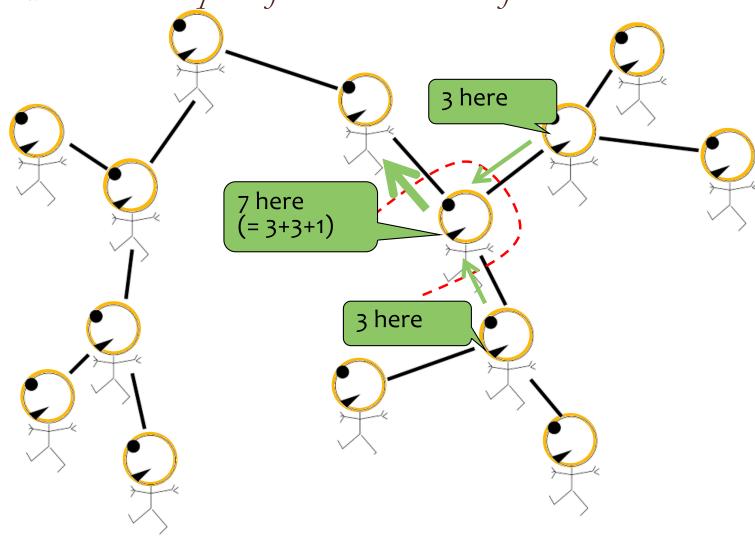
Count the soldiers

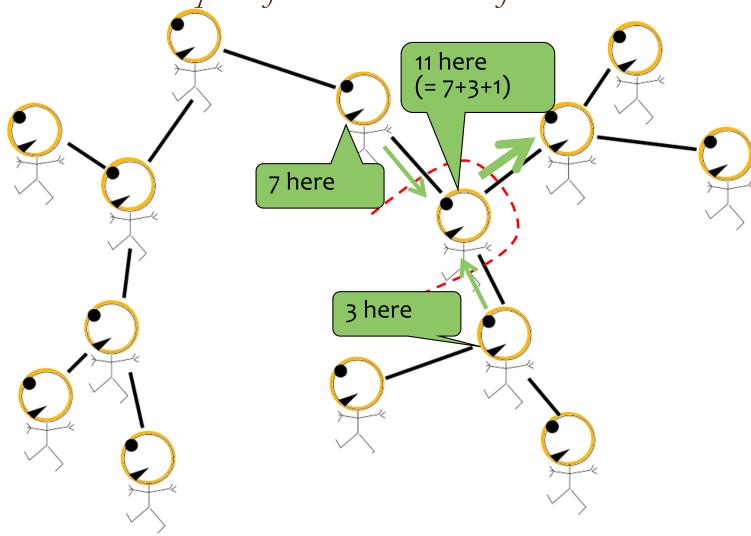


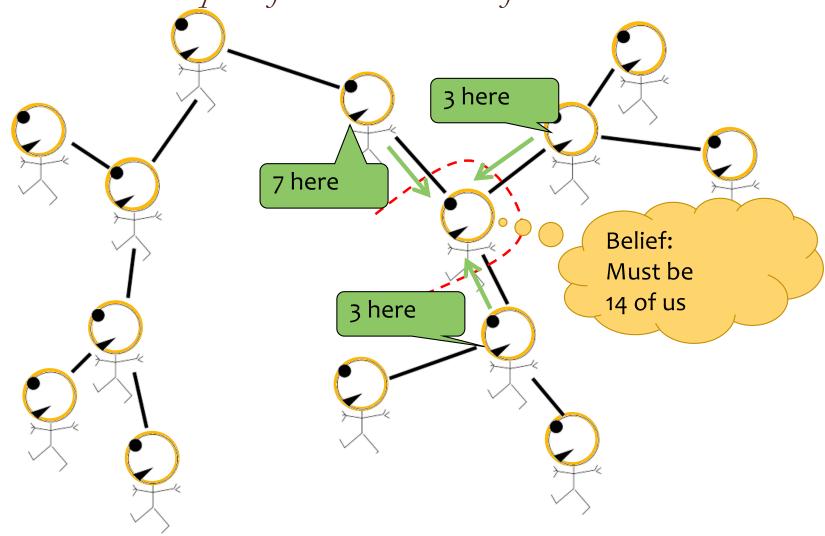


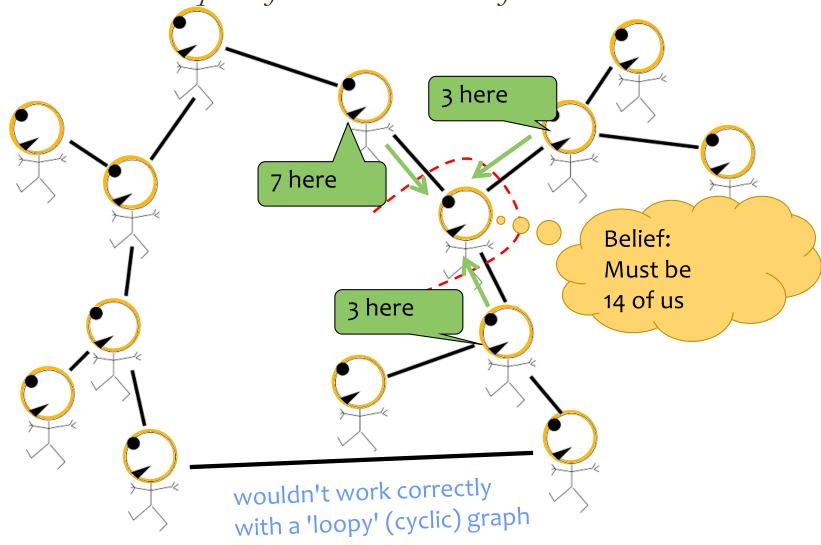












INFERENCE FOR HMMS

Inference

Question:

01

A= True B=toxic C= False 75% 2.3% -22%

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

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

Recall:

Emission matrix, A, where $P(X_t = k | Y_t = j) = A_{i,k}, \forall t, k$ Transition matrix, **B**, where $P(Y_t = k | Y_{t-1} = j) = B_{i,k}, \forall t, k$ Initial probs, C, where $P(Y_1 = k) = C_k, \forall k$

Inference

Question: Q2 A=True B=toxic C=False

True or False: The probability of the observations

in an HMM is given by:

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

Recall:

Emission matrix, A, where $P(X_t = k | Y_t = j) = A_{i,k}, \forall t, k$ Transition matrix, **B**, where $P(Y_t = k | Y_{t-1} = j) = B_{i,k}, \forall t, k$ Initial probs, C, where $P(Y_1 = k) = C_k, \forall k$

Inference for HMMs

Whiteboard

- Three Inference Problems for an HMM
 - 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

THE SEARCH SPACE FOR FORWARD-BACKWARD

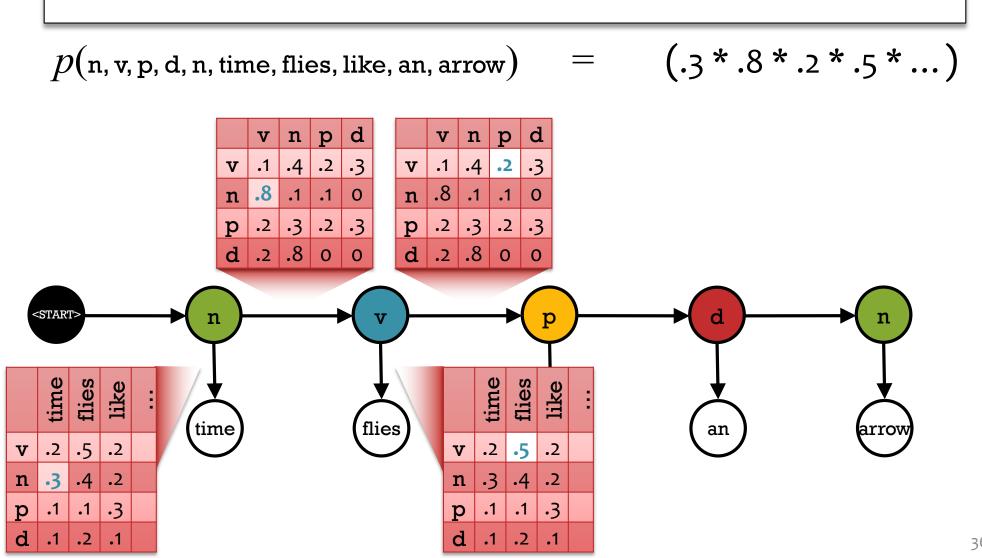
Dataset for Supervised Part-of-Speech (POS) Tagging

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

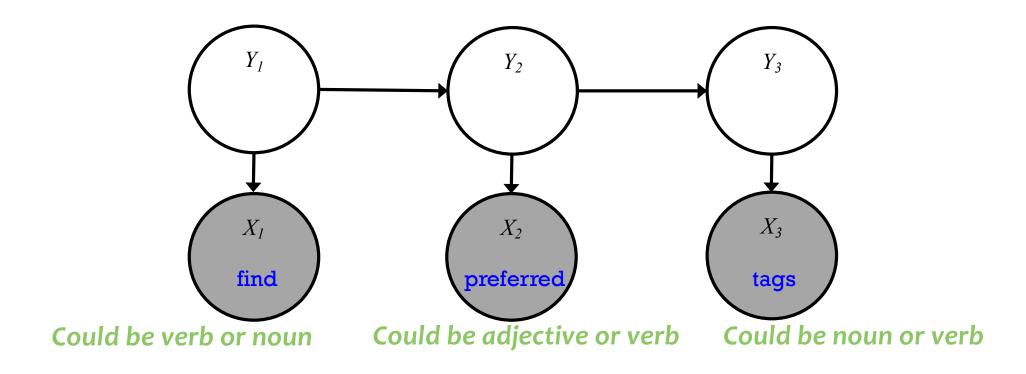
Sample 1:	n	v flies	p like	an	$\begin{array}{c c} & y \\ \hline & x \\ \end{array}$	
Sample 2:	n	n	like	an	$\begin{array}{c c} & & \\ & & \\ \hline & & \\ & & \\ \end{array}$	
Sample 3:	n	fly	with	heir	$\begin{array}{c c} & & \\ & &$	
Sample 4:	with	n	you	will	$\begin{array}{c c} & & \\ & &$	

Example: HMM for POS Tagging

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



Example: HMM for POS Tagging



Inference for HMMs

Whiteboard

- Brute Force Evaluation
- Forward-backward search space

HOW IS EFFICIENT COMPUTATION EVEN POSSIBLE?

How is efficient computation even possible?

- The short answer is dynamic programming!
- The key idea is this:
 - We first come up with a recursive definition for the quantity we want to compute
 - We then observe that many of the recursive intermediate terms are reused across timesteps and tags
 - We then perform bottom-up dynamic programming by running the recursion in reverse, storing the intermediate quantities along the way!
- This enables us to search the exponentially large space in polynomial time!

Derivation of Forward Algorithm

THE FORWARD-BACKWARD ALGORITHM

Inference for HMMs

Whiteboard

Forward-backward algorithm (edge weights version)

Forward-Backward Algorithm

Definitions

$$\alpha_t(k) \triangleq p(x_1, \dots, x_t, y_t = k)$$

$$\beta_t(k) \triangleq p(x_{t+1}, \dots, x_T \mid y_t = k)$$

Assume

$$y_0 = \mathsf{START}$$

$$y_{T+1} = \mathsf{END}$$

1. Initialize

$$\alpha_0(\mathsf{START}) = 1 \qquad \qquad \alpha_0(k) = 0, \ \forall k \neq \mathsf{START}$$

$$\beta_{T+1}(\mathsf{END}) = 1 \qquad \qquad \beta_{T+1}(k) = 0, \ \forall k \neq \mathsf{END}$$

2. Forward Algorithm

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

3. Backward Algorithm

for
$$t = T, ..., 1$$
:
for $k = 1, ..., K$:

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

- 4. Evaluation $p(\mathbf{x}) = \alpha_{T+1}(\mathsf{END})$
- 5. Marginals $p(y_t = k \mid \mathbf{x}) = \frac{\alpha_t(k)\beta_t(k)}{p(\mathbf{x})}$

Forward-Backward Algorithm

Definitions

$$\alpha_t(k) \triangleq p(x_1, \dots, x_t, y_t = k)$$

$$\beta_t(k) \triangleq p(x_{t+1}, \dots, x_T \mid y_t = k)$$

Assume

$$y_0 = \mathsf{START}$$

$$y_{T+1} = \mathsf{END}$$

1. Initialize

$$lpha_0({\sf START}) = 1$$
 $lpha_0(k) = 0, \, \forall k
eq {\sf START}$ $eta_{T+1}({\sf END}) = 1$ $eta_{T+1}(k) = 0, \, \forall k
eq {\sf END}$

2. Forward Algorithm

for
$$t = 1, ..., T$$
:

for
$$k = 1, ..., K$$
:

 $O(K^2T)$

O(K) kward Algorithm

for
$$t = T, ..., 1$$
:

Brute force $O(K^T)$

or
$$k = 1, ..., K$$
:

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

- 4. Evaluation $p(\mathbf{x}) = \alpha_{T+1}(\mathsf{END})$
- 5. Marginals $p(y_t = k \mid \mathbf{x}) = \frac{\alpha_t(k)\beta_t(k)}{n(\mathbf{x})}$