

#### 10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

# **Hidden Markov Models**

Matt Gormley Lecture 20 Mar. 30, 2020

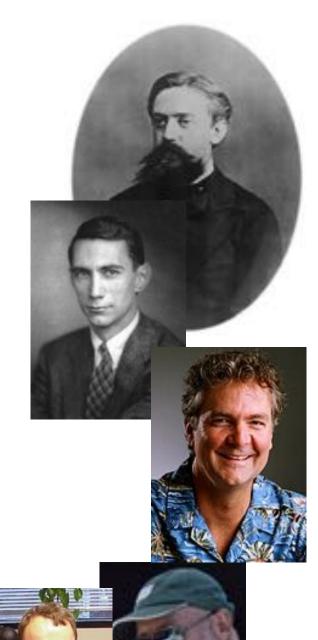
#### Reminders

- 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

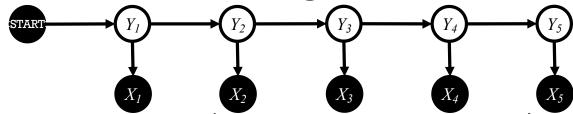
- Markov chains: Andrey Markov (1906)
  - Random walks and Brownian motion
- Used in Shannon's work on information theory (1948)
- Baum-Welsh learning algorithm: late 60's, early 70's.
  - Used mainly for speech in 60s-70s.
- Late 80's and 90's: David Haussler (major player in learning theory in 80's) 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

• ...

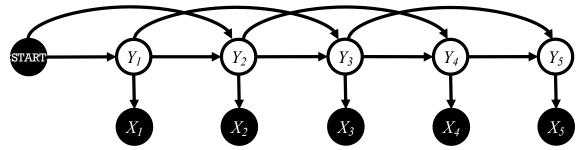


## Higher-order HMMs

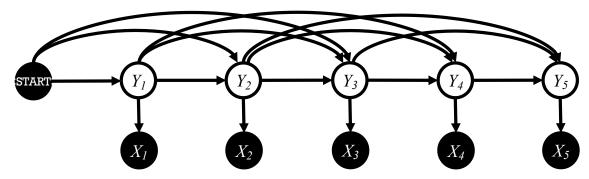
• 1st-order HMM (i.e. bigram HMM)



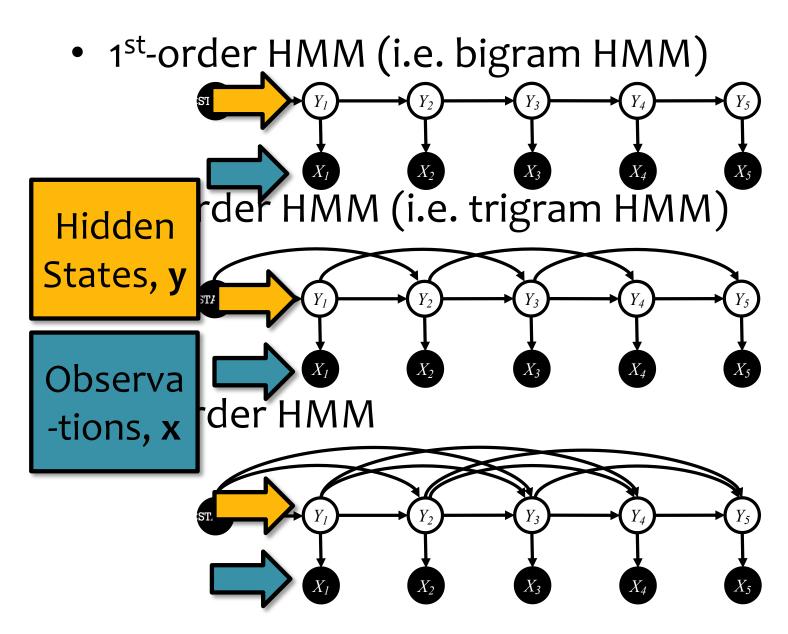
• 2<sup>nd</sup>-order HMM (i.e. trigram HMM)



• 3<sup>rd</sup>-order HMM

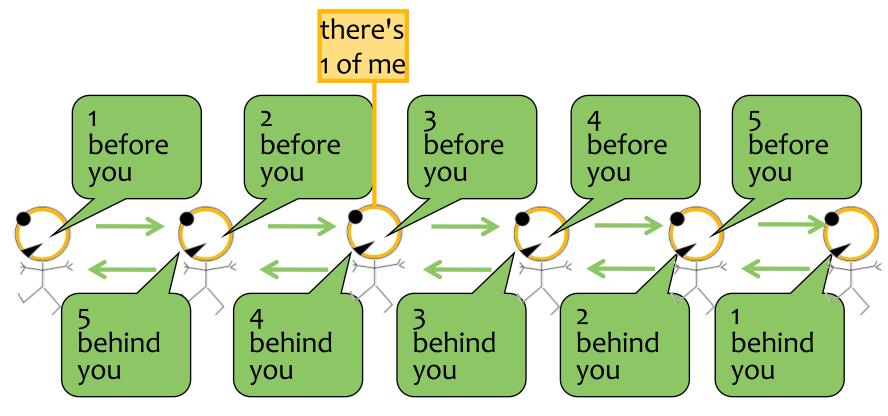


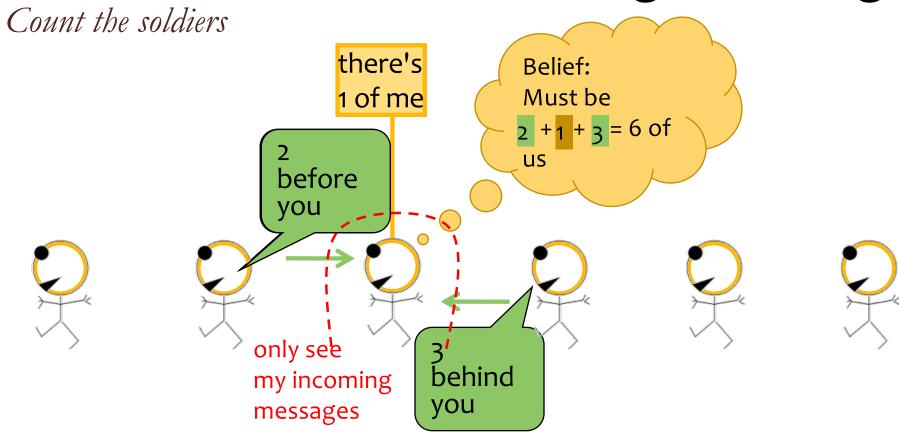
## Higher-order HMMs

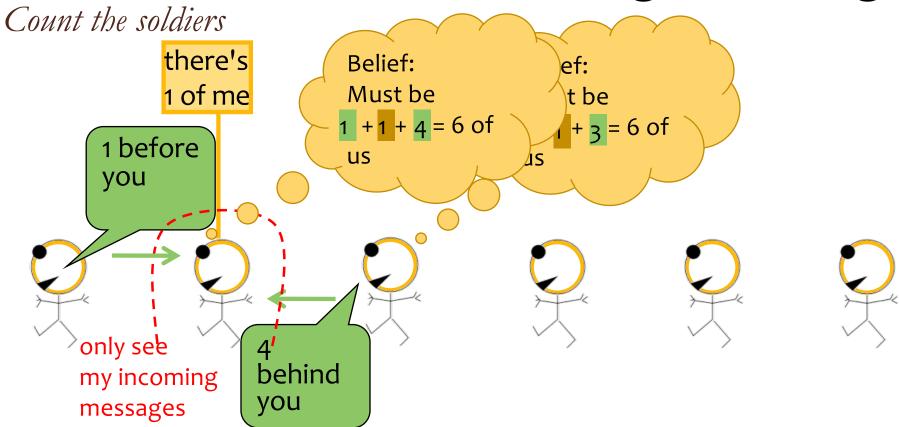


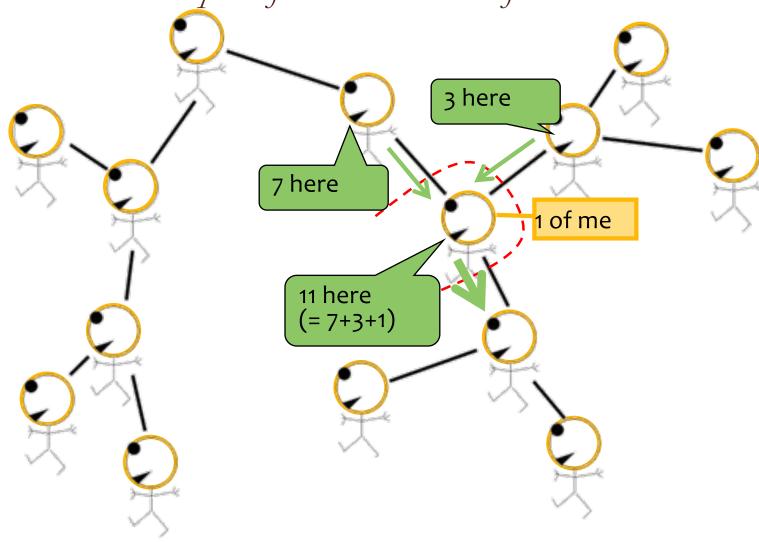
#### **BACKGROUND: MESSAGE PASSING**

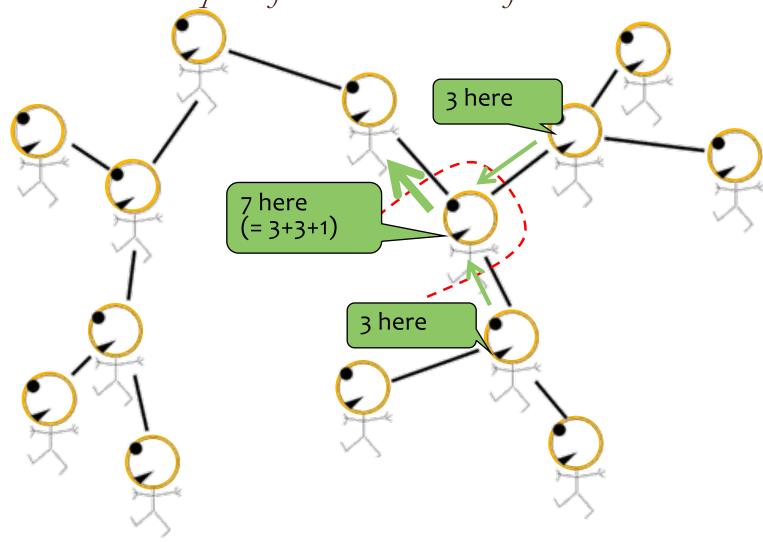
Count the soldiers

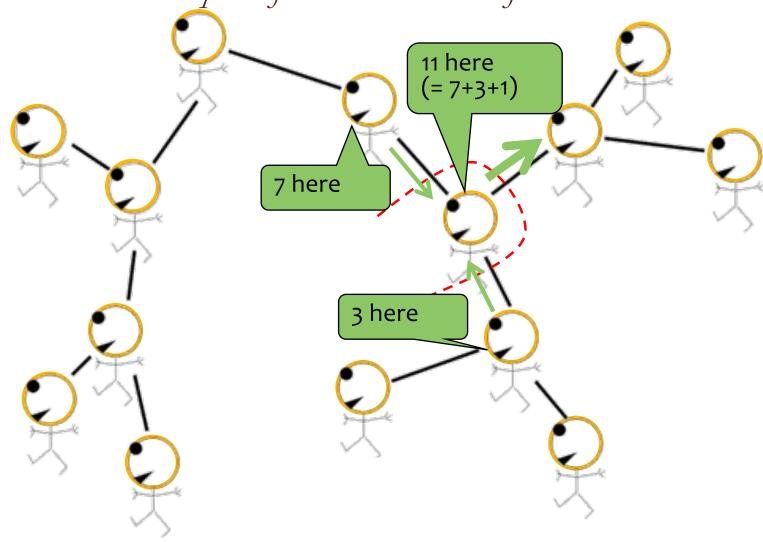


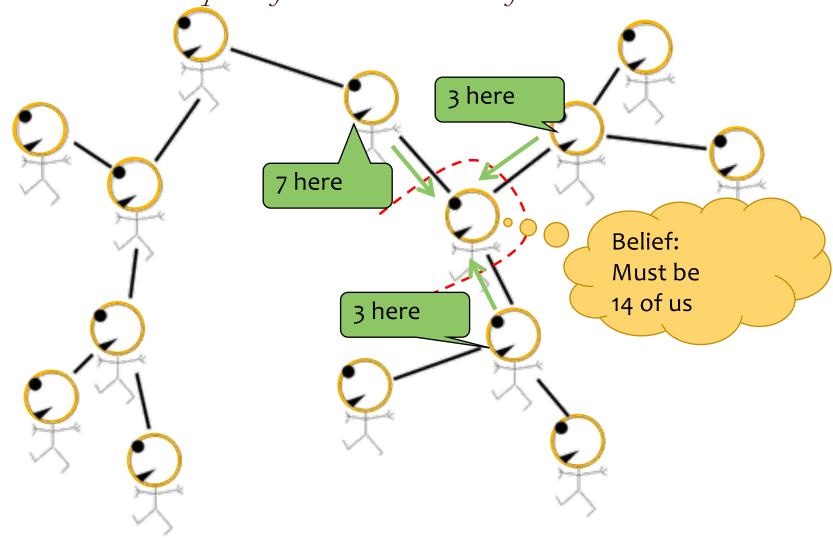


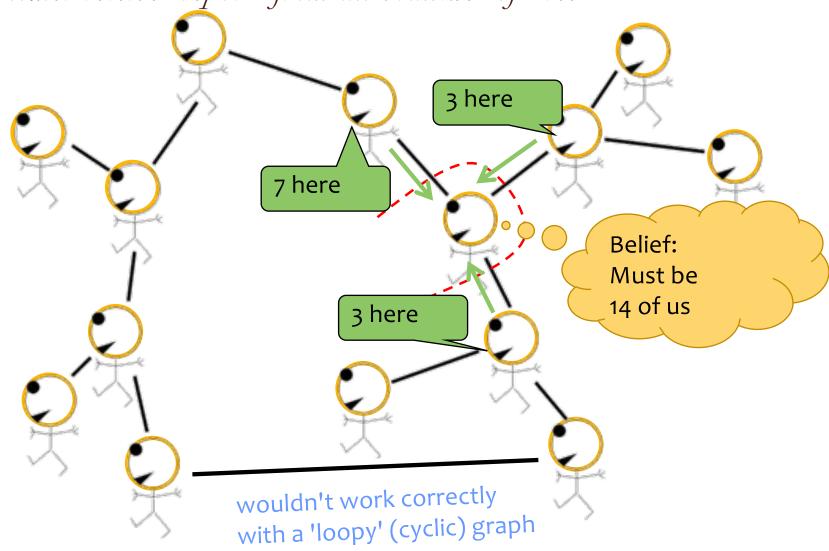












# THE FORWARD-BACKWARD ALGORITHM

## Inference

#### **Question:**

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

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

#### **Recall:**

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  $P(Y_1 = k) = C_k, \forall k$ 

### Inference

#### **Question:**

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

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

#### **Recall:**

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  $P(Y_1 = k) = C_k, \forall k$ 

## Inference

#### **Question:**

True or False: Suppose each hidden state takes K values. The marginal probability of a hidden state  $y_t$  given the observations x is given by:

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

#### **Recall:**

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  $P(Y_1 = k) = C_k, \forall k$ 

#### 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

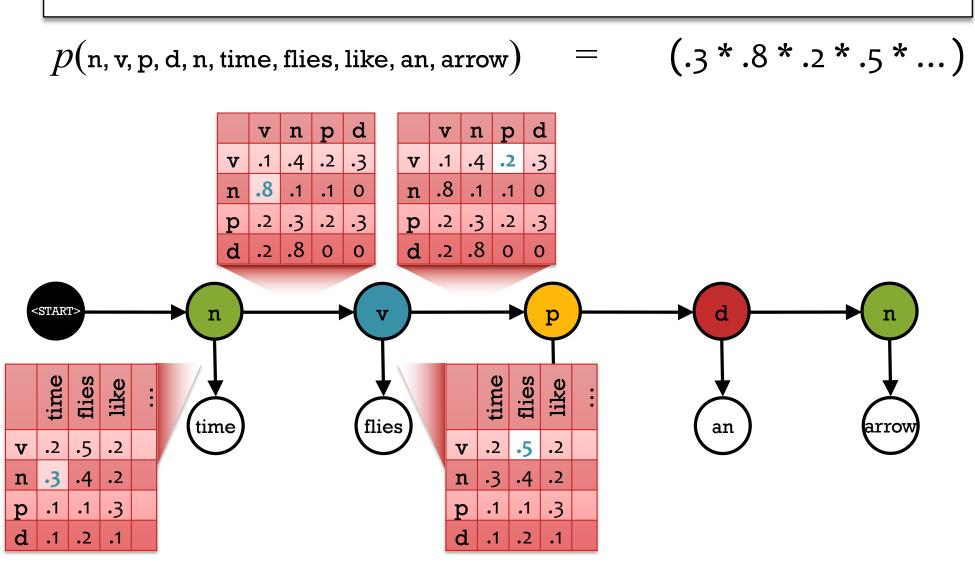
# Dataset for Supervised Part-of-Speech (POS) Tagging

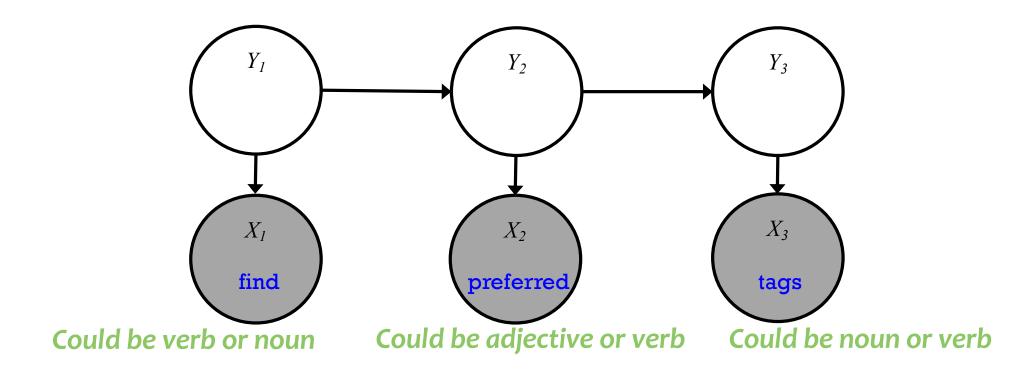
Data:  $\mathcal{D} = \{oldsymbol{x}^{(n)}, oldsymbol{y}^{(n)}\}_{n=1}^N$ 

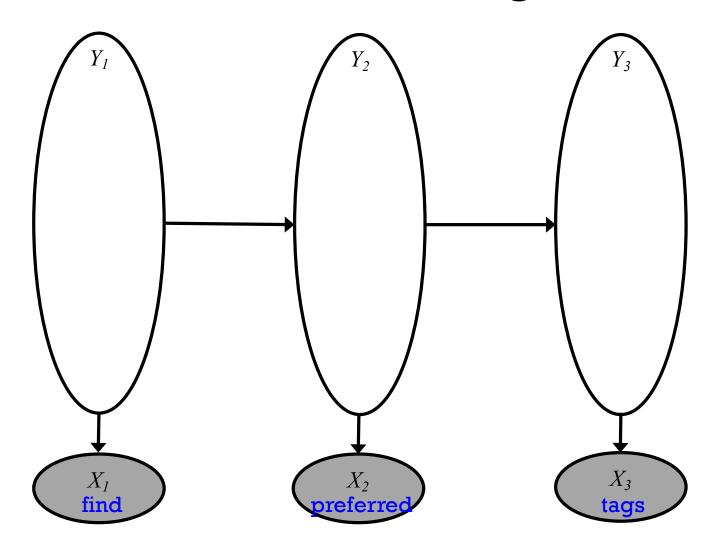
Sample 1:	n	flies	p like	dan	$ \begin{array}{c c} & & \\ & $
Sample 2:	n	n	like	d	$\begin{array}{c c}                                    $
Sample 3:	n	fily	with	heir	$ \begin{array}{c c}                                    $
Sample 4:	with	n	you	will	

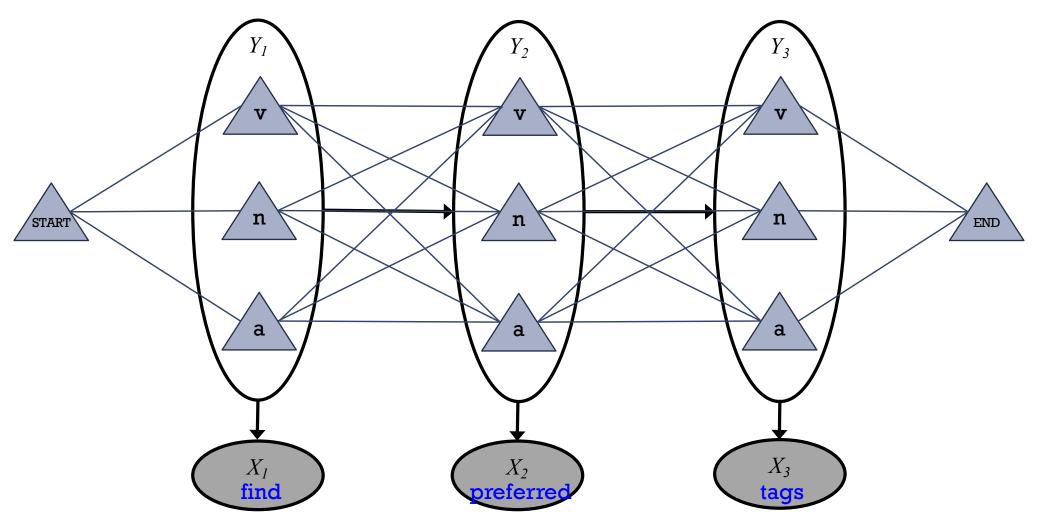
#### Hidden Markov Model

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

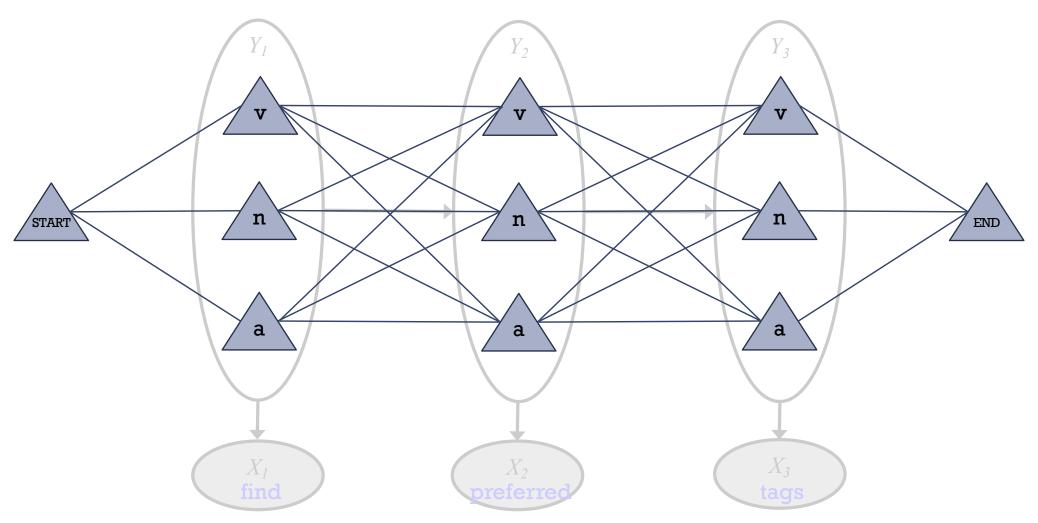




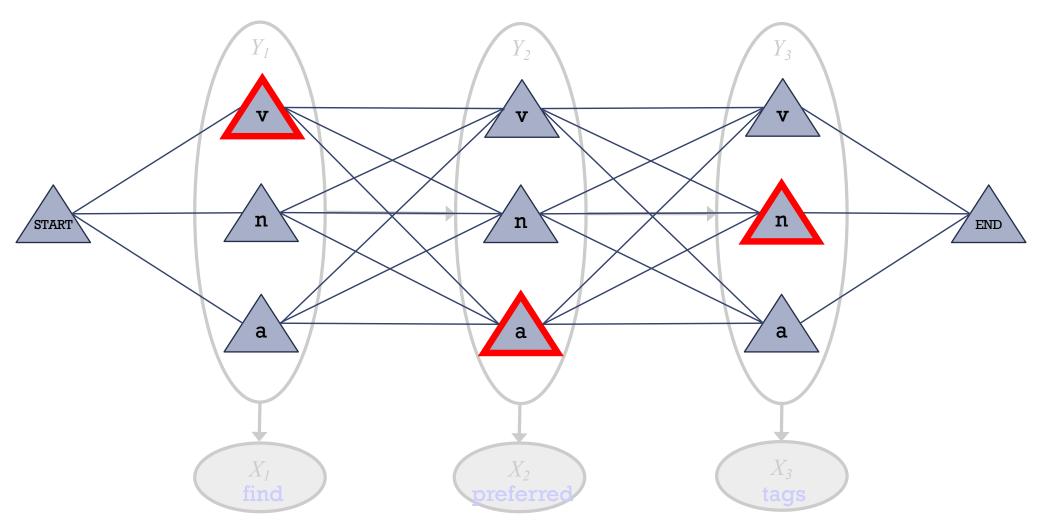




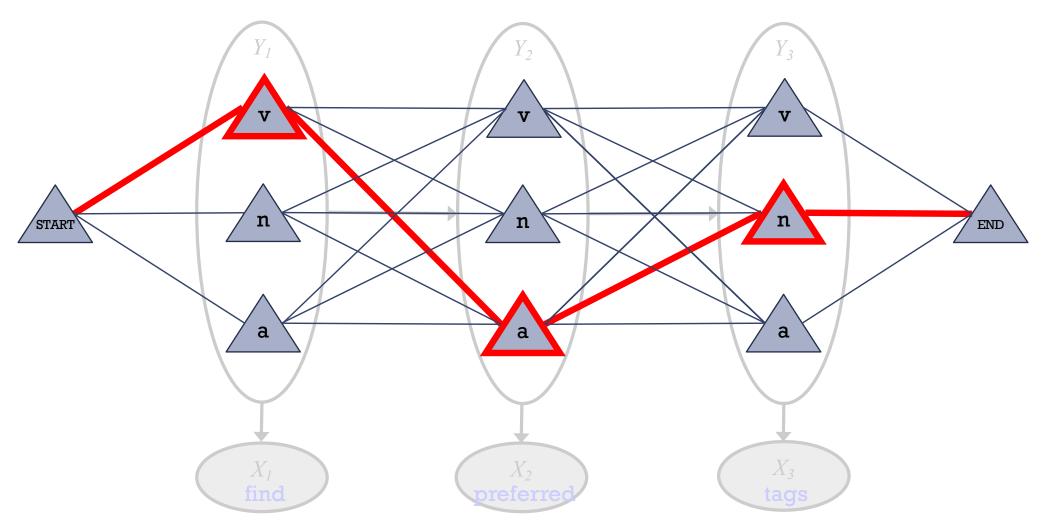
• Let's show the possible values for each variable



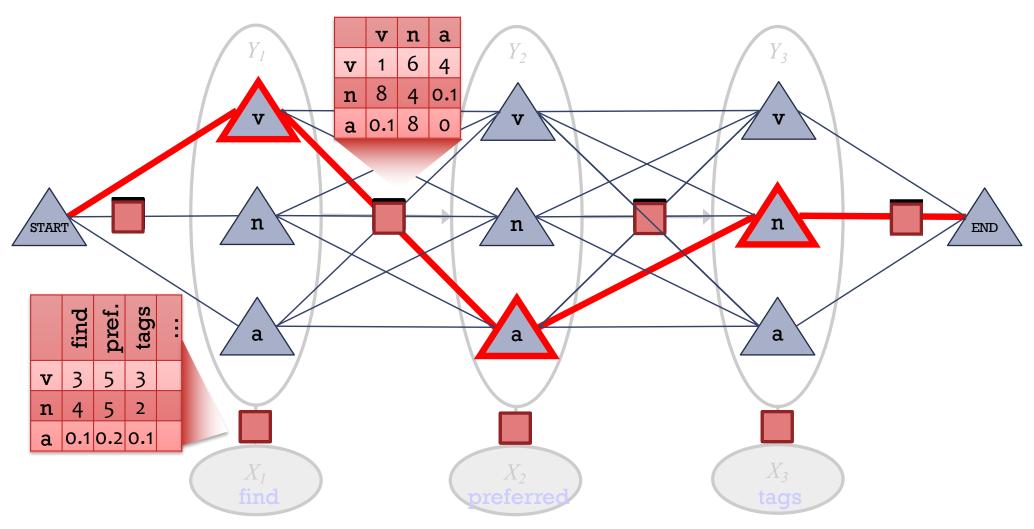
• Let's show the possible values for each variable



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

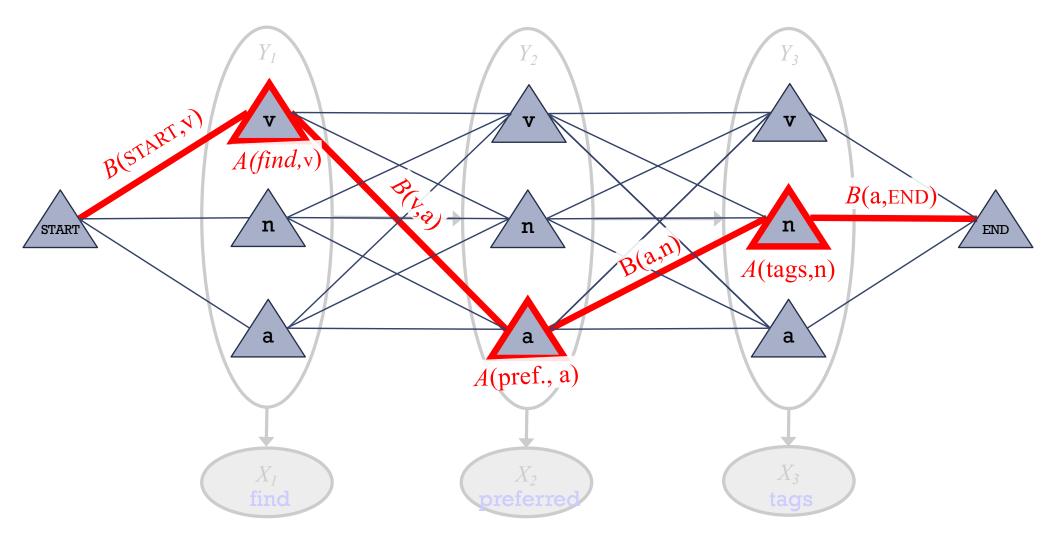


- Let's show the possible values for each variable
- One possible assignment
- And what the 7 transition / emission factors think of it ...



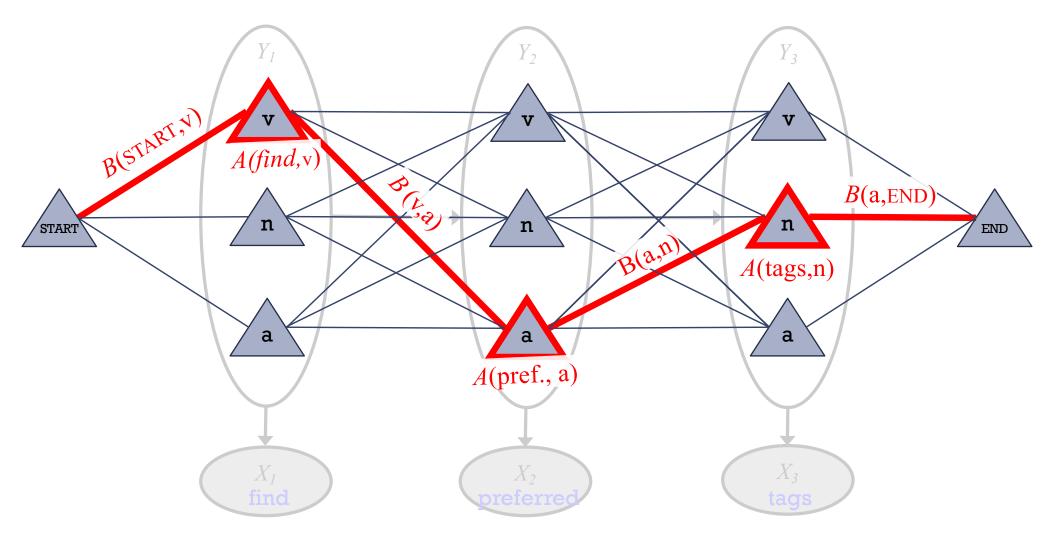
- Let's show the possible *values* for each variable
- One possible assignment
- And what the 7 transition / emission factors think of it ...

### Viterbi Algorithm: Most Probable Assignment

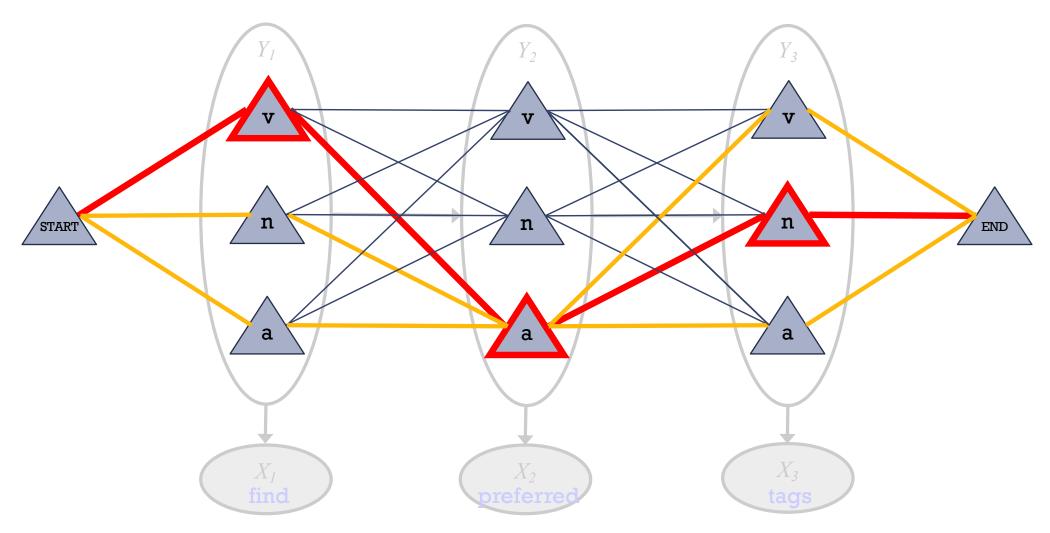


- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * product of 7 numbers$
- Numbers associated with edges and nodes of path
- Most probable assignment = path with highest product

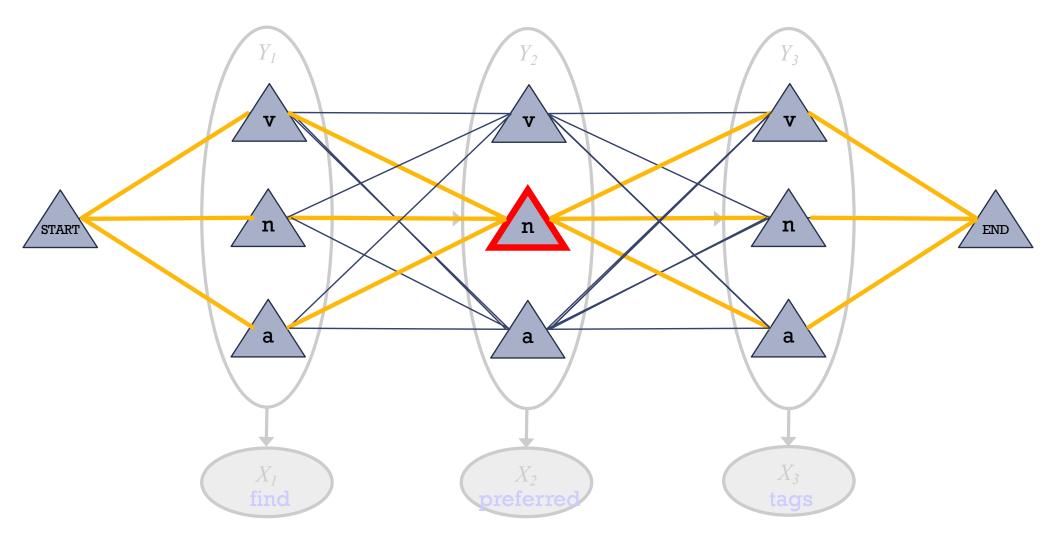
### Viterbi Algorithm: Most Probable Assignment



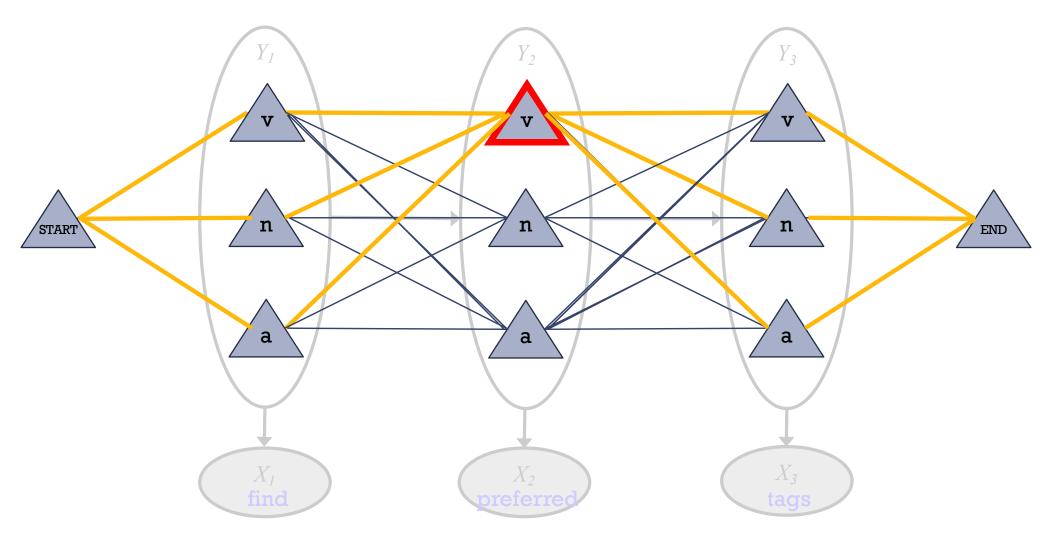
• So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * product weight of one path$ 



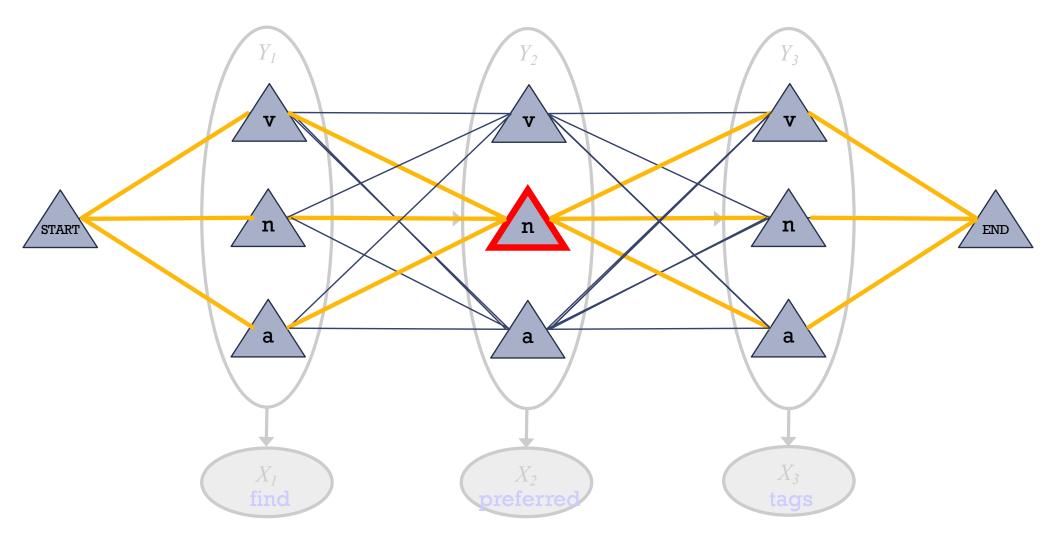
- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/\mathbf{Z}) * \text{product weight of one path}$
- Marginal probability  $p(Y_2 = a) = (1/Z) * total weight of all paths through a$



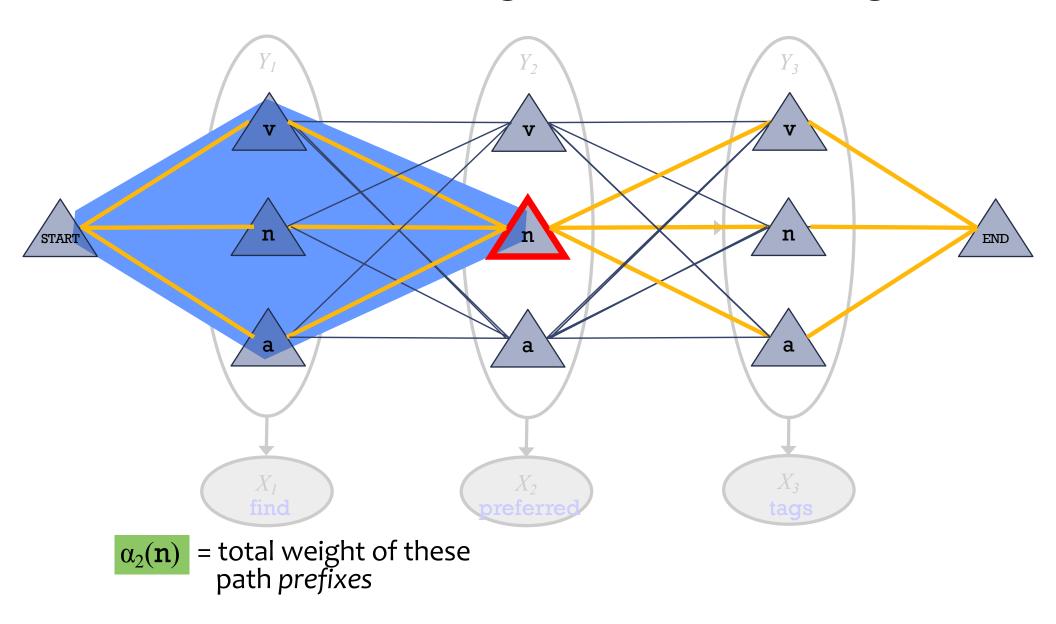
- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/\mathbf{Z}) * \text{product weight of one path}$
- Marginal probability  $p(Y_2 = n)$ = (1/Z) \* total weight of all paths through

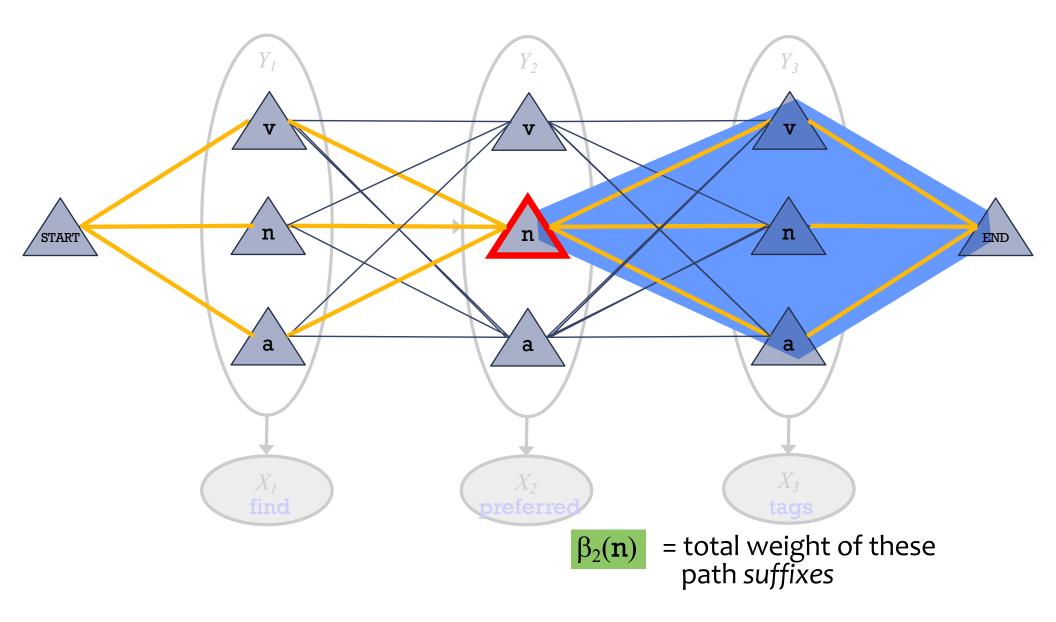


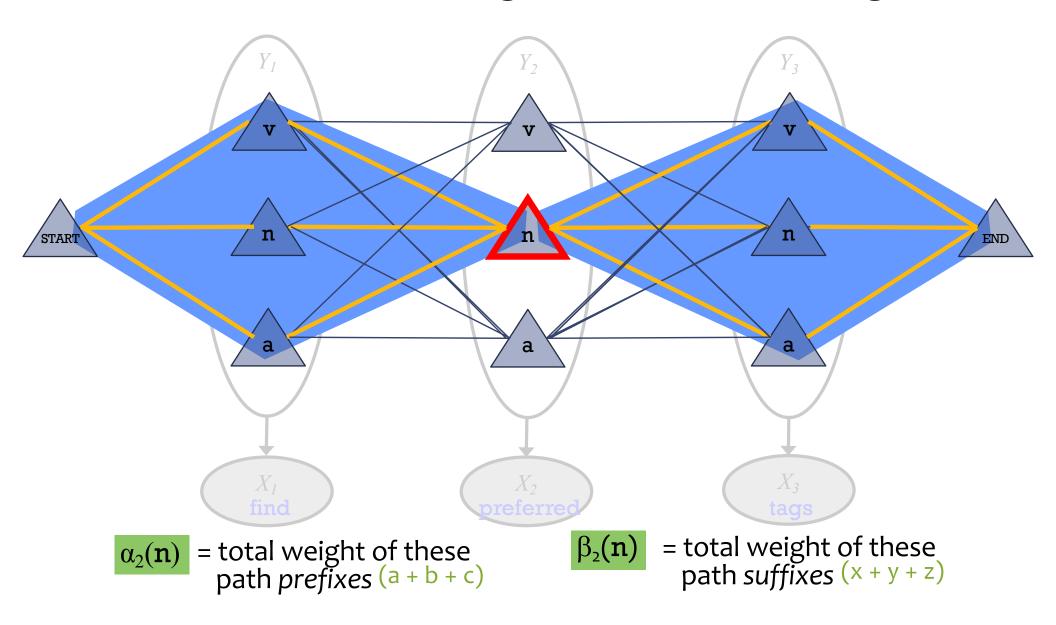
- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/\mathbf{Z}) * \text{product weight of one path}$
- Marginal probability  $p(Y_2 = v)$ = (1/Z) \* total weight of all paths through



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





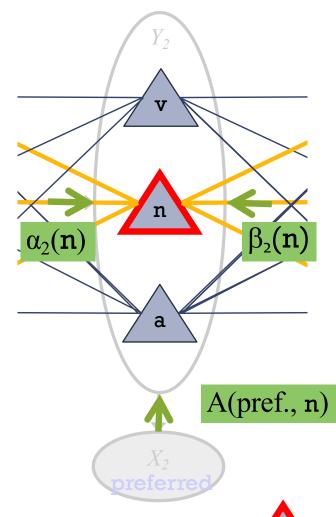


Product gives  $\frac{ax+ay+az+bx+by+bz+cx+cy+cz}{ax+ay+az+bx+by+bz+cx+cy+cz} = total weight of paths$ 

Oops! The weight of a path through a state also includes a weight at that state.

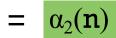
So  $\alpha(\mathbf{n}) \cdot \beta(\mathbf{n})$  isn't enough.

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

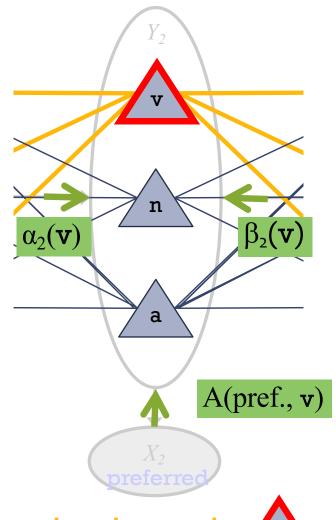


"belief that  $Y_2 = \mathbf{n}$ "

total weight of all paths through



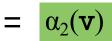
 $\alpha_2(\mathbf{n})$  A(pref.,  $\mathbf{n}$ )  $\beta_2(\mathbf{n})$ 



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

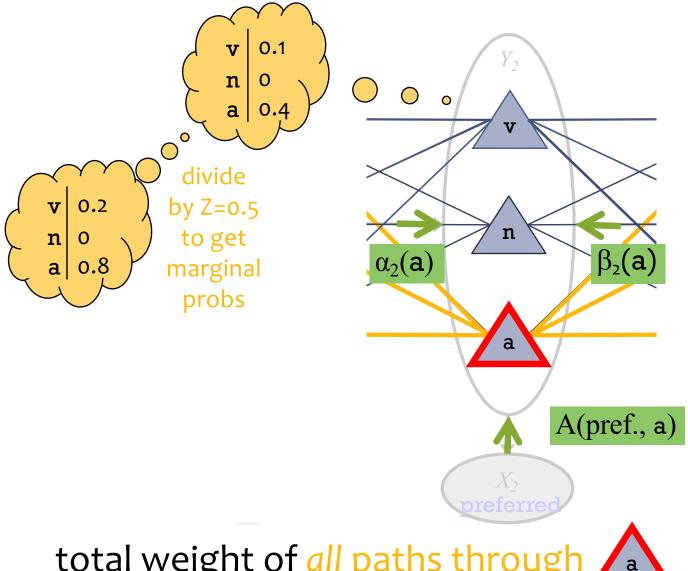
"belief that  $Y_2 = \mathbf{n}$ "

total weight of all paths through



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





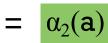
"belief that  $Y_2 = \mathbf{v}$ "

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

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

sum = Z(total weight of all paths)

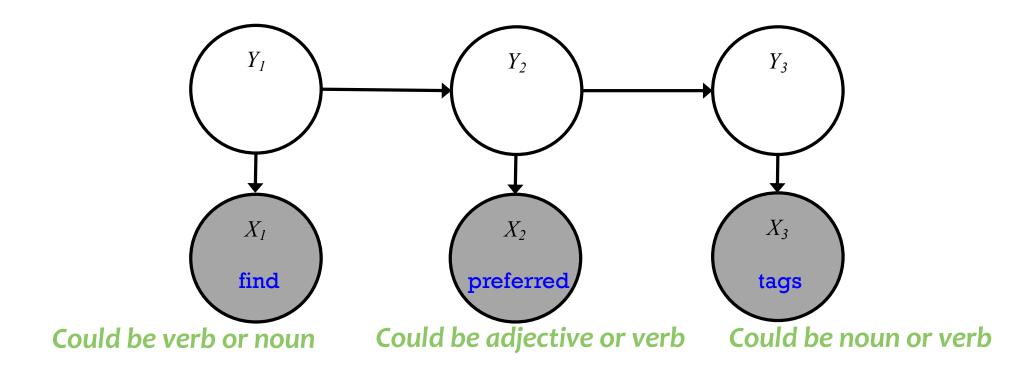
total weight of all paths through



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



## Forward-Backward Algorithm



#### Inference for HMMs

#### Whiteboard

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

# Forward-Backward Algorithm

Define: 
$$\propto_{\ell}(k) \triangleq p(x_1, ..., x_{\ell}, y_{\ell} = k)$$
 $p_{\ell}(k) \triangleq p(x_{\ell+1}, ..., x_{\ell}, y_{\ell} = k)$ 
 $p_{\ell}(k) \triangleq p(x_{\ell+1}, ..., x_{\ell+1}, y_{\ell} = k)$ 
 $p_{\ell}(k) = p(x_{\ell+1}, y_{\ell} = k)$ 
 $p_{\ell}(k) = p(x_{\ell+1}, y_{\ell+1} = k)$ 
 $p_{\ell$ 

## Derivation of Forward Algorithm

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

Derivation:
$$\alpha_{T}(\varepsilon_{NO}) = p(x_1, ..., x_{\tau}, y_{\tau} = \varepsilon_{NO})$$

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

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

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

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

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

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

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

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

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

## Viterbi Algorithm

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

"back points"  $\longrightarrow$   $b_{\xi}(k) \triangleq \alpha_{y_1, \dots, y_{t-1}} p(x_1, \dots, x_t, y_1, \dots, y_{t-1}, y_t=k)$ 

Assume  $y_0 = START$ 

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

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

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

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

(3) Compute Most Probable Assignment

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

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

Follow the

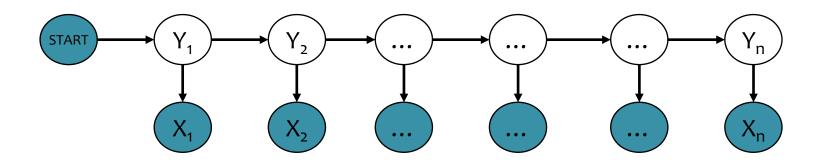
"back points"

#### 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
  - 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

$$egin{aligned} h_{m{ heta}}(m{x}) &= \underset{\hat{m{y}}}{\operatorname{argmin}} & \mathbb{E}_{m{y} \sim p_{m{ heta}}(\cdot | m{x})}[\ell(\hat{m{y}}, m{y})] \ &= \underset{\hat{m{y}}}{\operatorname{argmin}} & \sum_{m{y}} p_{m{ heta}}(m{y} \mid m{x})\ell(\hat{m{y}}, m{y}) \end{aligned}$$

## 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:

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

which is exactly the Viterbi decoding problem!

## 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 **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.

## 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