# Hidden Markov Models (Part II)

Matt Gormley & Henry Chai
Lecture 19
Oct. 29, 2021

# Reminders

- **Homework 6: Learning Theory / Generative Models**
  - **Out: Thu, Oct. 21**
  - **Due: Thu, Oct. 28 at 11:59pm**
  - **Same collaboration policy as Homework 3**
    - **Opt-in to homework groups on Piazza**
  - **IMPORTANT: you may only use 2 grace days on Homework 6**
    - **Last posible moment to submit HW6: Sat, Oct. 30 at 11:59pm**
- **Midterm Exam 2**
  - **Tue, Nov. 2, 6:30pm – 8:30pm**
- **Practice for Exam 2**
  - **Practice problems released on course website**
    - **(Tentatively) Out: Thu, Oct. 21**
  - **Mock Exam 2**
    - **(Tentatively) Out: Thu, Oct. 28**
    - **Due Sun, Oct. 31 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
   $$\ell(\theta) = \log p(x^{(1)}|\theta) + \ldots + \log p(x^{(N)}|\theta)$$
3. Compute partial derivatives (i.e. gradient)
   $$\partial\ell(\theta)/\partial\theta_1 = \ldots$$
   $$\partial\ell(\theta)/\partial\theta_2 = \ldots$$
   $$\ldots$$
   $$\partial\ell(\theta)/\partial\theta_M = \ldots$$
4. Set derivatives to zero and solve for $\theta$
   $$\partial\ell(\theta)/\partial\theta_m = 0 \text{ for all } m \in \{1, \ldots, M\}$$
   $\theta^{MLE}$ = solution to system of $M$ equations and $M$ variables
5. Compute the second derivative and check that $\ell(\theta)$ is concave down
   at $\theta^{MLE}$

# 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 $x^{(i)} \in \{1, \ldots, M\}$ and $x^{(i)} \sim \text{Categorical}(\boldsymbol{\phi})$.

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

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

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

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

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

## HMM Parameters:

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

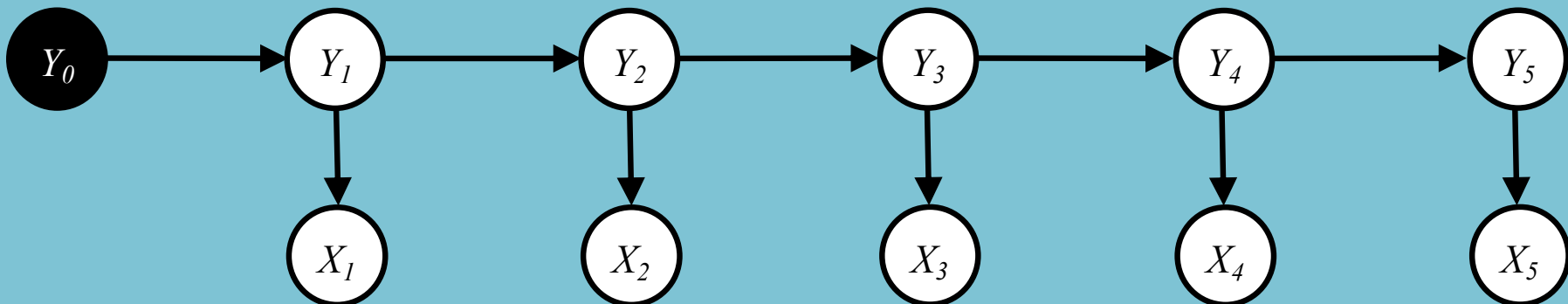Transition matrix, $\mathbf{B}$, where $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

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

# Training HMMs

*Whiteboard*

– (Supervised) Likelihood for an HMM
– Maximum Likelihood Estimation (MLE) for HMM

# Supervised Learning for HMMs

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



Data: $D = \{(\vec{x}^{(i)}, \vec{y}^{(i)})\}_{i=1}^{N}$ 　 $\vec{x} = [x_1, \ldots, x_T]^T$
$\vec{y} = [y_1, \ldots, y_T]^T$

Likelihood:

$$\ell(A,B,C) = \sum_{i=1}^{N} \log p(\vec{x}^{(i)}, \vec{y}^{(i)} \mid A, B, C)$$

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

MLE:

$$\hat{A}, \hat{B}, \hat{C} = \operatorname*{argmax}_{A,B,C} \ell(A,B,C)$$

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

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

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

Can solve in closed form, which yields...

$$\hat{C}_k = \frac{\#(y_1^{(i)} = k)}{N} \quad \forall i, k$$

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

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

9

# Hidden Markov Model

**HMM Parameters:**

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

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

**Assumption:** $y_0 = \text{START}$

**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.



11

# Hidden Markov Model

**Joint Distribution:**

$$y_0 = \text{START}$$

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

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



12

# Supervised Learning for HMMs

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



$$D = \{(\vec{x}^{(i)}, \vec{y}^{(i)})\}_{i=1}^{N}$$

Likelihood:
$$\ell(A,B) = \sum_{i=1}^{N} \log p(\vec{x}^{(i)}, \vec{y}^{(i)})$$

$$= \sum_{i=1}^{N}\left[\sum_{t=1}^{T} \log p(y_t^{(i)} | y_{t-1}^{(i)}, B) + \log p(x_t^{(i)} | y_t^{(i)}, A)\right]$$

MLE:
$$\hat{A}, \hat{B} = \text{argmax} \; \ell(A,B)$$

$$\hat{A} = \text{argmax} \sum_{i=1}^{N}\left[\sum_{t=1}^{T} \log p(x_t^{(i)} | y_t^{(i)}, A)\right]$$

$$\hat{\vec{B}} = \text{argmax} \sum_{i=1}^{N}\left[\sum_{t=1}^{T} \log p(y_t^{(i)} | y_{t-1}^{(i)}, B)\right]$$

← can solve in closed form to set...

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

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

# TO HMMS AND BEYOND...

# Unsupervised Learning for HMMs

- Unlike **discriminative** models p(y|x), **generative** models p(x,y) can maximize the likelihood of the data D = {x$^{(1)}$, x$^{(2)}$, ..., x$^{(N)}$} where we don't observe any y's.
- This **unsupervised learning** setting can be achieved by finding parameters that maximize the **marginal likelihood**
- We optimize using the **Expectation-Maximization** algorithm

Since we don't observe **y**, we define the marginal probability:

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}} p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}) \tag{1}$$

The log-likelihood of the data is thus:

$$\ell(\boldsymbol{\theta}) = \log \prod_{i=1}^{N} p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})$$

$$= \sum_{i=1}^{N} \log \sum_{\mathbf{y} \in \mathcal{Y}} p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \mathbf{y}) \tag{3}$$

Beyond the scope of today's lecture!

# 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

- 1<sup>st</sup>-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

- 1<sup>st</sup>-order HMM (i.e. bigram HMM)



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

3<sup>rd</sup>-order HMM

**Hidden States, y**

**Observa-tions, x**

18

# BACKGROUND: MESSAGE PASSING

# Great Ideas in ML: Message Passing

*Count the soldiers*

# Great Ideas in ML: Message Passing

*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

*Count the soldiers*

# 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)

23

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*

25

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# INFERENCE FOR HMMS

# Inference

**Question:**

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

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

**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*: 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, $\mathbf{A}$, 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

# THE SEARCH SPACE FOR FORWARD-BACKWARD

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

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

# 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.

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

# Example: HMM for POS Tagging



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

# Inference for HMMs

*Whiteboard*

- Brute Force Evaluation
- Forward-backward search space

# THE FORWARD-BACKWARD ALGORITHM

# 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

44

# 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



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

47

# Viterbi Algorithm: Most Probable Assignment



- So p(**v a n**) = (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 ▲a

49

# Forward-Backward Algorithm: Finds Marginals



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

50

# 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 ▲v

# Forward-Backward Algorithm: Finds Marginals



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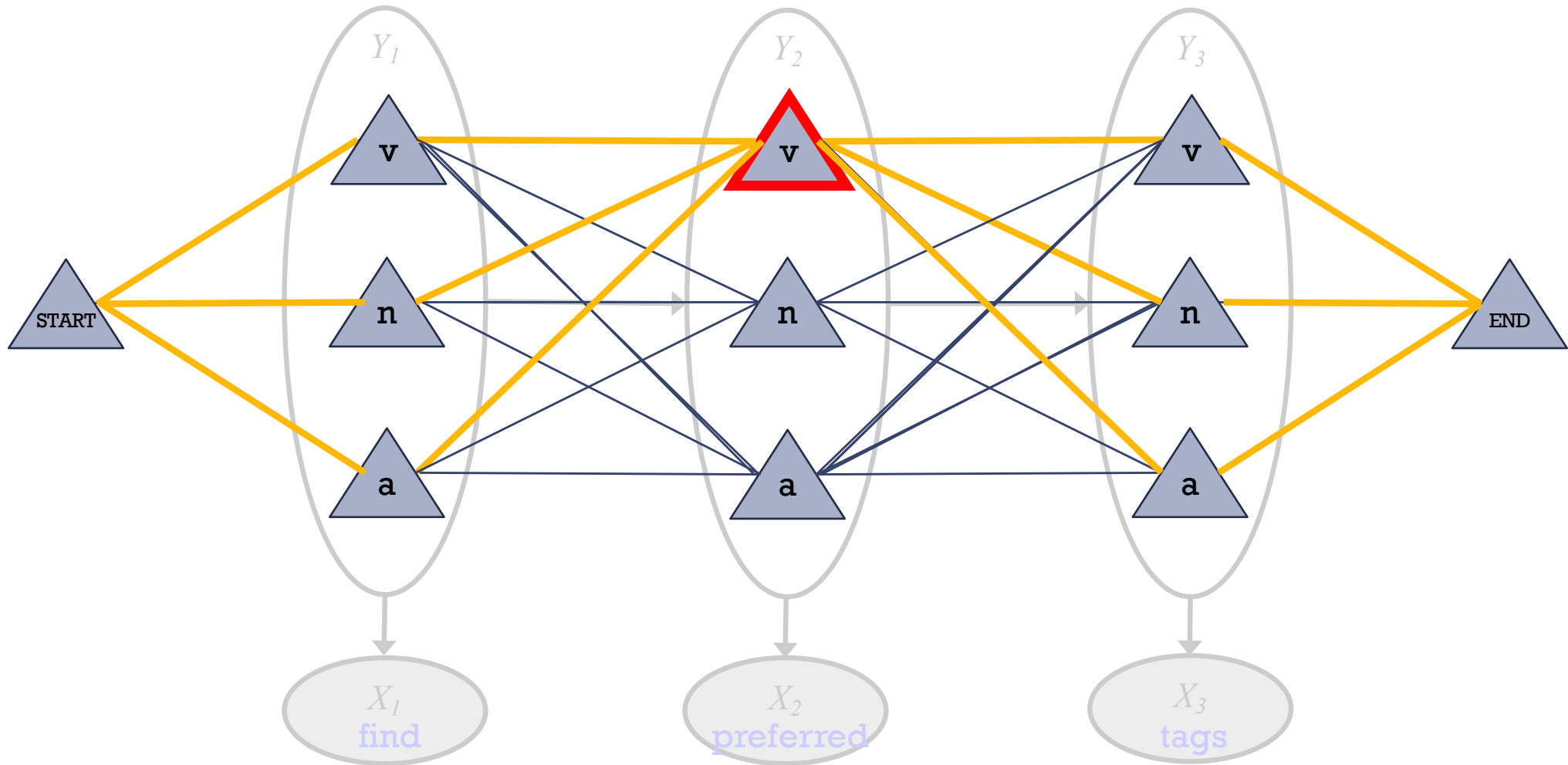     = (1/Z) * total weight of *all* paths through ⧊ **n**

52

# Forward-Backward Algorithm: Finds Marginals



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

(found by dynamic programming: matrix-vector products)

# Forward-Backward Algorithm: Finds Marginals



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

(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)

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

# Forward-Backward Algorithm: Finds Marginals

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.

$Y_2$

v

n

$\alpha_2(\mathbf{n})$     $\beta_2(\mathbf{n})$

a

A(pref., n)

$X_2$
preferred

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

total weight of *all* paths through  n

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

# 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 $\mathbf{v}$

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

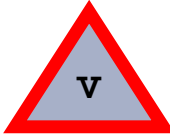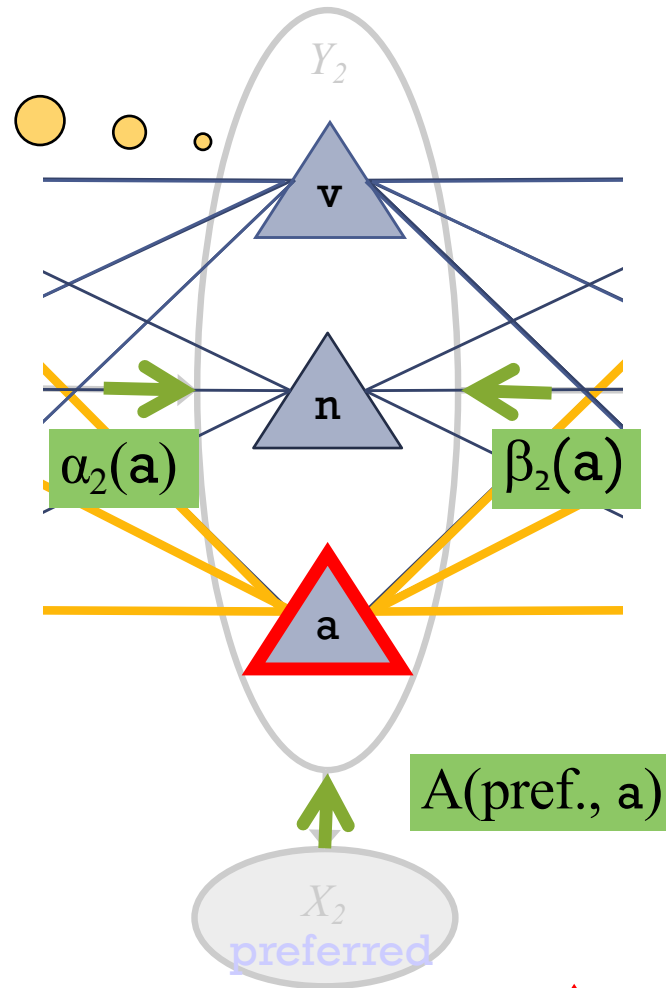# Forward-Backward Algorithm: Finds Marginals

$$\begin{array}{c|c} \mathbf{v} & 0.1 \\ \mathbf{n} & 0 \\ \mathbf{a} & 0.4 \end{array}$$

$$\begin{array}{c|c} \mathbf{v} & 0.2 \\ \mathbf{n} & 0 \\ \mathbf{a} & 0.8 \end{array}$$

divide
by Z=0.5
to get
marginal
probs

$Y_2$

v

n

$\alpha_2(\mathbf{a})$    $\beta_2(\mathbf{a})$

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)

total weight of *all* paths through ▲ a

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

58

# Forward-Backward Algorithm

# Forward-Backward Algorithm
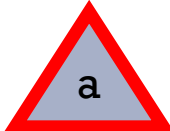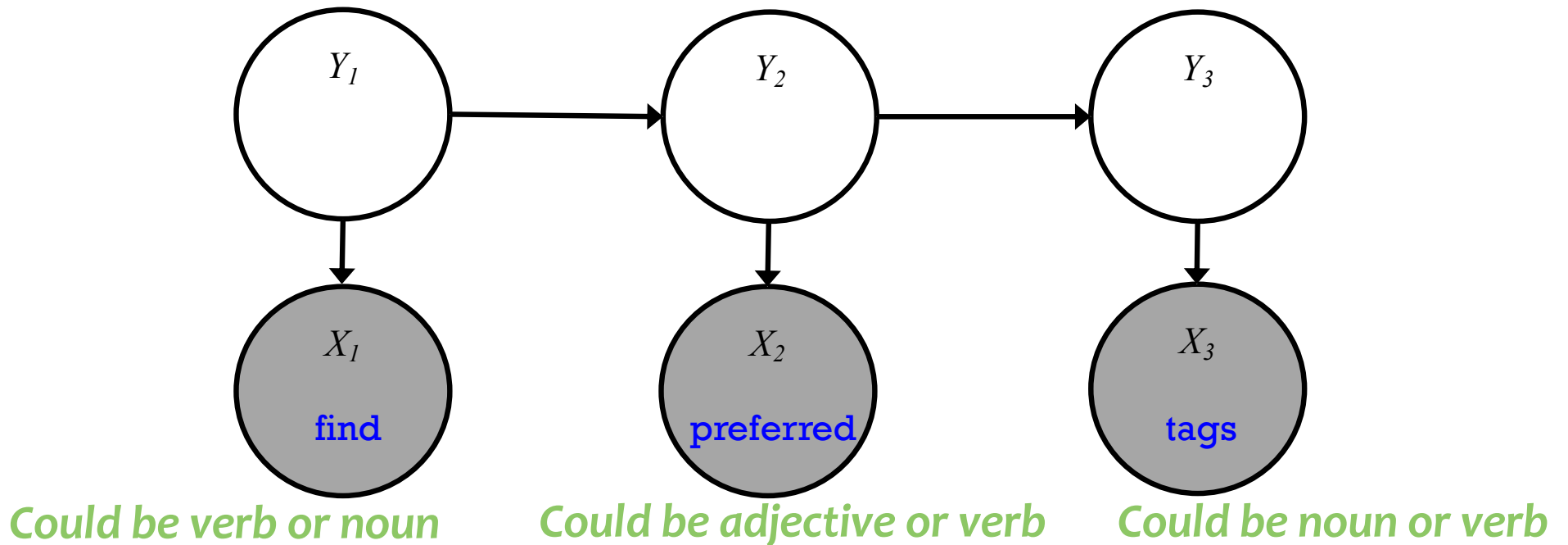
Define: $\alpha_t(k) \triangleq p(x_1, ..., x_t, y_t = k)$

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

Assume $\quad y_0 = \text{START}$

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

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

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

② For $t = 1, ..., T$ :

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

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

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

③ For $t = T, ..., T$ :

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

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

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

⑤ Compute $p(y_t = k | \vec{x}) = \dfrac{\alpha_t(k) \, \beta_t(k)}{P(\vec{x})}$ $\qquad$ [Marginals]