



10-301/601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

HMMs + Bayesian Networks

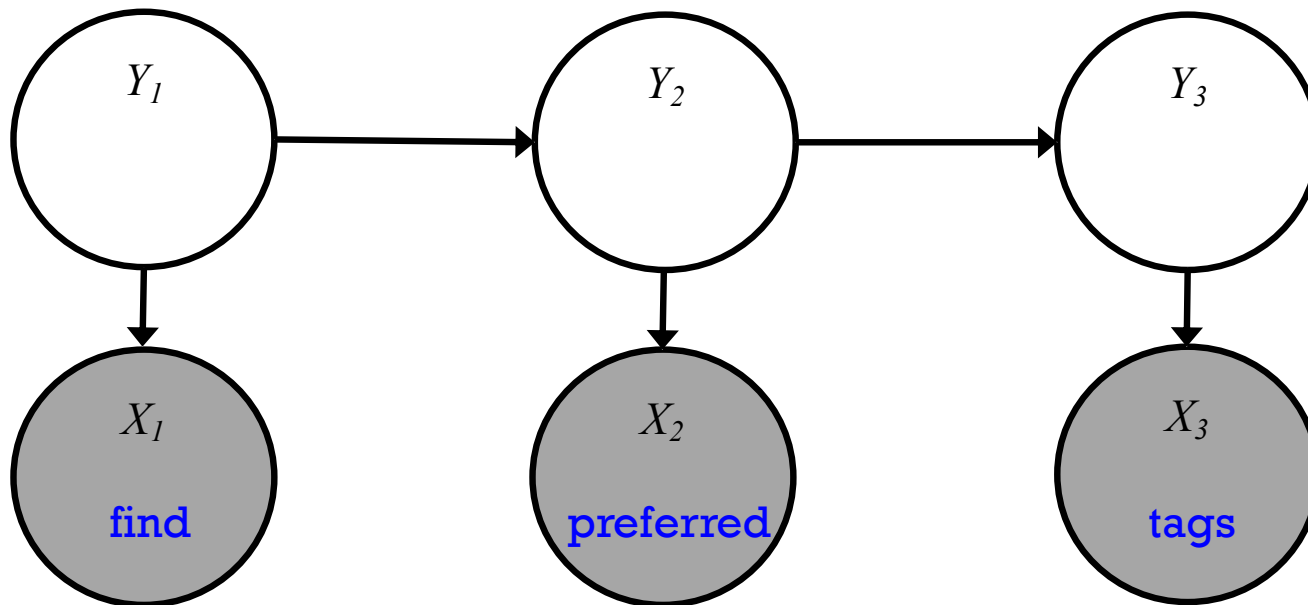
Matt Gormley
Lecture 20
Mar. 30, 2022

Reminders

- **Exam 2 (Thu, Mar 3rd)**
 - Thu, Mar. 31, 6:30pm – 8:30pm
- **Practice for Exam 2**
 - Practice problems released on course website
 - Out: Fri, Mar. 25
 - Mock Exam 2
 - Out: Fri, Mar. 25
 - Due Wed, Mar. 30 at 11:59pm
- **Homework 7: HMMs**
 - Out: Fri, Apr. 1
 - Due: Tue, Apr. 12 at 11:59pm

**EXAMPLE: FORWARD-BACKWARD
ON THREE WORDS**

Forward-Backward Algorithm

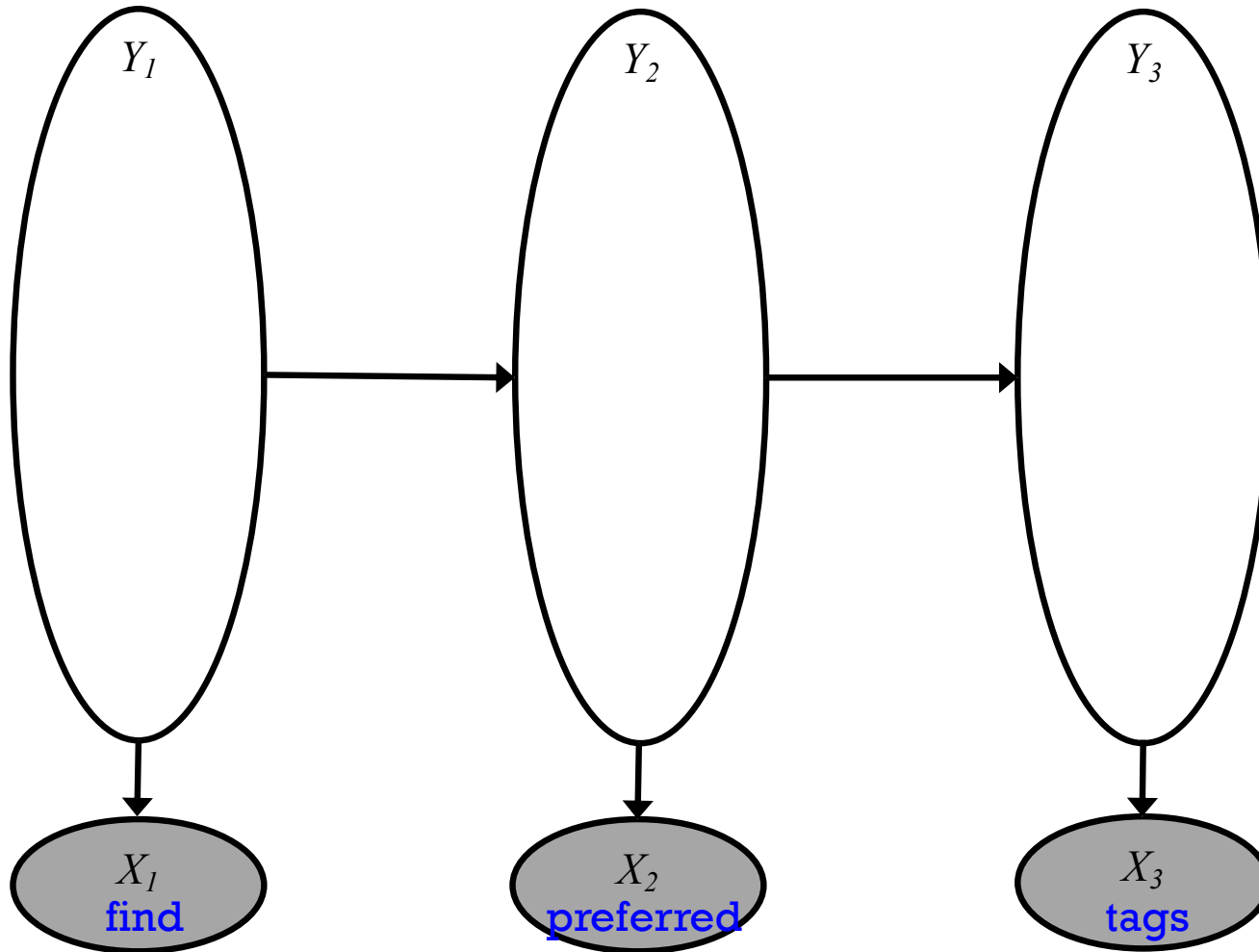


Could be verb or noun

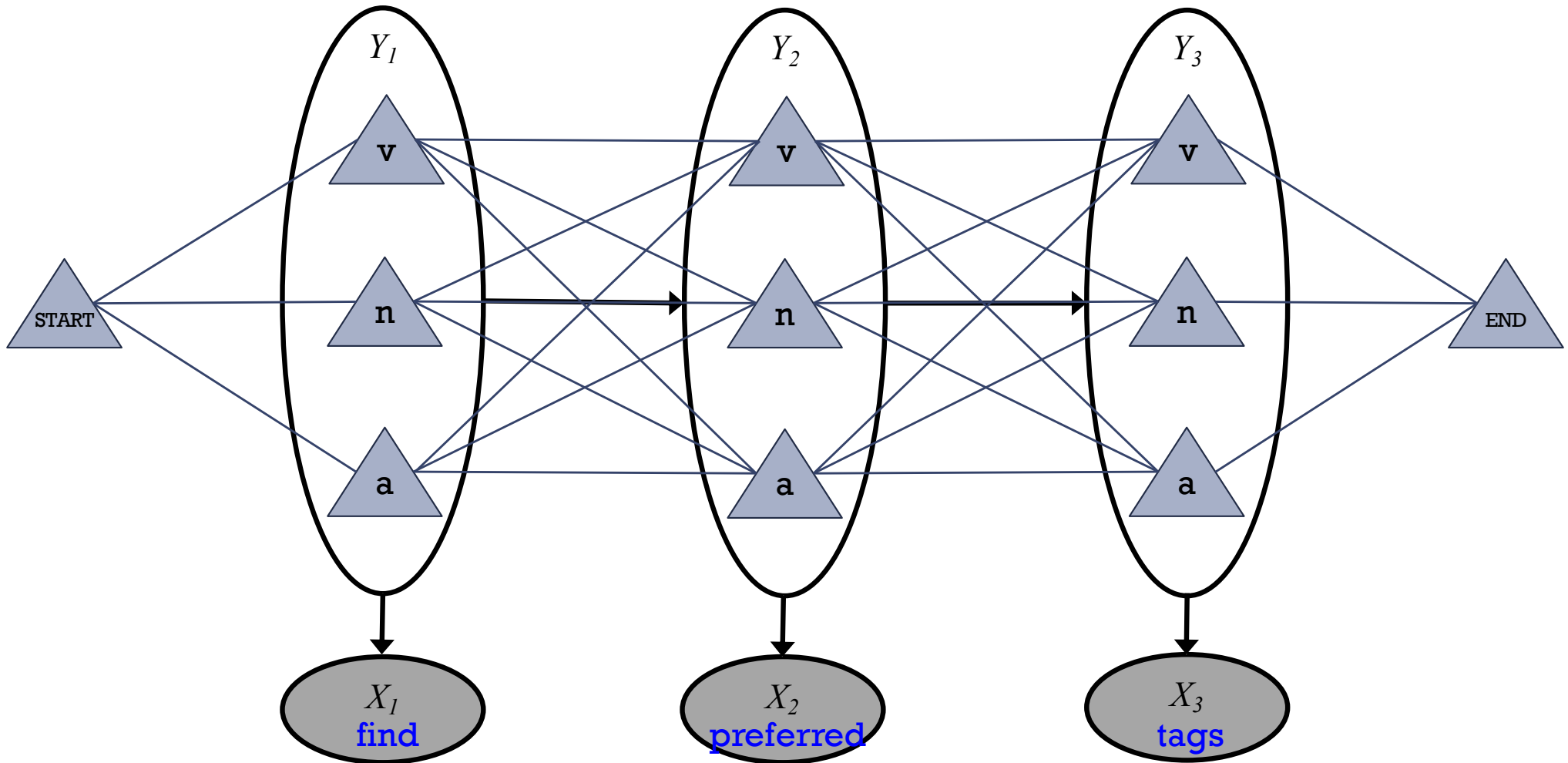
Could be adjective or verb

Could be noun or verb

Forward-Backward Algorithm

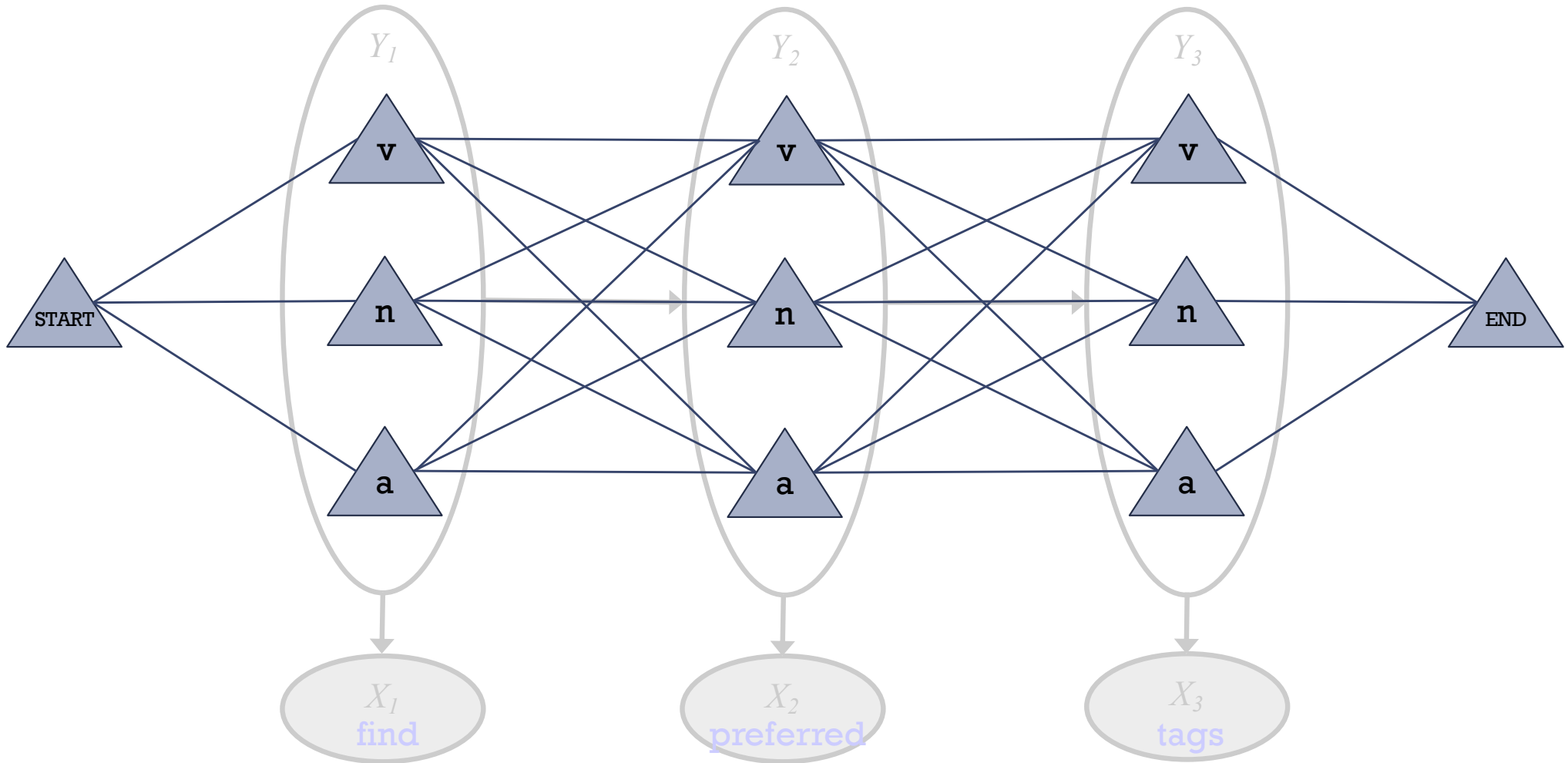


Forward-Backward Algorithm



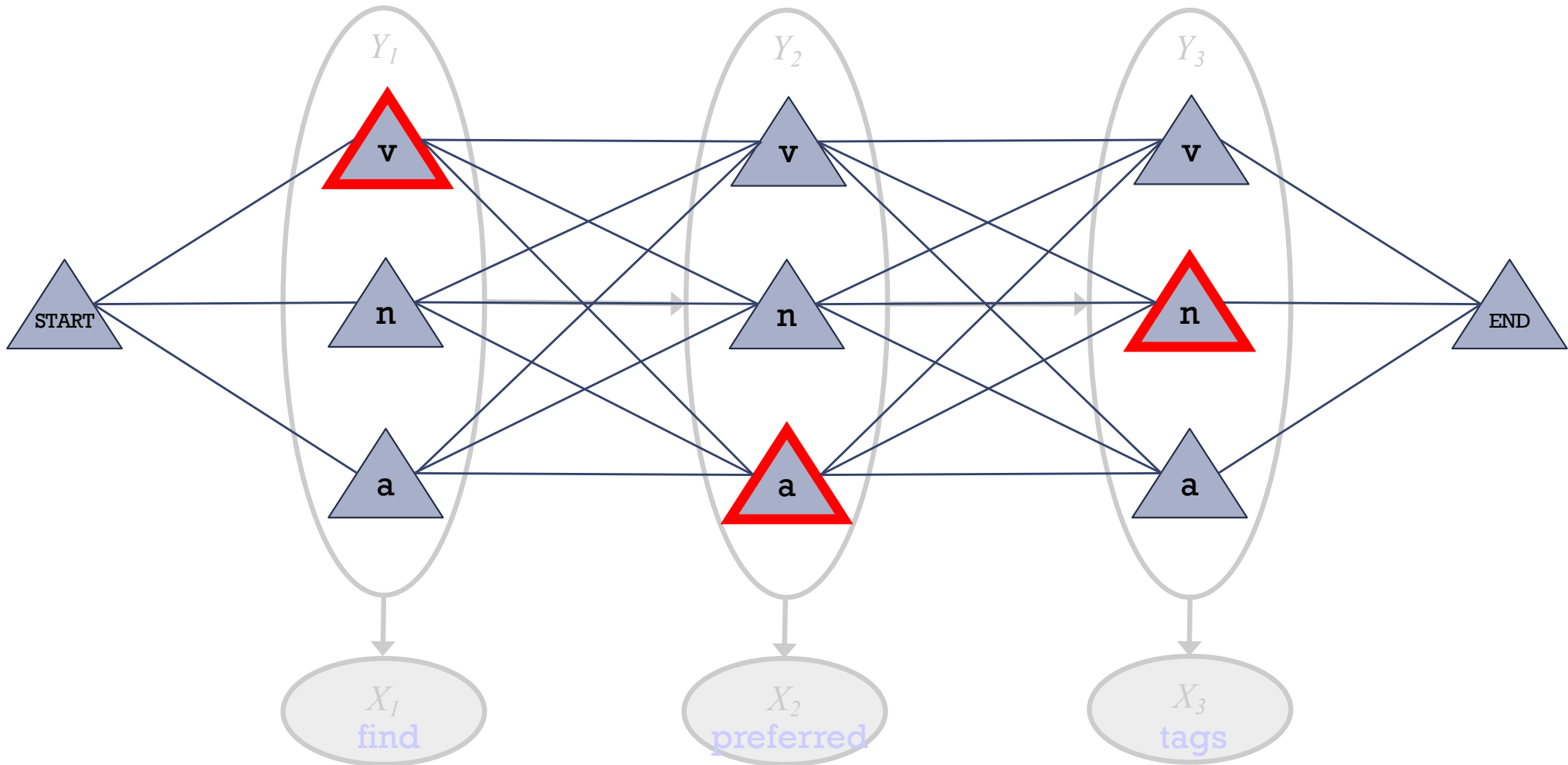
- Let's show the possible *values* for each variable

Forward-Backward Algorithm



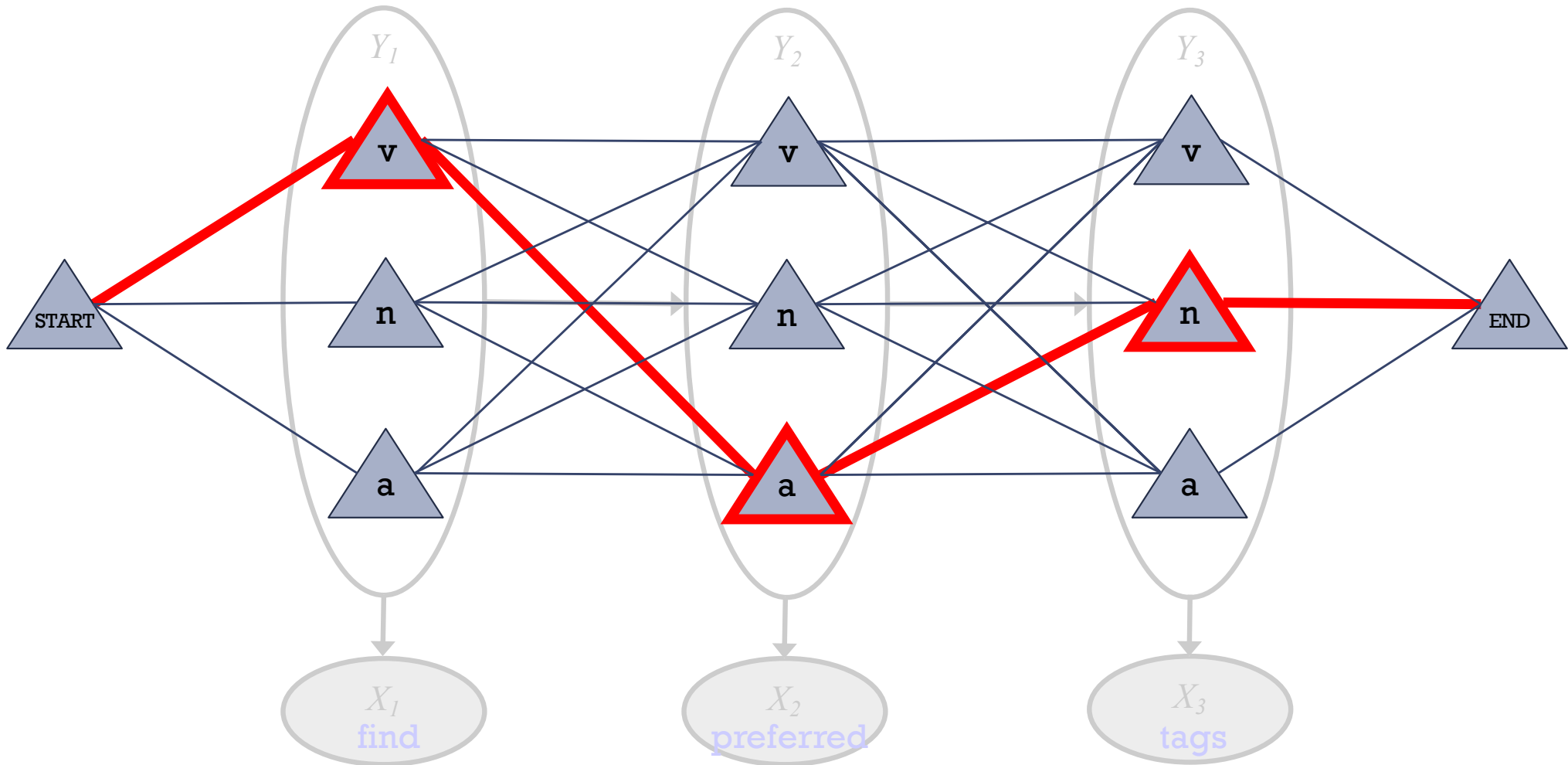
- Let's show the possible *values* for each variable

Forward-Backward Algorithm



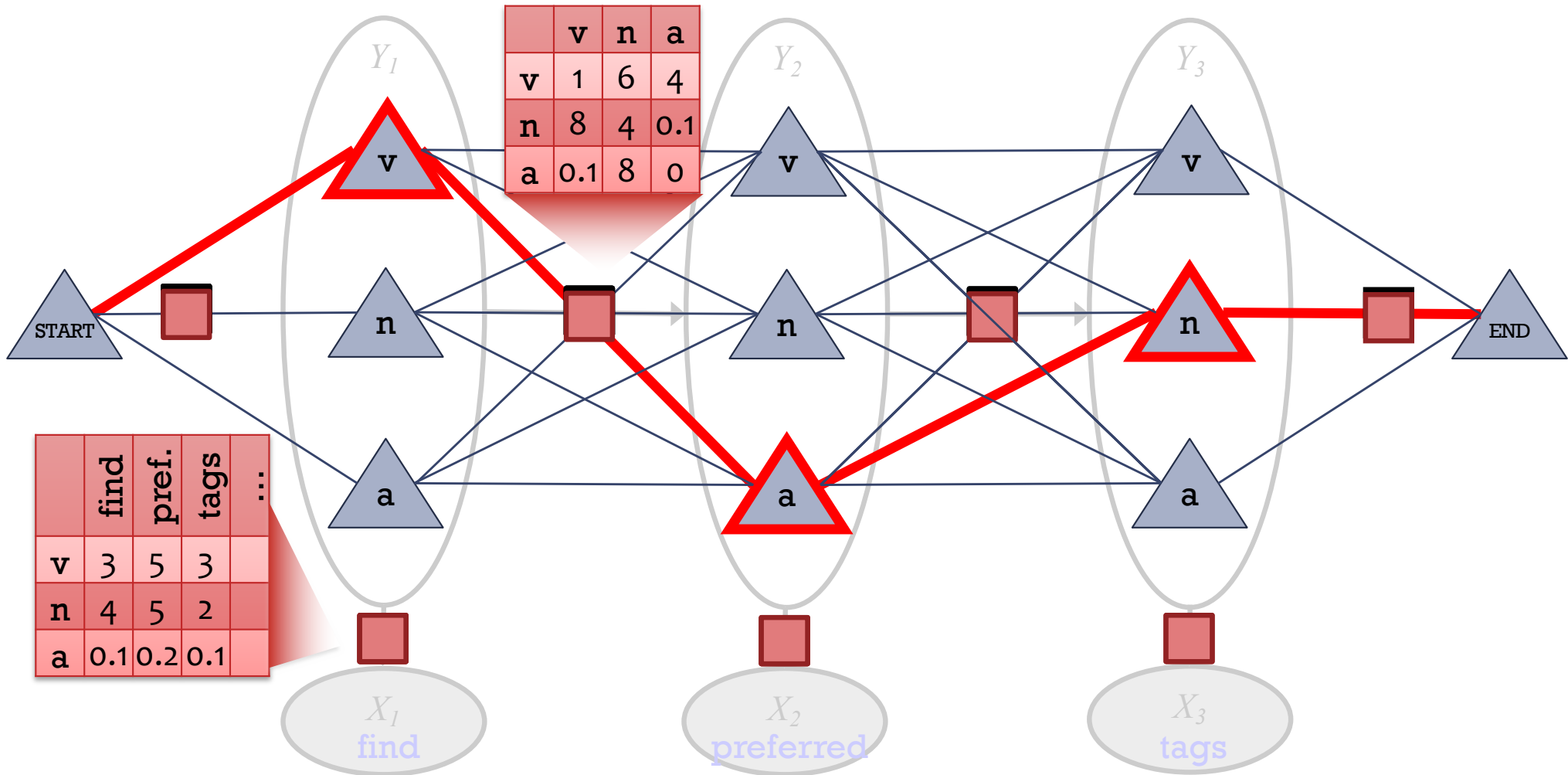
- Let's show the possible *values* for each variable
- One possible assignment

Forward-Backward Algorithm



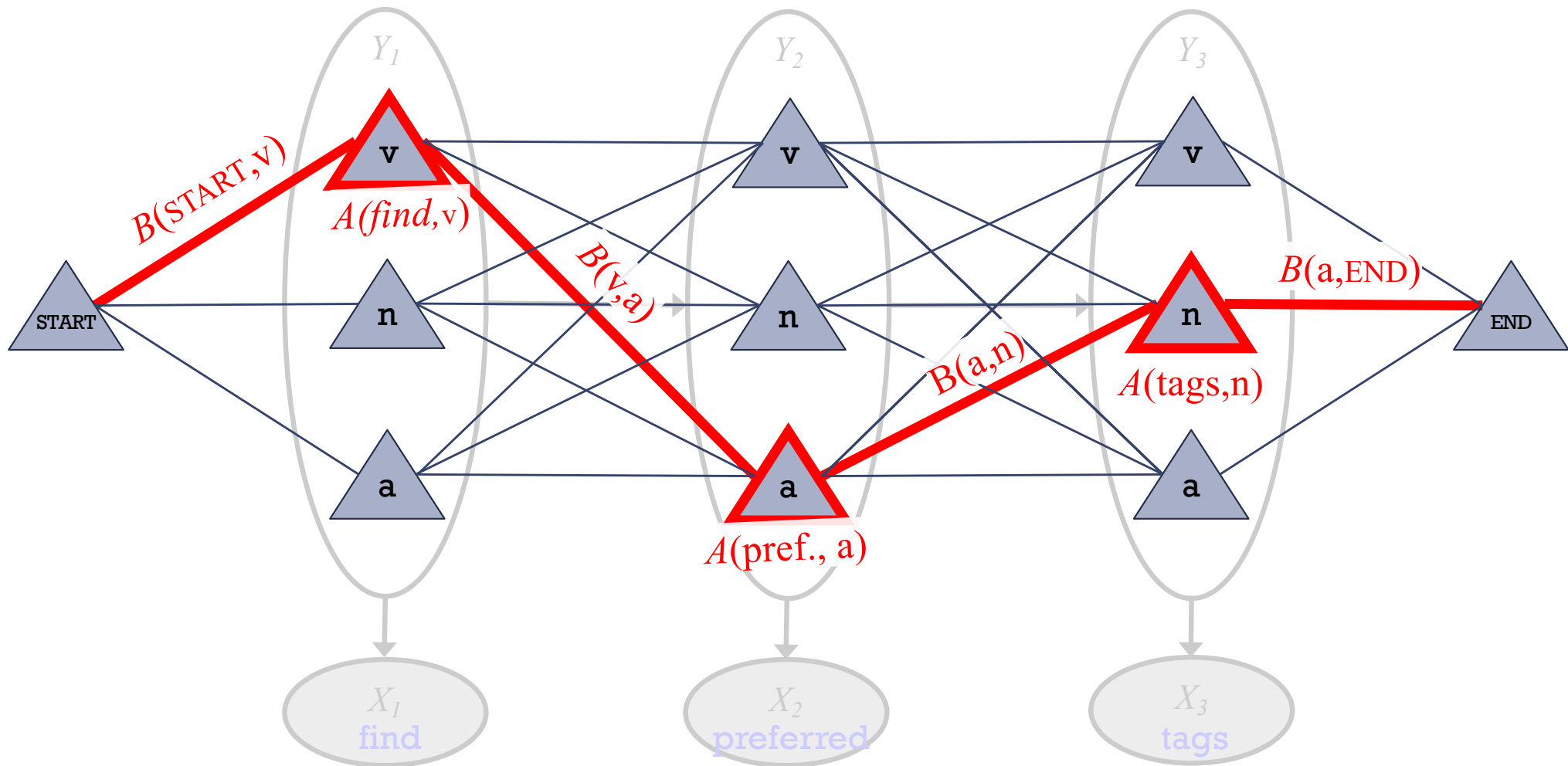
- Let's show the possible *values* for each variable
- One possible assignment
- And what the 7 transition / emission factors **think of it** ...

Forward-Backward Algorithm



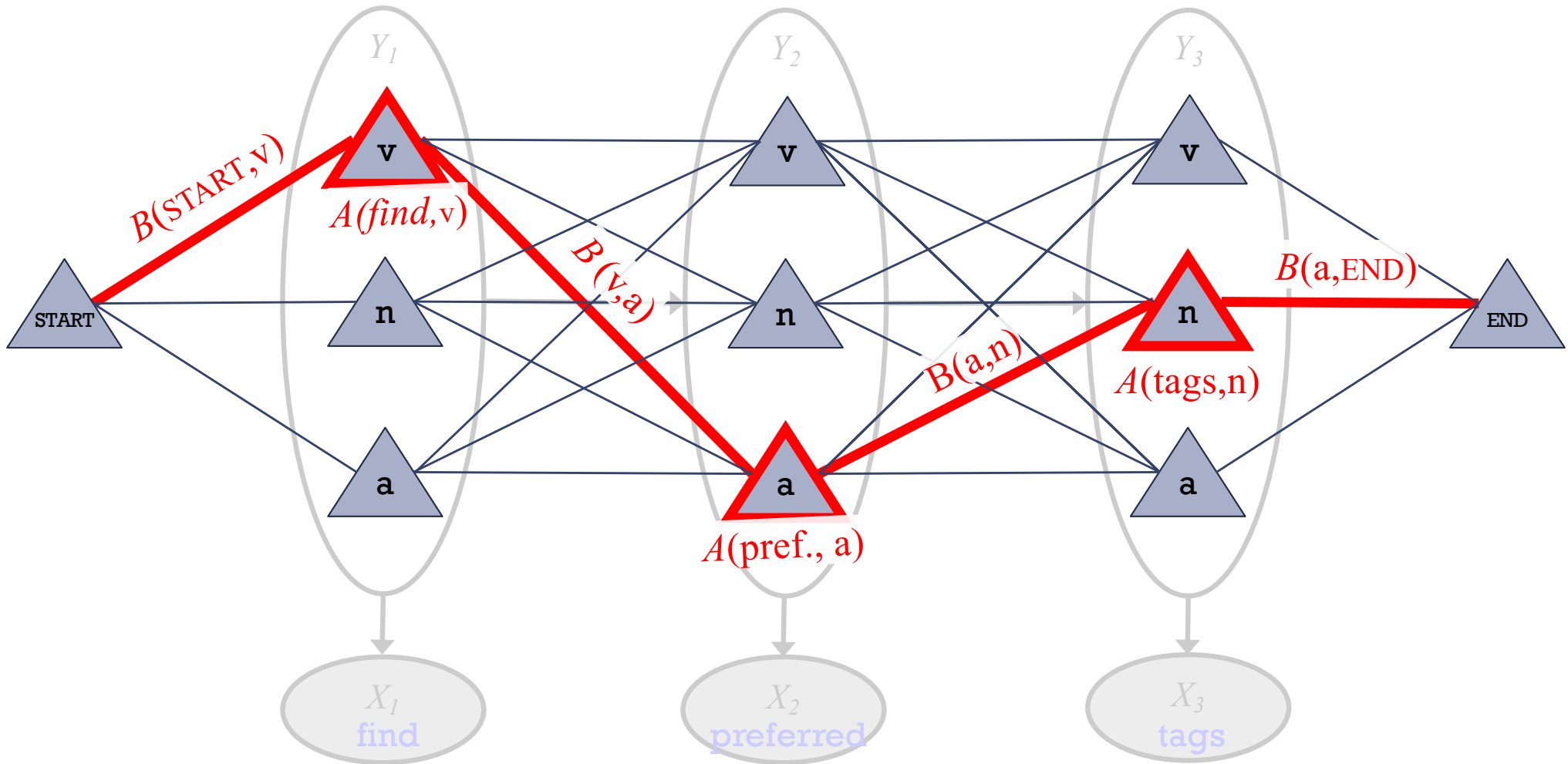
- Let's show the possible *values* for each variable
- One possible assignment
- And what the 7 transition / emission factors **think of it** ...

Viterbi Algorithm: Most Probable Assignment



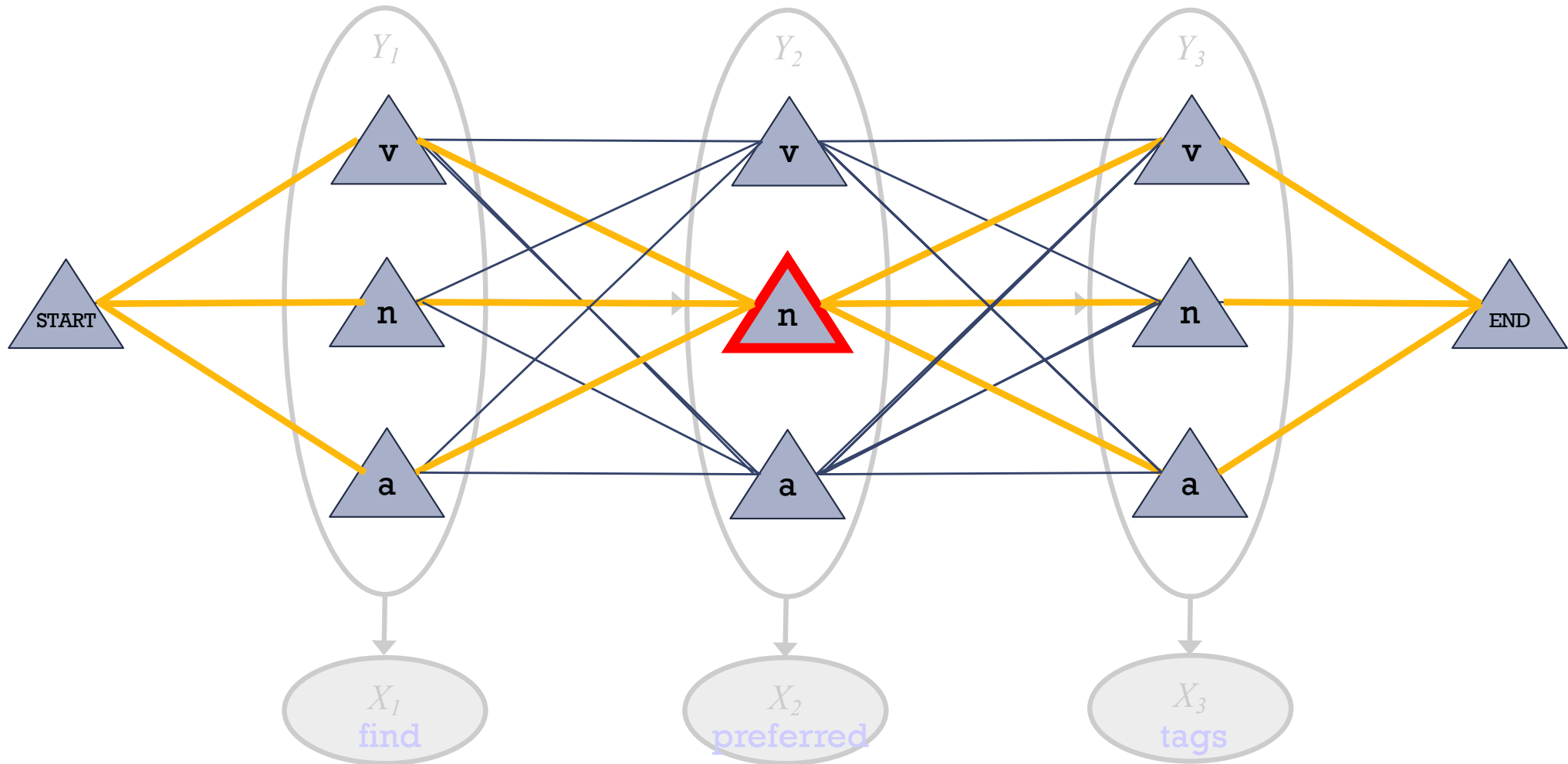
- So $p(v \ a \ n) = (1/Z) * \text{product of 7 numbers}$
- Numbers associated with edges and nodes of path
- Most probable assignment = **path with highest product**

Viterbi Algorithm: Most Probable Assignment



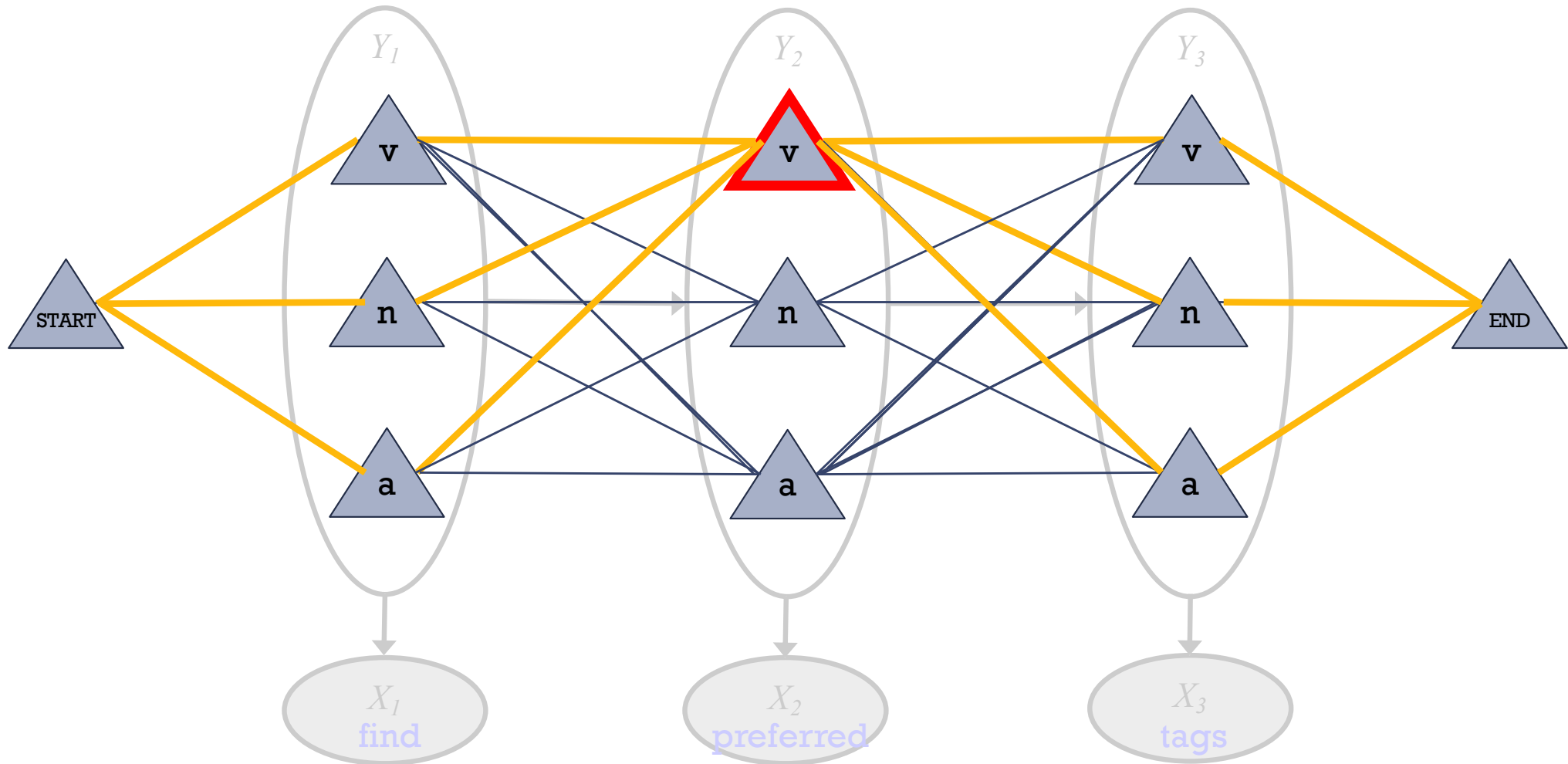
- So $p(v a n) = (1/Z) * \text{product weight of one path}$

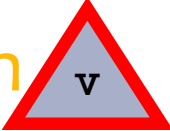
Forward-Backward Algorithm: Finds Marginals



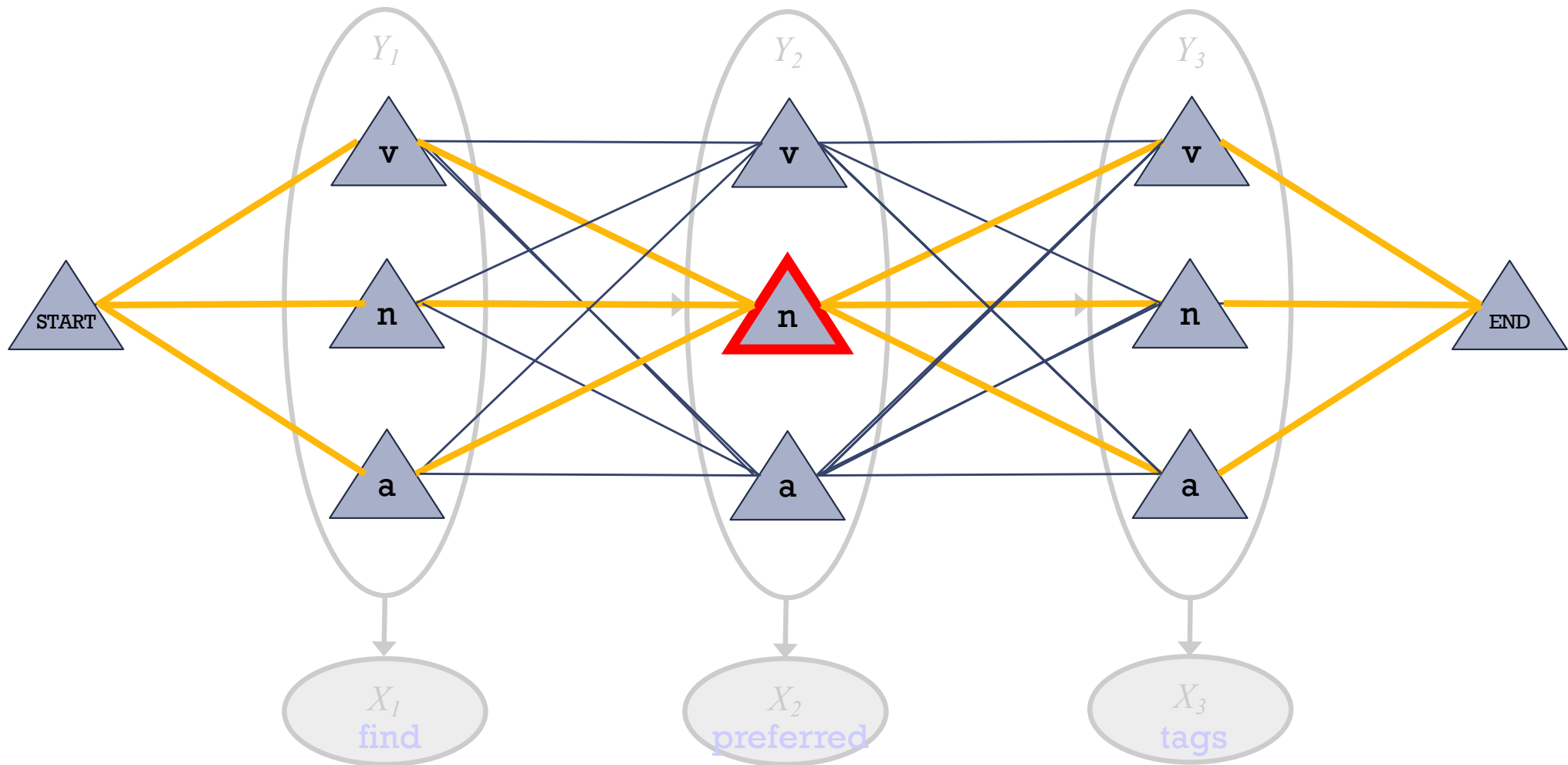
- So $p(\mathbf{v} \ \mathbf{a} \ \mathbf{n}) = (1/Z) * \text{product weight of one path}$
- Marginal probability $p(Y_2 = \mathbf{n}) = (1/Z) * \text{total weight of all paths through } \mathbf{n}$

Forward-Backward Algorithm: Finds Marginals



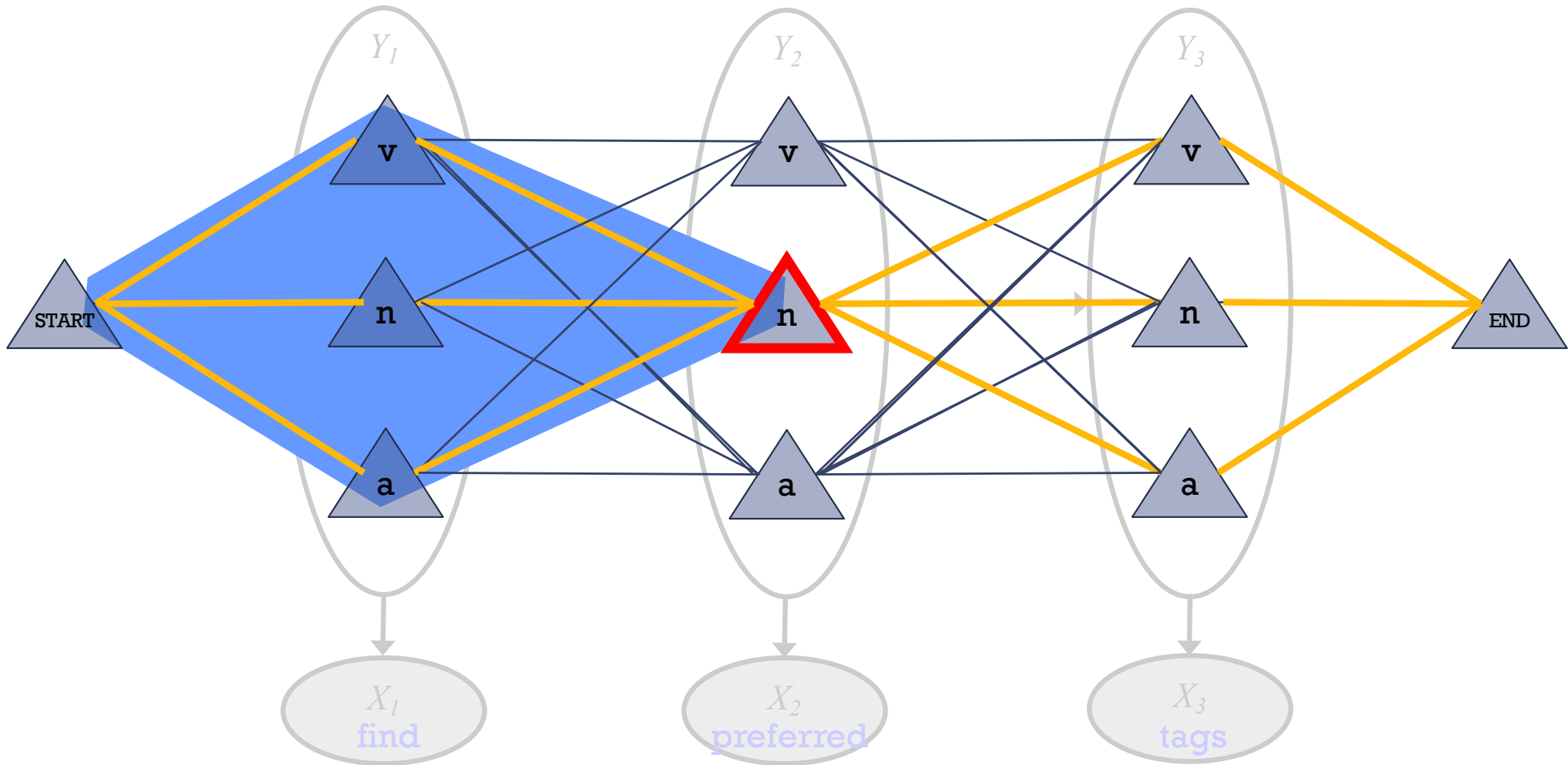
- So $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * \text{product weight of one path}$
- Marginal probability $p(Y_2 = \mathbf{v}) = (1/Z) * \text{total weight of all paths through}$ 

Forward-Backward Algorithm: Finds Marginals



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

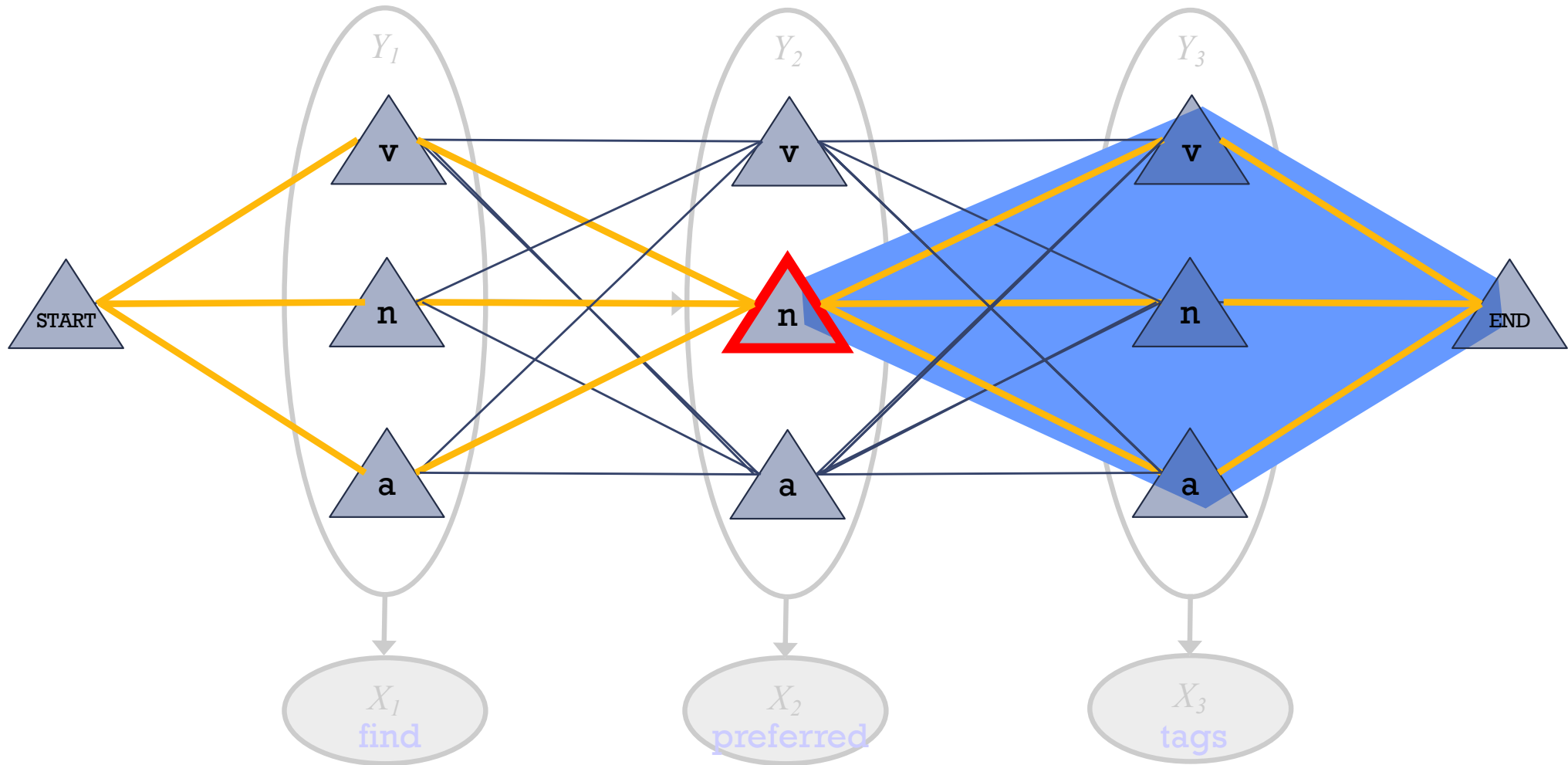
Forward-Backward Algorithm: Finds Marginals



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

(found by dynamic programming: matrix-vector products)

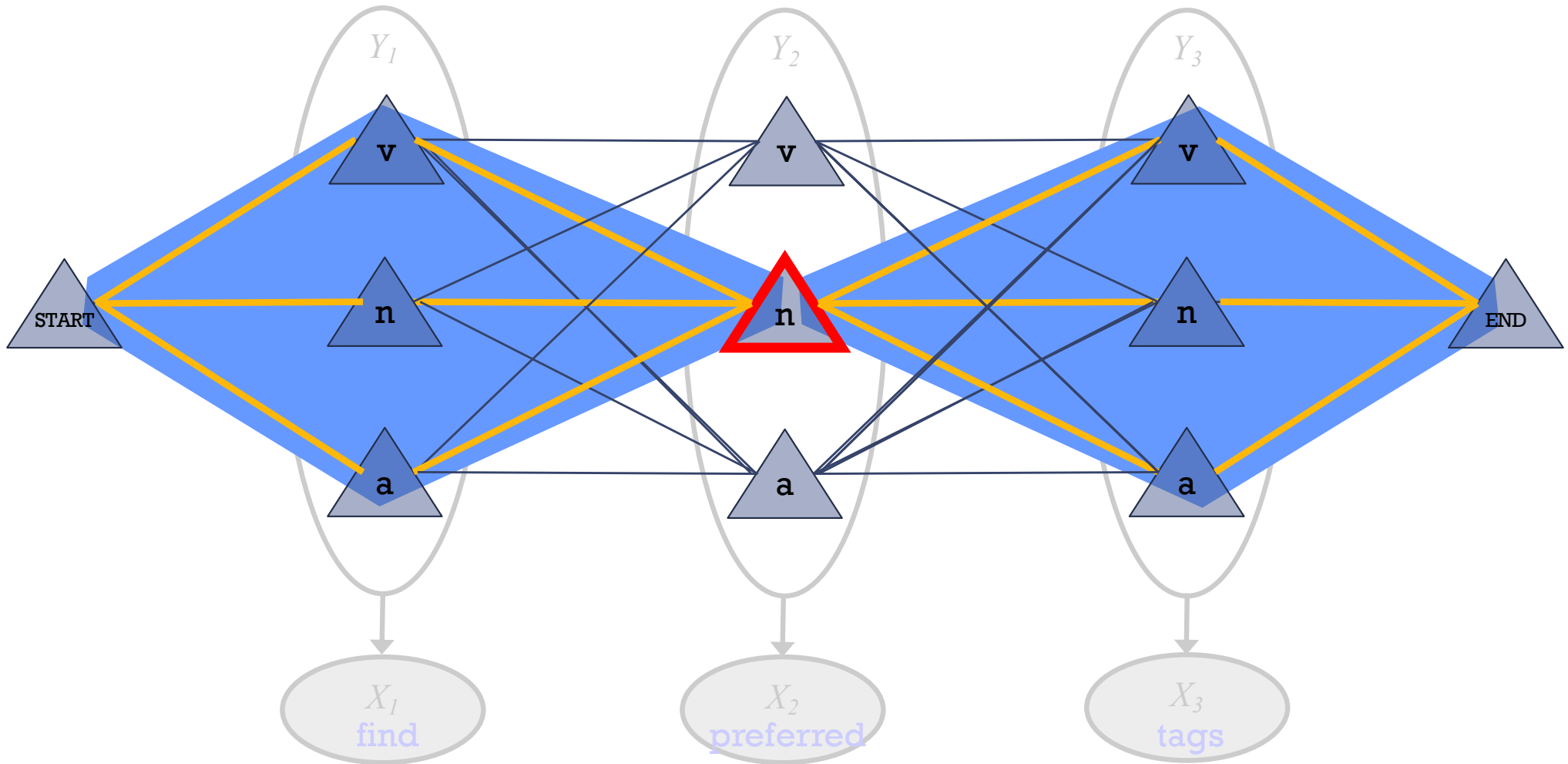
Forward-Backward Algorithm: Finds Marginals



$\beta_2(\mathbf{n})$ = 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 $(x + y + z)$

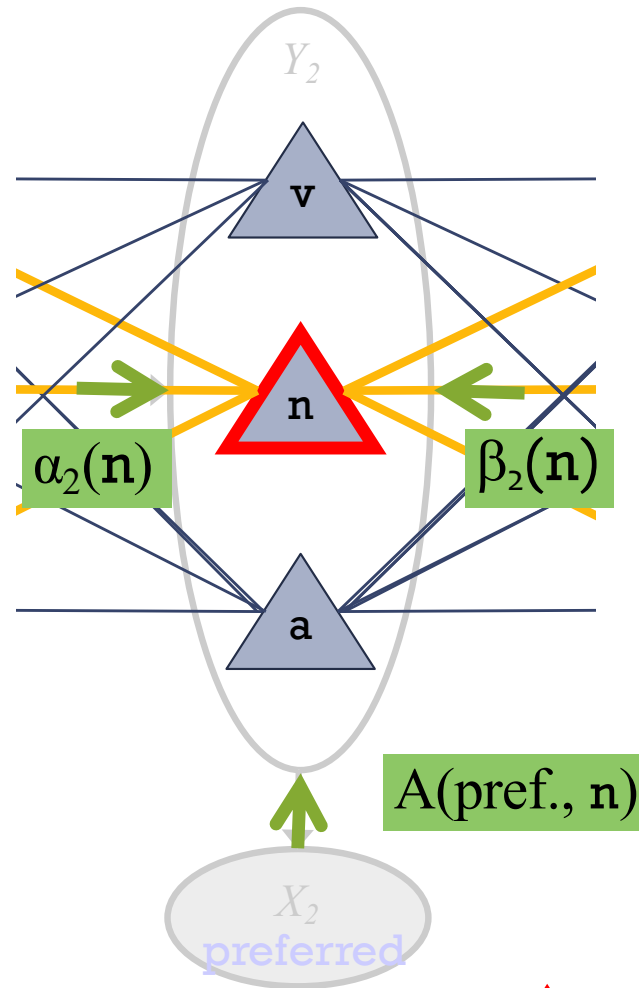
Product gives $ax+ay+az+bx+by+bz+cx+cy+cz$ = total weight of paths ²¹

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.

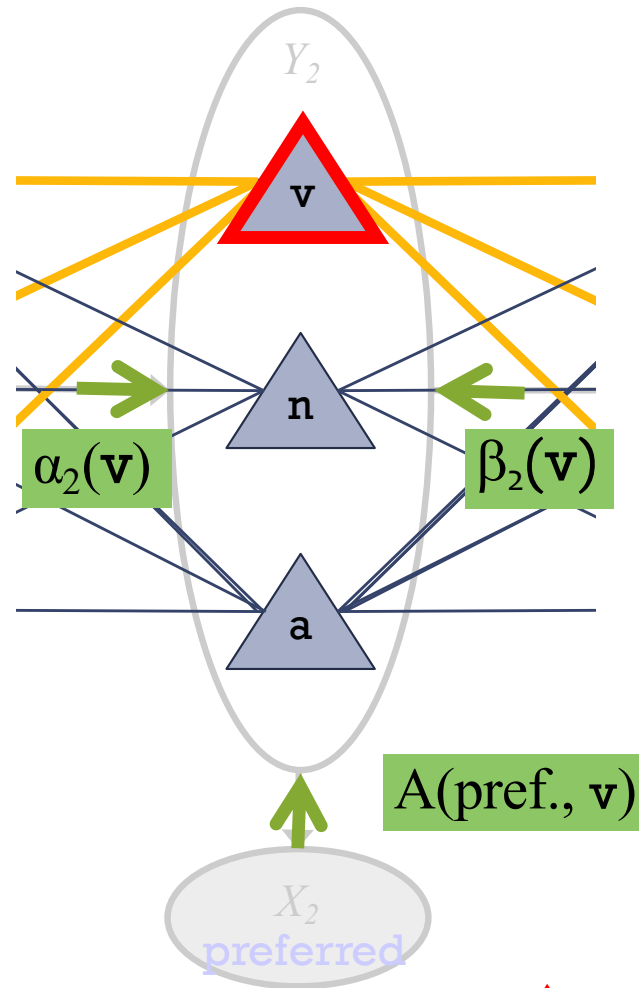


“belief that $Y_2 = \mathbf{n}$ ”

total weight of *all paths through* 

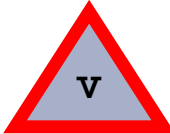
$$= \alpha_2(\mathbf{n}) A(\text{pref.}, \mathbf{n}) \beta_2(\mathbf{n})$$

Forward-Backward Algorithm: Finds Marginals



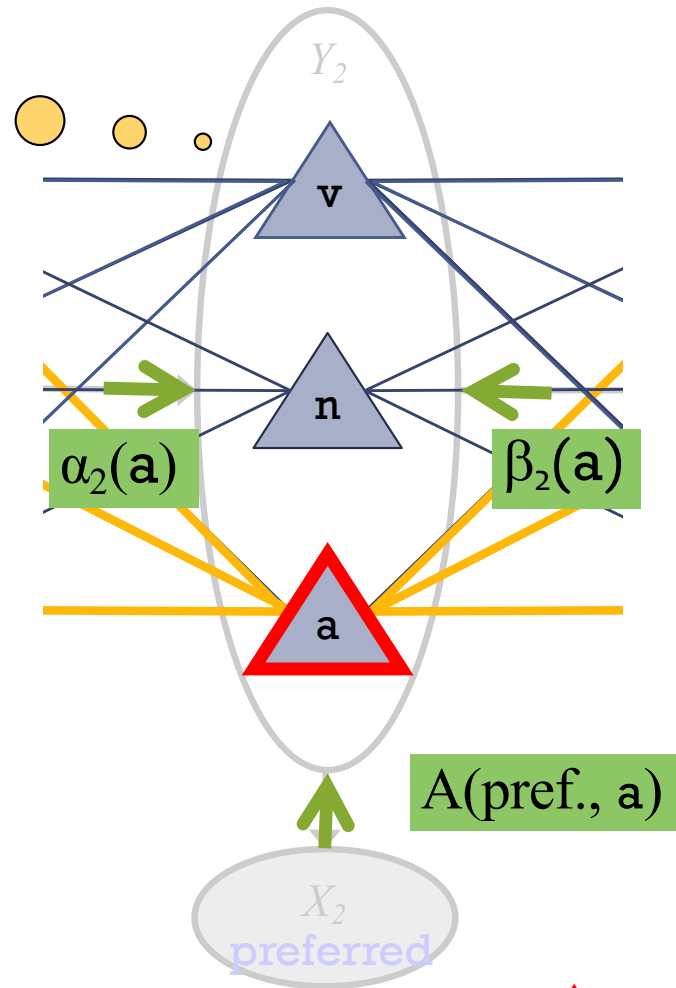
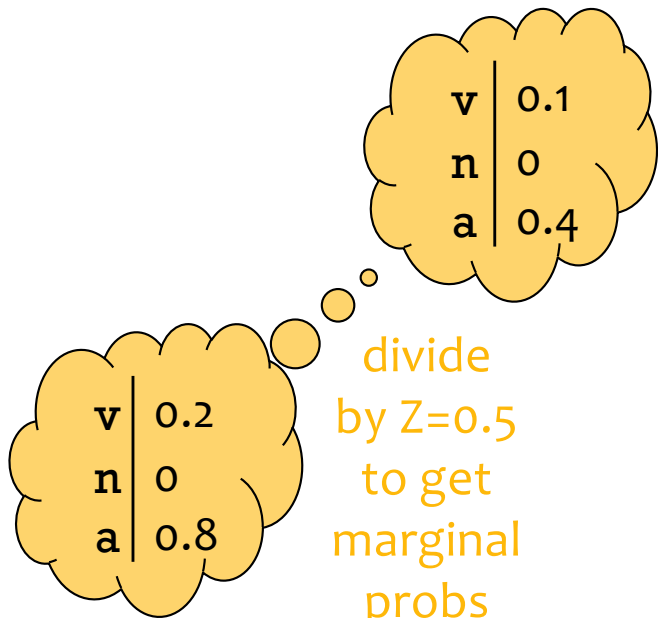
“belief that $Y_2 = \mathbf{v}$ ”

“belief that $Y_2 = \mathbf{n}$ ”

total weight of *all paths through* 

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

Forward-Backward Algorithm: Finds Marginals

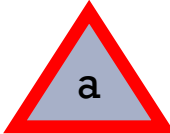


“belief that $Y_2 = \mathbf{v}$ ”

“belief that $Y_2 = \mathbf{n}$ ”

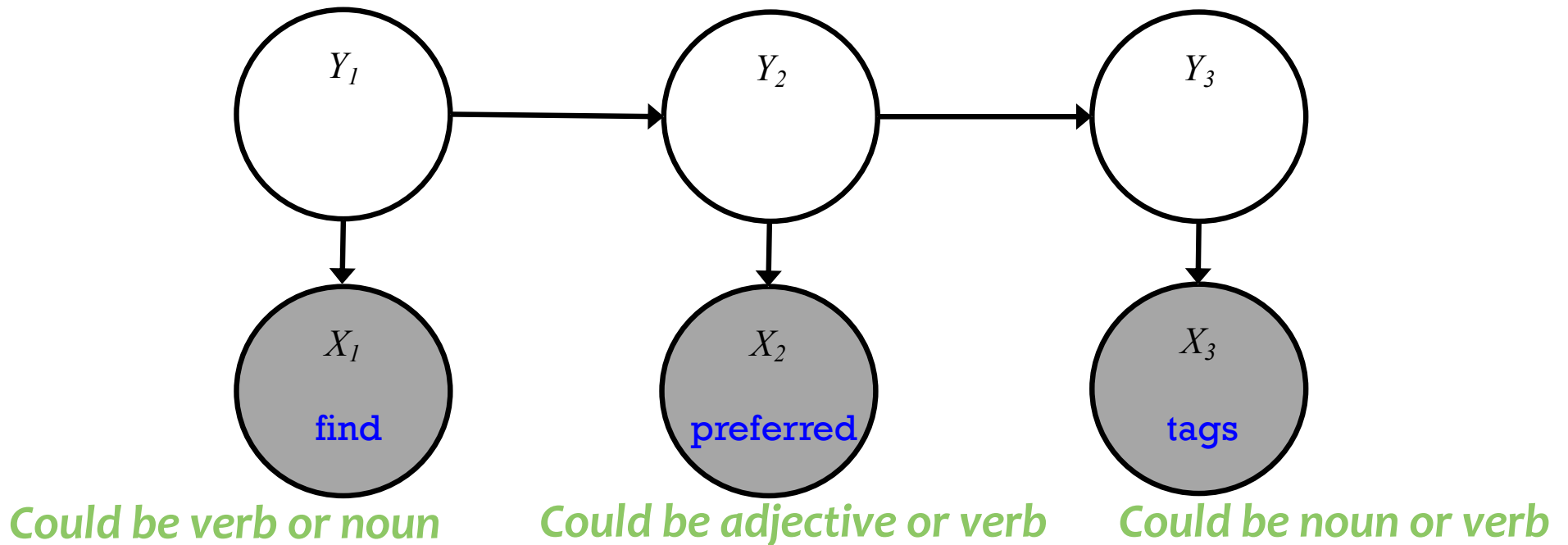
“belief that $Y_2 = \mathbf{a}$ ”

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

total weight of *all* paths through 

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

Forward-Backward Algorithm



Forward-Backward Algorithm

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

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

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

② For $t = 1, \dots, T$:

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

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

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

③ For $t = T, \dots, 1$:

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

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

$O(K)$

$O(K^2T)$

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

[Evaluation]

⑤ Compute $p(y_t = k | \vec{x}) = \frac{\alpha_t(k) \beta_t(k)}{p(\vec{x})}$

[Marginals]

Brute force
 algorithm
 would be
 $O(K^T)$

THE VITERBI ALGORITHM

Inference for HMMs

Whiteboard

- Viterbi algorithm
(edge weights version)

Viterbi Algorithm

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

"back pointers" $\rightarrow b_t(k) \triangleq \text{argmax}_{y_1, \dots, y_{t-1}} p(x_1, \dots, x_t, y_1, \dots, y_{t-1}, y_t = k)$

Assume $y_0 = \text{START}$

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

② For $t = 1, \dots, T$:

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

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

$$b_t(k) = \text{argmax}_{j \in \{1, \dots, K\}} p(x_t | y_t = k) \omega_{t-1}(j) p(y_t = k | y_{t-1} = j)$$

$O(K^2T)$

③ Compute Most Probable Assignment

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

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

$$\hat{y}_t = b_{t+1}(\hat{y}_{t+1})$$

[Decoding]

Brute force algorithm would be $O(K^T)$

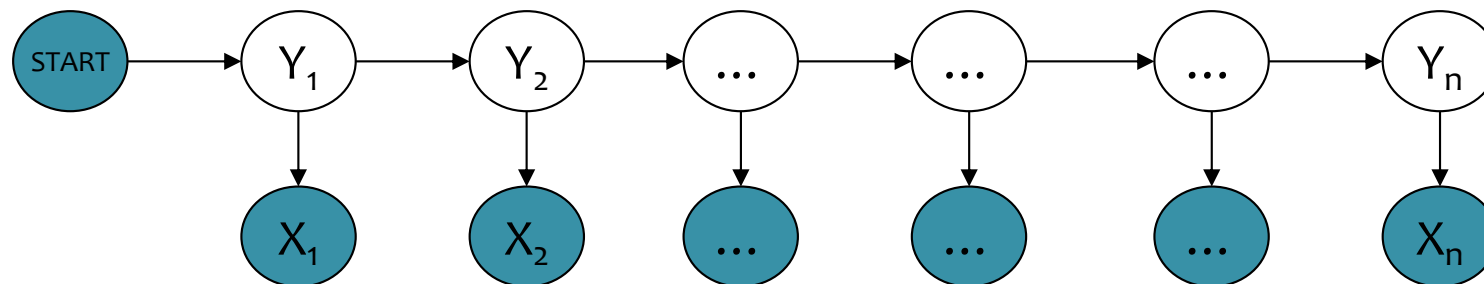
follow the "back pointers"

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

$$\begin{aligned} h_{\theta}(\mathbf{x}) &= \operatorname{argmin}_{\hat{\mathbf{y}}} \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\cdot|\mathbf{x})} [\ell(\hat{\mathbf{y}}, \mathbf{y})] \\ &= \operatorname{argmin}_{\hat{\mathbf{y}}} \sum_{\mathbf{y}} p_{\theta}(\mathbf{y} | \mathbf{x}) \ell(\hat{\mathbf{y}}, \mathbf{y}) \end{aligned}$$

Minimum Bayes Risk Decoding

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

Consider some example loss functions:

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

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

The MBR decoder is:

$$\begin{aligned} h_{\theta}(\mathbf{x}) &= \operatorname{argmin}_{\hat{\mathbf{y}}} \sum_{\mathbf{y}} p_{\theta}(\mathbf{y} | \mathbf{x}) (1 - \mathbb{I}(\hat{\mathbf{y}}, \mathbf{y})) \\ &= \operatorname{argmax}_{\hat{\mathbf{y}}} p_{\theta}(\hat{\mathbf{y}} | \mathbf{x}) \end{aligned}$$

which is exactly the Viterbi decoding problem!

Minimum Bayes Risk Decoding

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

Consider some example loss functions:

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

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

The MBR decoder is:

$$\hat{y}_i = h_{\theta}(\mathbf{x})_i = \operatorname{argmax}_{\hat{y}_i} p_{\theta}(\hat{y}_i | \mathbf{x})$$

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

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

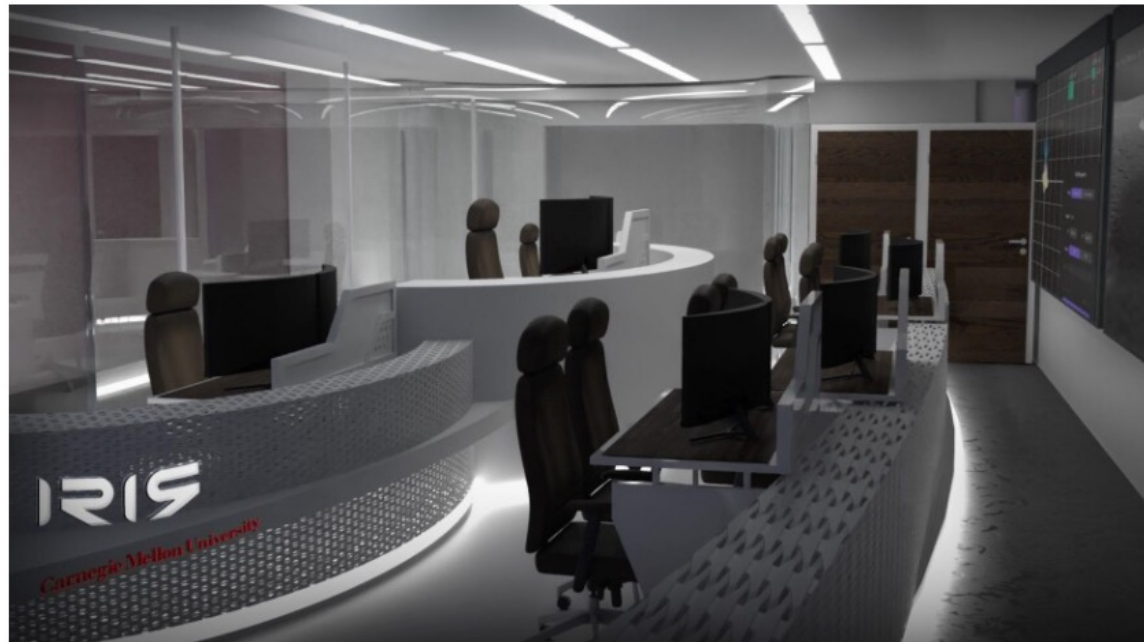
Example: CMU Mission Control

≡ 90.5 WESA
Pittsburgh's NPR News Station

▶ WESA
Morning Edition

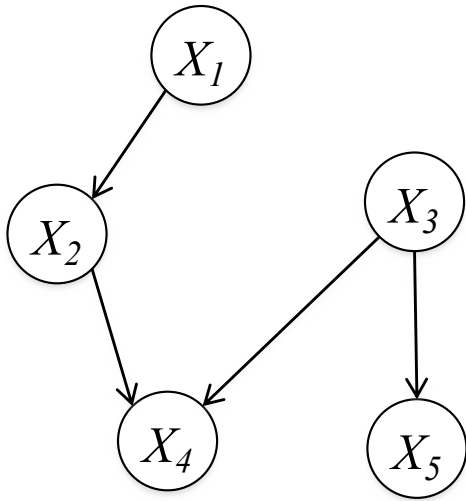
Pittsburgh's first mission control center to land at CMU ahead of 2022 lunar rover launch

90.5 WESA | By [Kiley Koscinski](#)
Published March 29, 2022 at 4:44 PM EDT



Courtesy Of Carnegie Mellon University

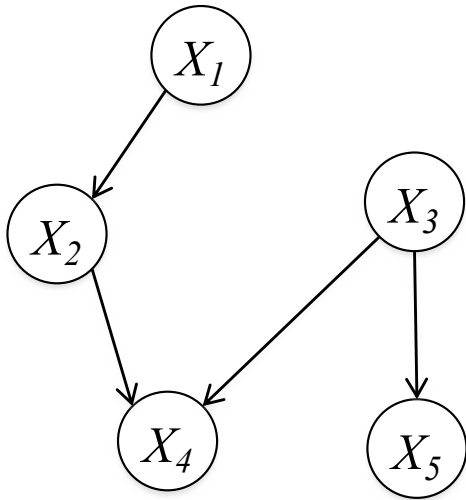
Bayesian Network



$$\begin{aligned} p(X_1, X_2, X_3, X_4, X_5) = & \\ & p(X_5|X_3)p(X_4|X_2, X_3) \\ & p(X_3)p(X_2|X_1)p(X_1) \end{aligned}$$

Bayesian Network

Definition:



$$P(X_1, \dots, X_T) = \prod_{t=1}^T P(X_t \mid \text{parents}(X_t))$$

- 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)
 - ...

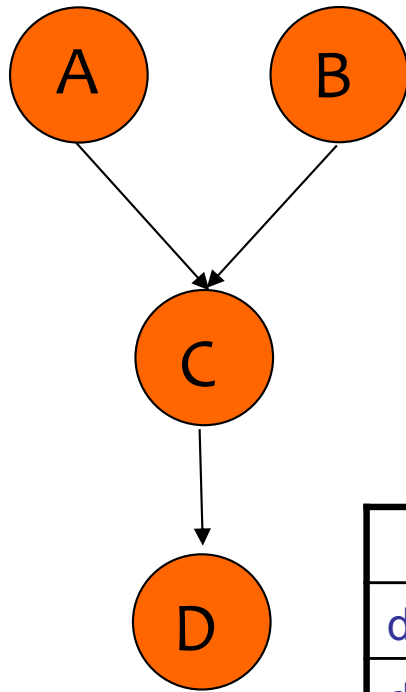
Quantitative Specification

Example: Conditional probability tables (CPTs)
for discrete random variables

a^0	0.75
a^1	0.25

b^0	0.33
b^1	0.67

$$P(a,b,c,d) = P(a)P(b)P(c|a,b)P(d|c)$$



	a^0b^0	a^0b^1	a^1b^0	a^1b^1
c^0	0.45	1	0.9	0.7
c^1	0.55	0	0.1	0.3

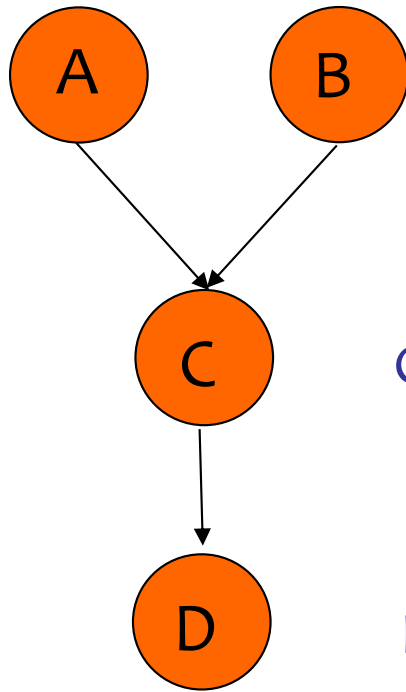
	c^0	c^1
d^0	0.3	0.5
d^1	0.7	0.5

Quantitative Specification

Example: Conditional probability density functions (CPDs)
for continuous random variables

$$A \sim N(\mu_a, \Sigma_a) \quad B \sim N(\mu_b, \Sigma_b)$$

$$P(a,b,c,d) = P(a)P(b)P(c|a,b)P(d|c)$$



$$C \sim N(A+B, \Sigma_c)$$

$$D \sim N(\mu_d + C, \Sigma_d)$$

$P(D|C)$

D

C

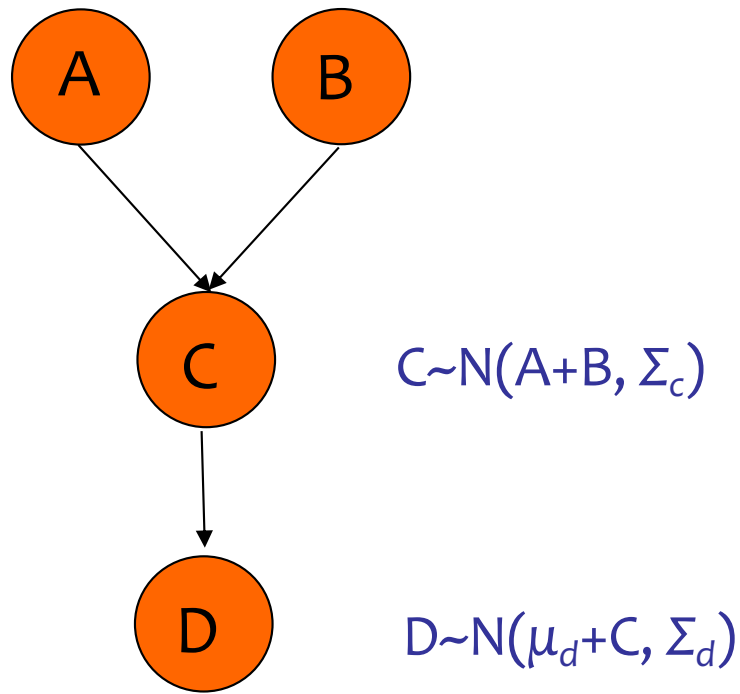
Quantitative Specification

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

a^0	0.75
a^1	0.25

b^0	0.33
b^1	0.67

$$P(a,b,c,d) = P(a)P(b)P(c|a,b)P(d|c)$$

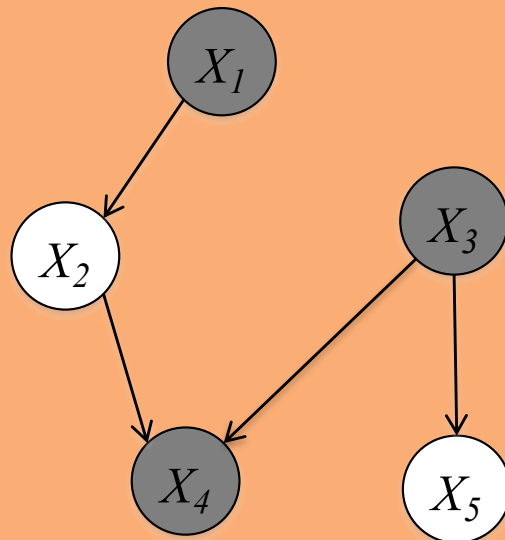


Observed Variables

- 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
2. Linear Regression
3. Bernoulli Naïve Bayes
4. Gaussian Naïve Bayes
5. 1D Gaussian

Answer:

