

GRAPHICAL MODELS

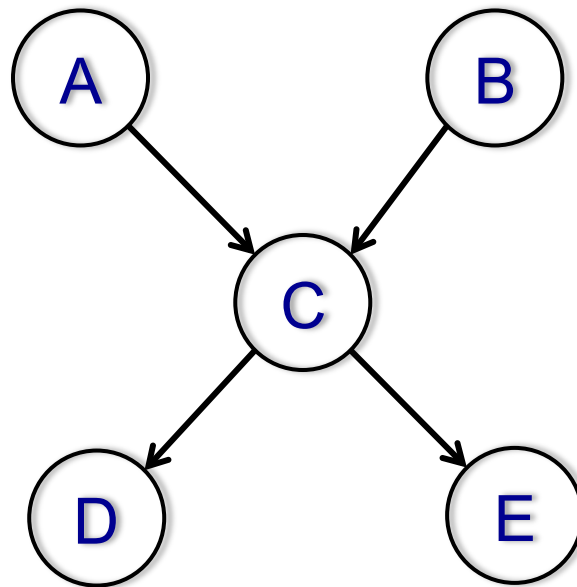
TRU CAO

**HO CHI MINH CITY UNIVERSITY OF TECHNOLOGY
AND JOHN VON NEUMANN INSTITUTE**

OUTLINE

- **Bayesian Networks (revisited)**
- **Naïve Bayes Classifier (revisited)**
- **Tree Augmented Naïve Bayes Model**
- **Hidden Markov Model**

BAYESIAN NETWORKS



BAYESIAN NETWORKS

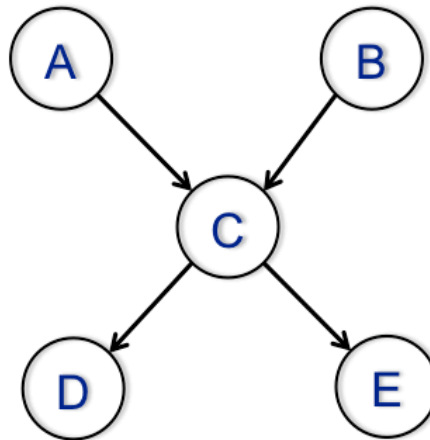
- **Advantages of graphical modeling:**

- Conditional independence:

$$p(D \mid C, E, A, B) = p(D \mid C)$$

- Factorization:

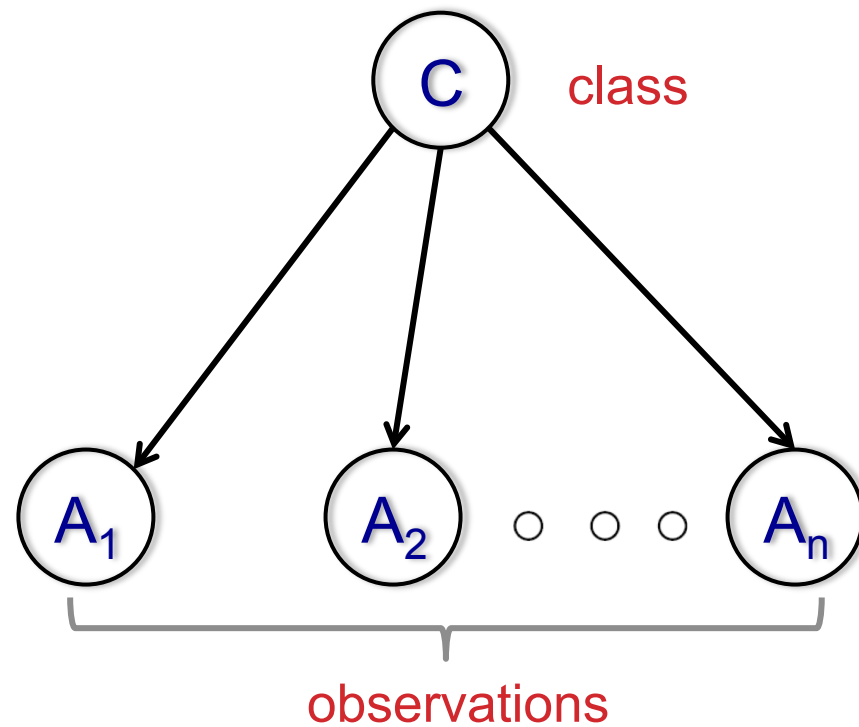
$$p(A, B, C, D, E) = p(D \mid C) \cdot p(E \mid C) \cdot p(C \mid A, B) \cdot p(A) \cdot p(B)$$



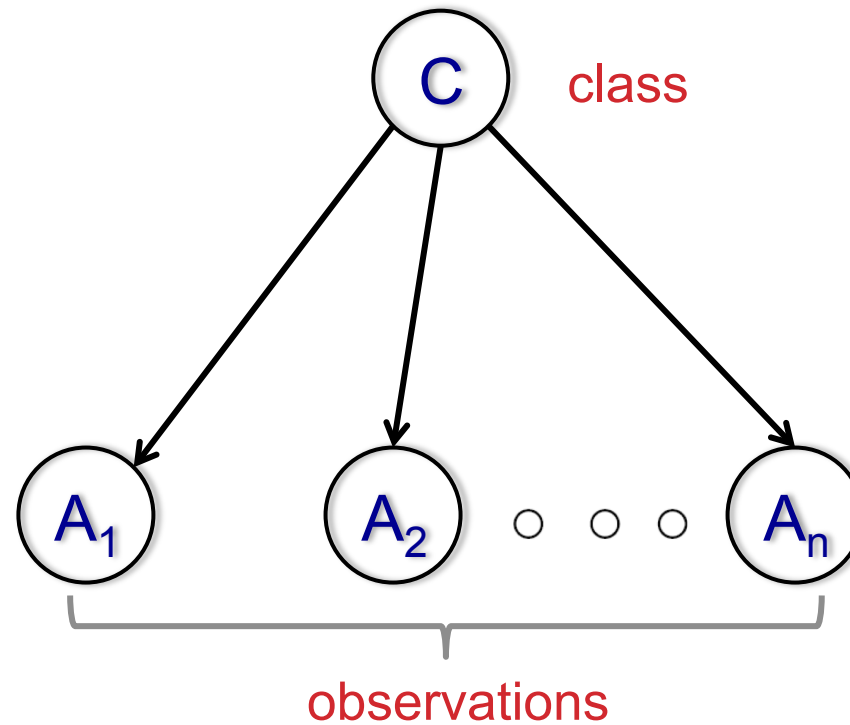
NAÏVE BAYES CLASSIFIER

- Each instance x is described by a conjunction of attribute values $\langle a_1, a_2, \dots, a_n \rangle$.
- It is to assign the most probable class c to an instance.
- $c_{NB} = \operatorname{argmax}_{c \in C} p(a_1, a_2, \dots, a_n | c) \cdot p(c)$
 $= \operatorname{argmax}_{c \in C} \prod_{i=1, n} p(a_i | c) \cdot p(c)$

NAÏVE BAYES CLASSIFIER



NAÏVE BAYES CLASSIFIER



Joint distribution: $p(C, A_1, A_2, \dots, A_n)$

NAÏVE BAYES CLASSIFIER

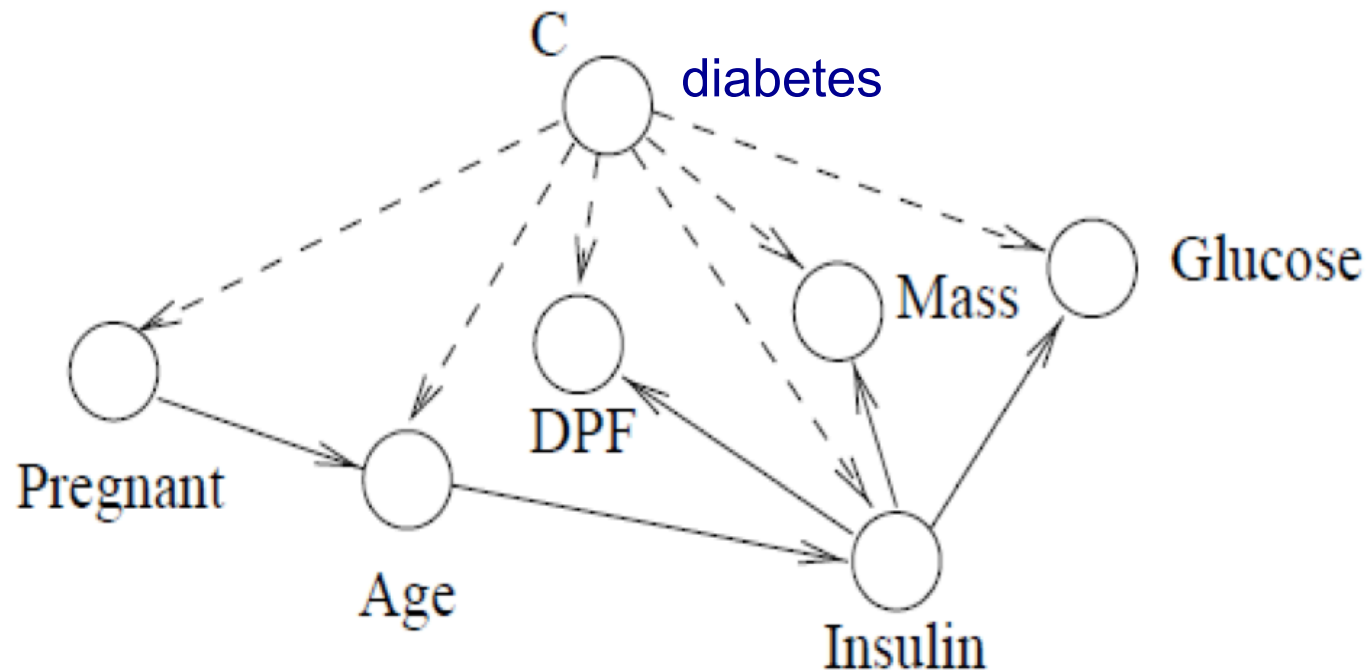
- **NB is a generative model:**
 - It models a joint distribution: $p(C, A)$
 - It can generate any distribution on C and A .

NAÏVE BAYES CLASSIFIER

- **NB is a generative model:**
 - It models a joint distribution: $p(C, A)$
 - It can generate any distribution on C and A .
- **In contrast to a discriminative model (e.g., CRF):**
 - Conditional distribution: $p(C | A)$
 - It discriminates C given A .

TREE AUGMENTED NB MODEL

- An extension of NB with dependence between attributes/observations:



HIDDEN MARKOV MODELS

- Introduction
- Example
- Independence assumptions
- Forward algorithm
- Viterbi algorithm
- Training
- Application to NER

HIDDEN MARKOV MODELS

- One of the most popular graphical models.
- Dynamic extension of Bayesian networks.
- Sequential extension of NB classifier.

HIDDEN MARKOV MODELS

- **Example:**
 - Your possible looking prior to the exam = {**tired**, **hungover**, **scared**, **fine**}.
 - Your possible activity last night = {**TV**, **pub**, **party**, **study**}.
 - Given a sequence of observations of your looking, guess what you did in previous nights.

HIDDEN MARKOV MODELS

- **Example:**

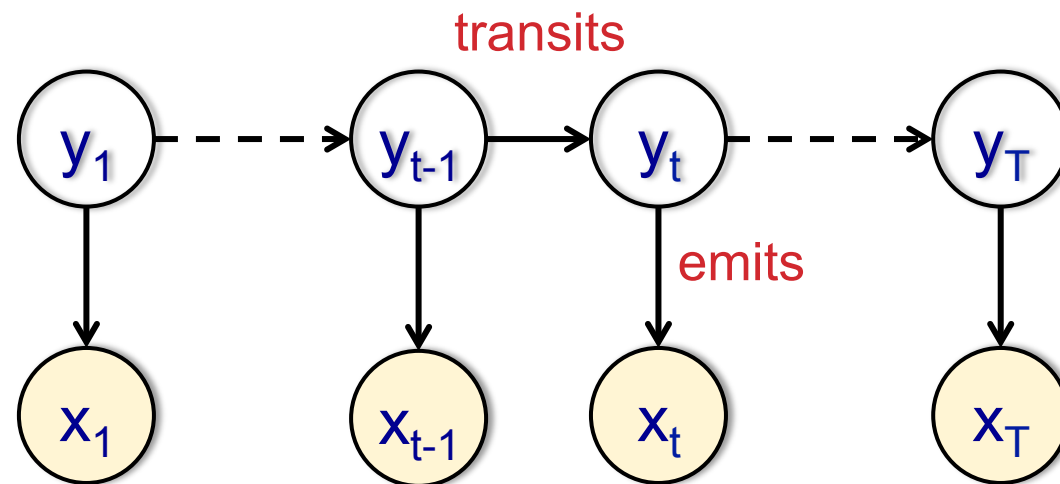
- Your possible looking prior to the exam = {**tired**, **hungover**, **scared**, **fine**}.
- Your possible activity last night = {**TV**, **pub**, **party**, **study**}.
- Given a sequence of observations of your looking, guess what you did in previous nights.

- **A model:**

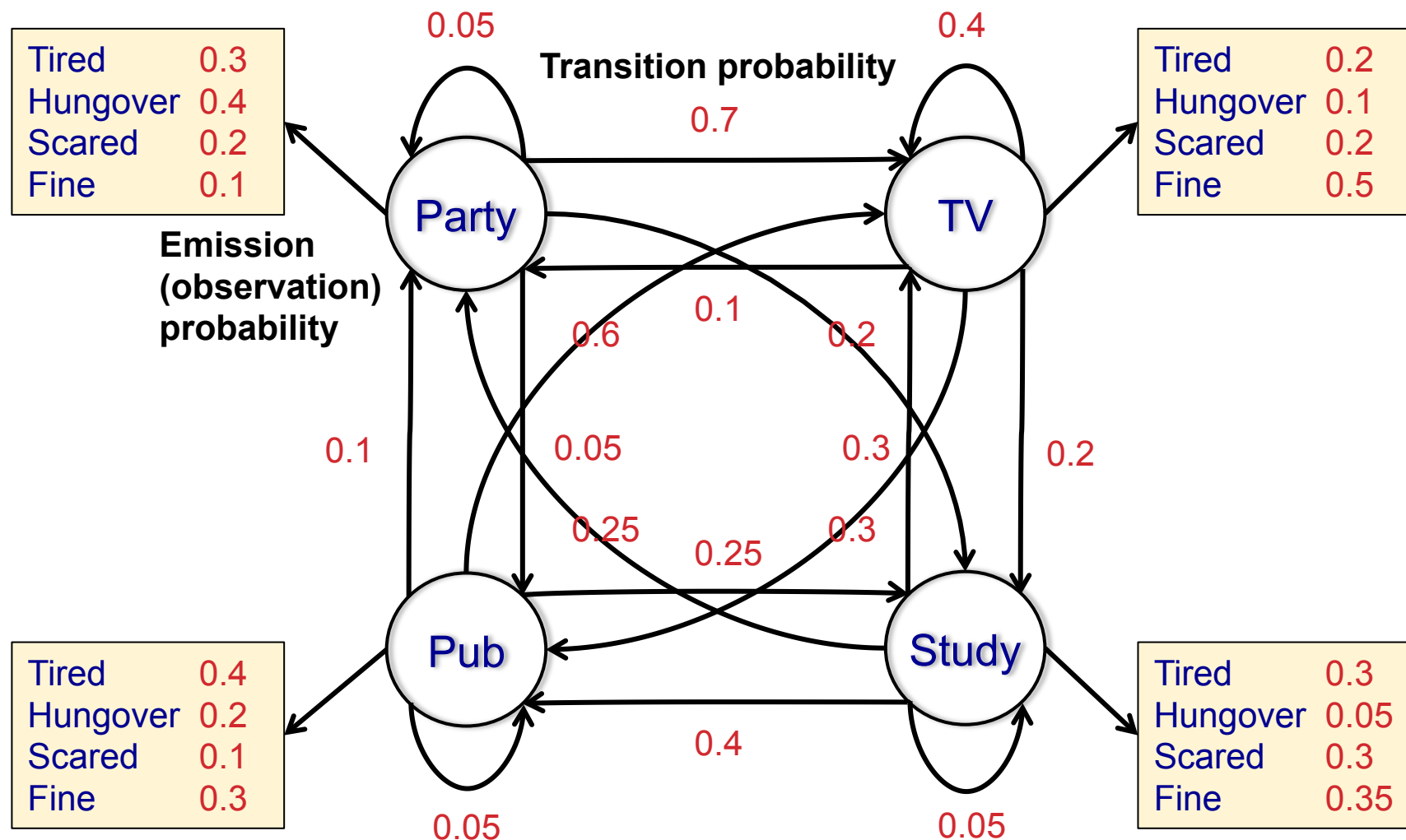
- Your looking depends on what you did in the night before.
- Your activity in a night depends on what you did in some previous nights.

HIDDEN MARKOV MODELS

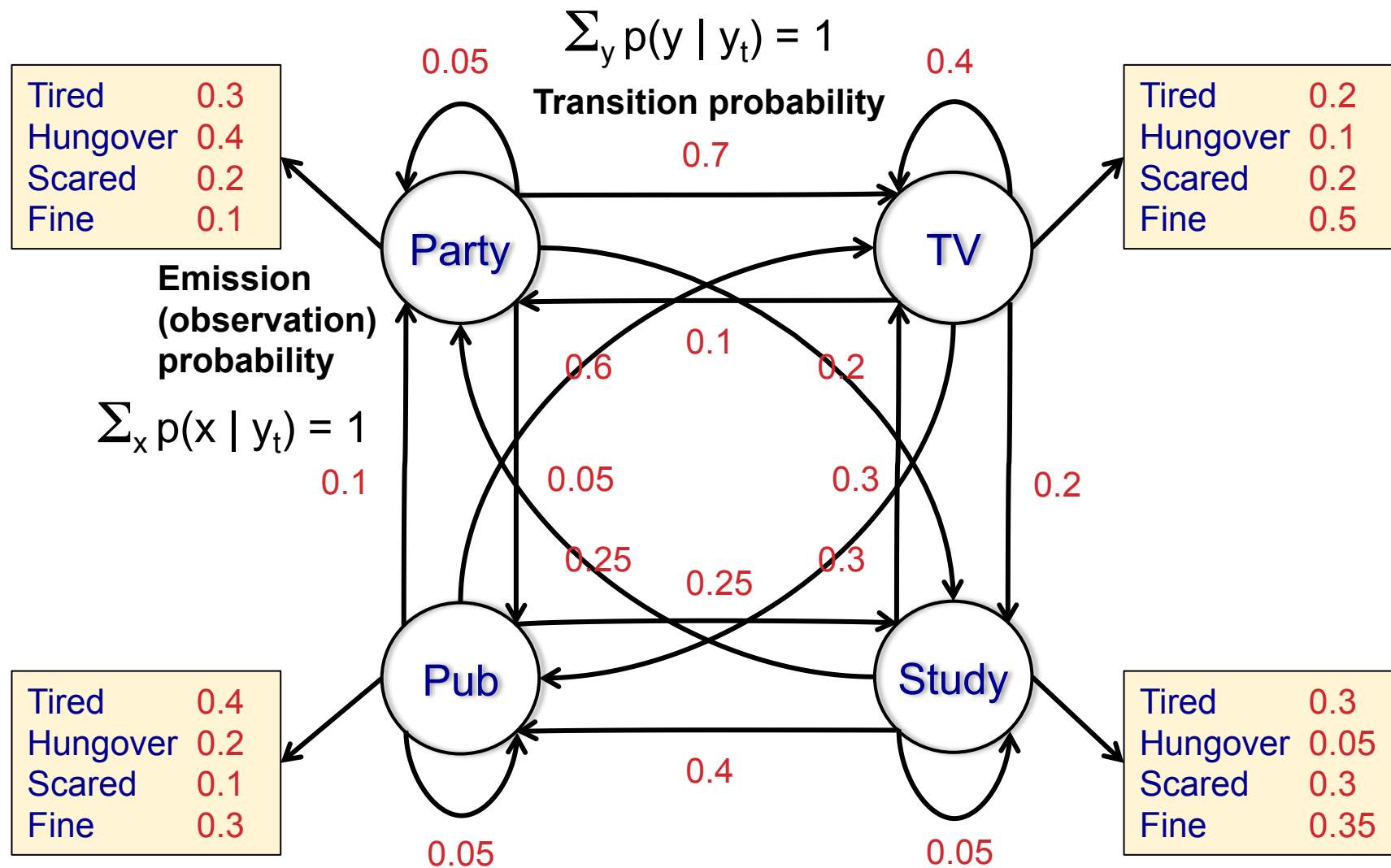
- A finite set of possible **observations**.
- A finite set of possible hidden **states**.
- To predict the most probable sequence of underlying states $\{y_1, y_2, \dots, y_T\}$ for a given sequence of observations $\{x_1, x_2, \dots, x_T\}$.



HIDDEN MARKOV MODELS



HIDDEN MARKOV MODELS



HIDDEN MARKOV MODELS

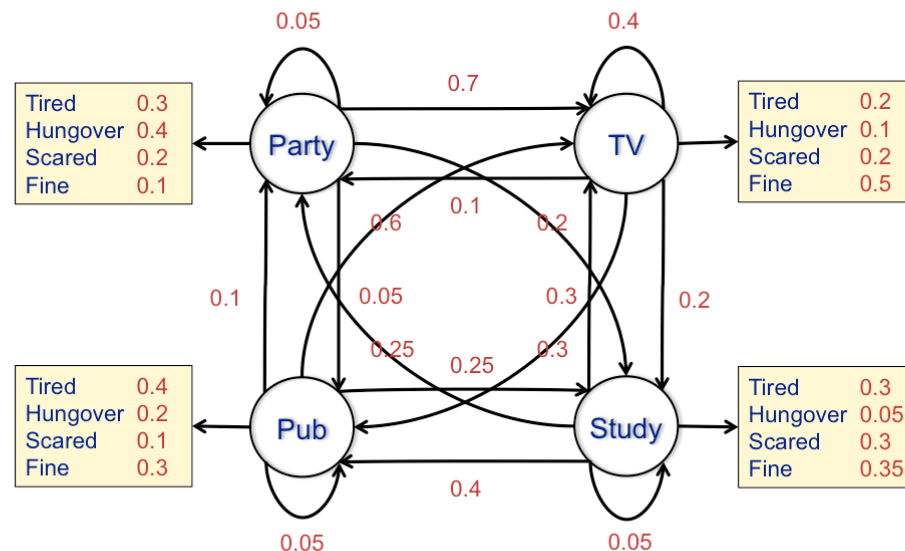
- **HMM conditional independence assumptions:**

- State at time t depends only on state at time $t - 1$.

$$p(y_t | y_{t-1}, Z) = p(y_t | y_{t-1})$$

- Observation at time t depends only on state at time t .

$$P(x_t | y_t, Z) = p(x_t | y_t)$$



HIDDEN MARKOV MODELS

- **HMM is a generative model:**

- Joint distributions:

$$p(Y, X) = p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T) = \prod_{t=1, T} p(x_t | y_t) \cdot p(y_t | y_{t-1})$$

HIDDEN MARKOV MODELS

- **HMM is a generative model:**

- Joint distributions:

$$p(Y, X) = p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T) = \prod_{t=1, T} p(x_t | y_t) \cdot p(y_t | y_{t-1})$$

$$p(y_1 | y_0) = p(y_1)$$

HIDDEN MARKOV MODELS

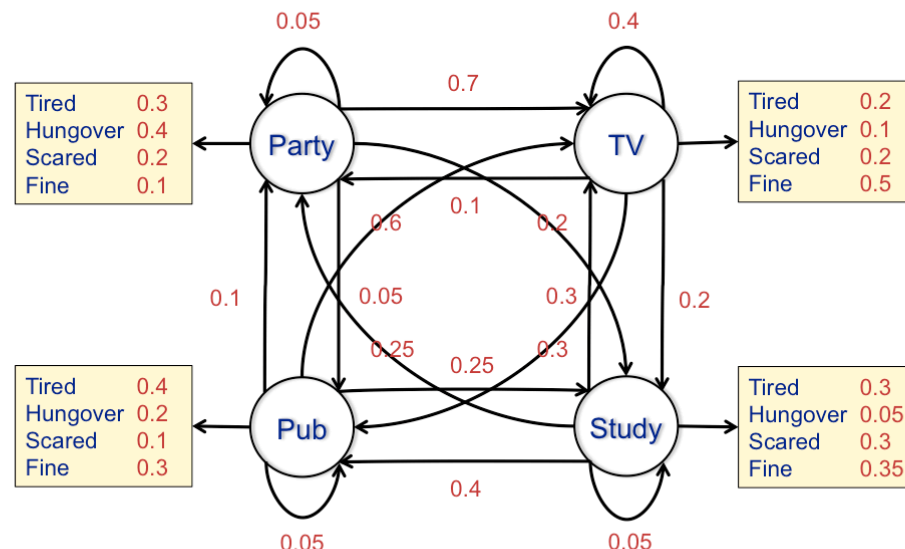
- **HMM is a generative model:**

- Joint distributions:

$$p(Y, X) = p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T) = \prod_{t=1, T} p(y_t | y_{t-1}) \cdot p(x_t | y_t)$$

$$p(y_1 | y_0) = p(y_1)$$

- It can generate any distribution on Y and X.



HIDDEN MARKOV MODELS

- **HMM is a generative model:**

- Joint distributions:

$$p(Y, X) = p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T) = \prod_{t=1, T} p(y_t | y_{t-1}) \cdot p(x_t | y_t)$$

$$p(y_1 | y_0) = p(y_1)$$

- It can generate any distribution on Y and X .
- **In contrast to a discriminative model (e.g., CRF):**
 - Conditional distributions: $p(Y | X)$
 - It discriminates Y given X .

HIDDEN MARKOV MODELS

- **Forward algorithm:**
 - To compute the joint probability of the state at time t being y_t and the sequence of observations in the first t steps being $\{x_1, x_2, \dots, x_t\}$:

$$\alpha_t(y_t) = p(y_t, x_1, x_2, \dots, x_t)$$

HIDDEN MARKOV MODELS

- **Forward algorithm:**

- To compute the joint probability of the state at time t being y_t and the sequence of observations in the first t steps being $\{x_1, x_2, \dots, x_t\}$:

$$\alpha_t(y_t) = p(y_t, x_1, x_2, \dots, x_t)$$

- Bayes' theorem gives:

$$\begin{aligned} & p(y_t \mid x_1, x_2, \dots, x_t) \\ &= p(y_t, x_1, x_2, \dots, x_t) / p(x_1, x_2, \dots, x_t) \\ &= \alpha_t(y_t) / p(x_1, x_2, \dots, x_t) \end{aligned}$$

HIDDEN MARKOV MODELS

- **Forward algorithm:**

- To compute the joint probability of the state at time t being y_t and the sequence of observations in the first t steps being $\{x_1, x_2, \dots, x_t\}$:

$$\alpha_t(y_t) = p(y_t, x_1, x_2, \dots, x_t)$$

- Bayes' theorem gives:

$$\begin{aligned} p(y_t \mid x_1, x_2, \dots, x_t) \\ &= p(y_t, x_1, x_2, \dots, x_t) / p(x_1, x_2, \dots, x_t) \\ &= \alpha_t(y_t) / p(x_1, x_2, \dots, x_t) \end{aligned}$$

- The highest $\alpha_t(y_t)$ is, the most likely y_t would be given the same $\{x_1, x_2, \dots, x_t\}$.

HIDDEN MARKOV MODELS

- **Forward algorithm:**

$$\begin{aligned}\alpha_t(y_t) &= p(y_t, x_1, x_2, \dots, x_t) \\ &= \sum_{y_{t-1}} p(y_t, y_{t-1}, x_1, x_2, \dots, x_t) \\ &= \sum_{y_{t-1}} p(x_t | y_t, y_{t-1}, x_1, x_2, \dots, x_{t-1}) \cdot p(y_t, y_{t-1}, x_1, x_2, \dots, x_{t-1}) \\ &= \sum_{y_{t-1}} p(x_t | y_t) \cdot p(y_t | y_{t-1}, x_1, x_2, \dots, x_{t-1}) \cdot p(y_{t-1}, x_1, x_2, \dots, x_{t-1}) \\ &= \sum_{y_{t-1}} p(x_t | y_t) \cdot p(y_t | y_{t-1}) \cdot p(y_{t-1}, x_1, x_2, \dots, x_{t-1}) \\ &= p(x_t | y_t) \sum_{y_{t-1}} p(y_t | y_{t-1}) \cdot \alpha_{t-1}(y_{t-1})\end{aligned}$$

HIDDEN MARKOV MODELS

- **Forward algorithm:**

$$\alpha_t(y_t)$$

$$\alpha_1(y_1) = p(y_1, x_1) = p(x_1 | y_1) \cdot p(y_1)$$

$$= p(y_t, x_1, x_2, \dots, x_t)$$

$$= \sum_{y_{t-1}} p(y_t, y_{t-1}, x_1, x_2, \dots, x_t)$$

$$= \sum_{y_{t-1}} p(x_t | y_t, y_{t-1}, x_1, x_2, \dots, x_{t-1}) \cdot p(y_t, y_{t-1}, x_1, x_2, \dots, x_{t-1})$$

$$= \sum_{y_{t-1}} p(x_t | y_t) \cdot p(y_t | y_{t-1}, x_1, x_2, \dots, x_{t-1}) \cdot p(y_{t-1}, x_1, x_2, \dots, x_{t-1})$$

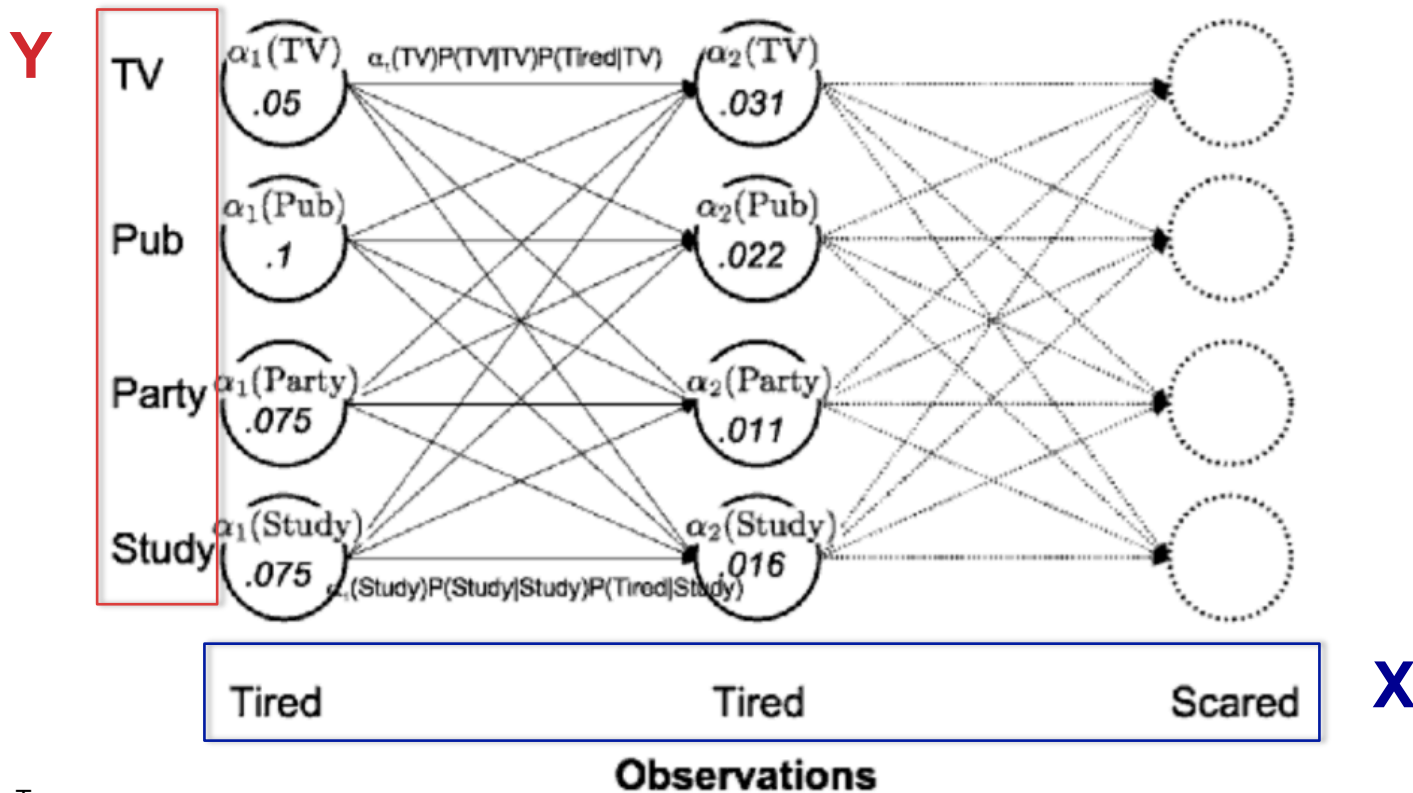
$$= \sum_{y_{t-1}} p(x_t | y_t) \cdot p(y_t | y_{t-1}) \cdot p(y_{t-1}, x_1, x_2, \dots, x_{t-1})$$

$$= p(x_t | y_t) \sum_{y_{t-1}} p(y_t | y_{t-1}) \cdot \alpha_{t-1}(y_{t-1})$$

HIDDEN MARKOV MODELS

- Forward algorithm:

$$\alpha_t(y_t) = p(x_t | y_t) \sum_{y_{t-1}} p(y_t | y_{t-1}) \cdot \alpha_{t-1}(y_{t-1})$$



HIDDEN MARKOV MODELS

- **Viterbi algorithm:**

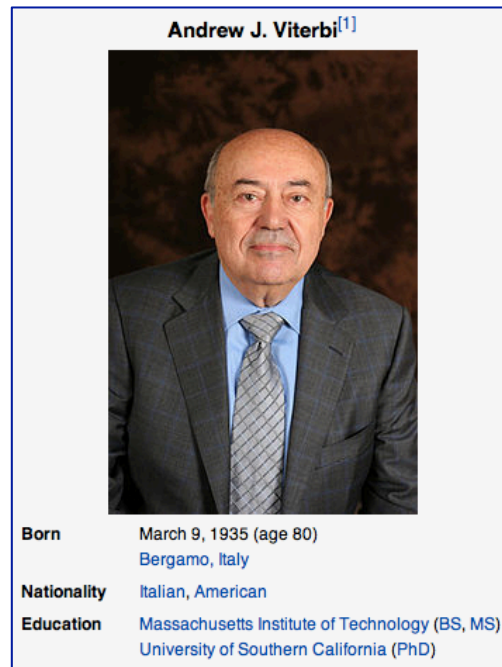
- To compute the most probable sequence of states $\{y_1, y_2, \dots, y_T\}$ given a sequence of observations $\{x_1, x_2, \dots, x_T\}$:

$$Y^* = \operatorname{argmax}_Y p(Y | X) = \operatorname{argmax}_Y p(Y, X)$$

HIDDEN MARKOV MODELS

- **Viterbi algorithm:**
 - To compute the most probable sequence of states $\{y_1, y_2, \dots, y_T\}$ given a sequence of observations $\{x_1, x_2, \dots, x_T\}$:

$$Y^* = \operatorname{argmax}_Y p(Y | X) = \operatorname{argmax}_Y p(Y, X)$$



10/28/16

Cao Hoang Tru

HIDDEN MARKOV MODELS

- **Viterbi algorithm:**

$$\begin{aligned} & \max_{y_{1:T}} p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T) \\ &= \max_{y_T} \max_{y_{1:T-1}} p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T) \\ &= \max_{y_T} \max_{y_{1:T-1}} \{p(x_T | y_T) \cdot p(y_T | y_{T-1}) \cdot p(y_1, \dots, y_{T-1}, x_1, \dots, x_{T-1})\} \\ &= \max_{y_T} \max_{y_{T-1}} \{p(x_T | y_T) \cdot p(y_T | y_{T-1}) \cdot \max_{y_{1:T-2}} p(y_1, \dots, y_{T-1}, x_1, \dots, x_{T-1})\} \end{aligned}$$

HIDDEN MARKOV MODELS

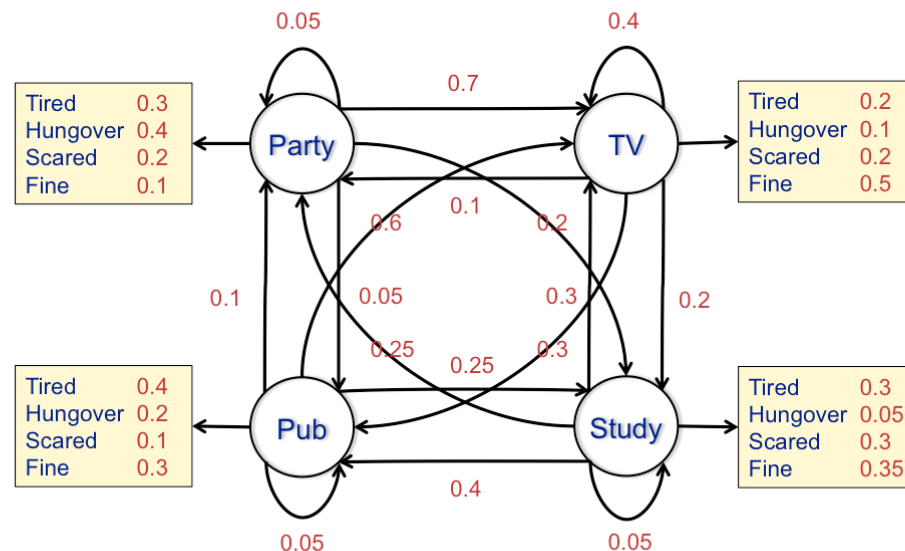
- Viterbi algorithm:

$$\max_{y_{1:T}} p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T)$$

$$= \max_{y_T} \max_{y_{1:T-1}} p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T)$$

$$= \max_{y_T} \max_{y_{1:T-1}} \{p(x_T | y_T) \cdot p(y_T | y_{T-1}) \cdot p(y_1, \dots, y_{T-1}, x_1, \dots, x_{T-1})\}$$

$$= \max_{y_T} \max_{y_{T-1}} \{p(x_T | y_T) \cdot p(y_T | y_{T-1}) \cdot \max_{y_{1:T-2}} p(y_1, \dots, y_{T-1}, x_1, \dots, x_{T-1})\}$$



HIDDEN MARKOV MODELS

- **Viterbi algorithm:**

$$\begin{aligned} & \max_{y_{1:T}} p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T) \\ &= \max_{y_T} \max_{y_{1:T-1}} p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T) \\ &= \max_{y_T} \max_{y_{1:T-1}} \{p(x_T | y_T) \cdot p(y_T | y_{T-1}) \cdot p(y_1, \dots, y_{T-1}, x_1, \dots, x_{T-1})\} \\ &= \max_{y_T} \max_{y_{T-1}} \{p(x_T | y_T) \cdot p(y_T | y_{T-1}) \cdot \max_{y_{1:T-2}} p(y_1, \dots, y_{T-1}, x_1, \dots, x_{T-1})\} \end{aligned}$$

- **Dynamic programming:**

- Compute

$$\operatorname{argmax}_{y_1} p(y_1, x_1) = \operatorname{argmax}_{y_1} p(x_1 | y_1) \cdot p(y_1)$$

- For each t from 2 to T , and for each state y_t , compute:

$$\operatorname{argmax}_{y_{1:t-1}} p(y_1, y_2, \dots, y_t, x_1, x_2, \dots, x_t)$$

- Select $\operatorname{argmax}_{y_{1:T}} p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T)$

HIDDEN MARKOV MODELS

- Dynamic programming:

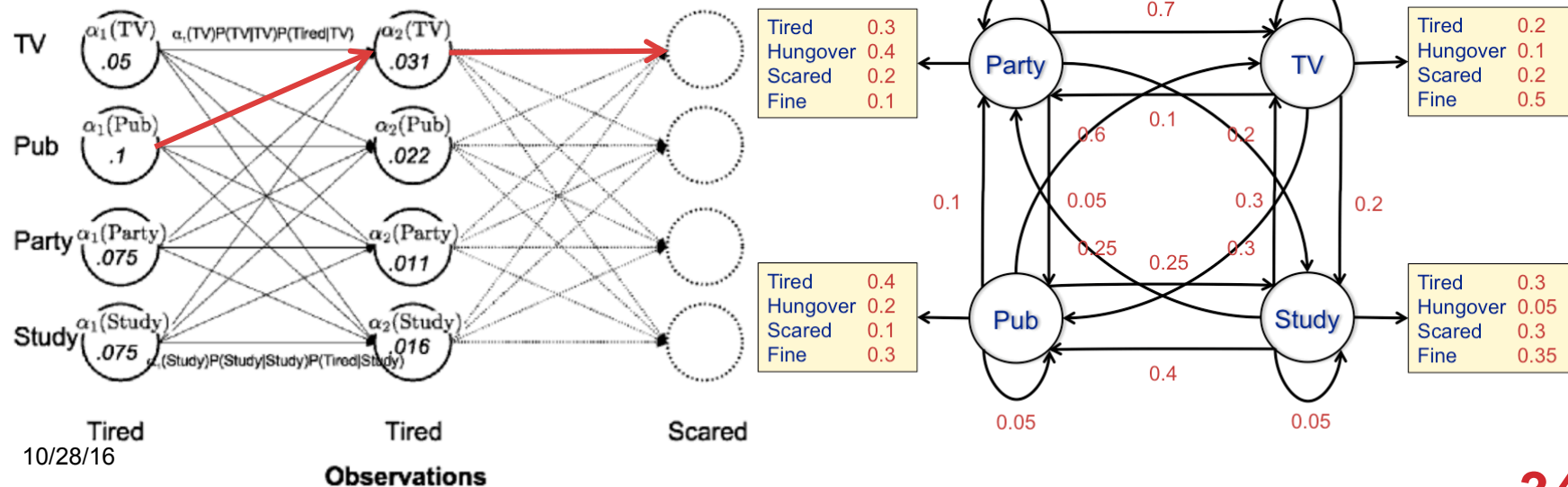
- Compute

$$\operatorname{argmax}_{y_1} p(y_1, x_1) = \operatorname{argmax}_{y_1} p(x_1 | y_1) \cdot p(y_1)$$

- For each t from 2 to T , and for each state y_t , compute:

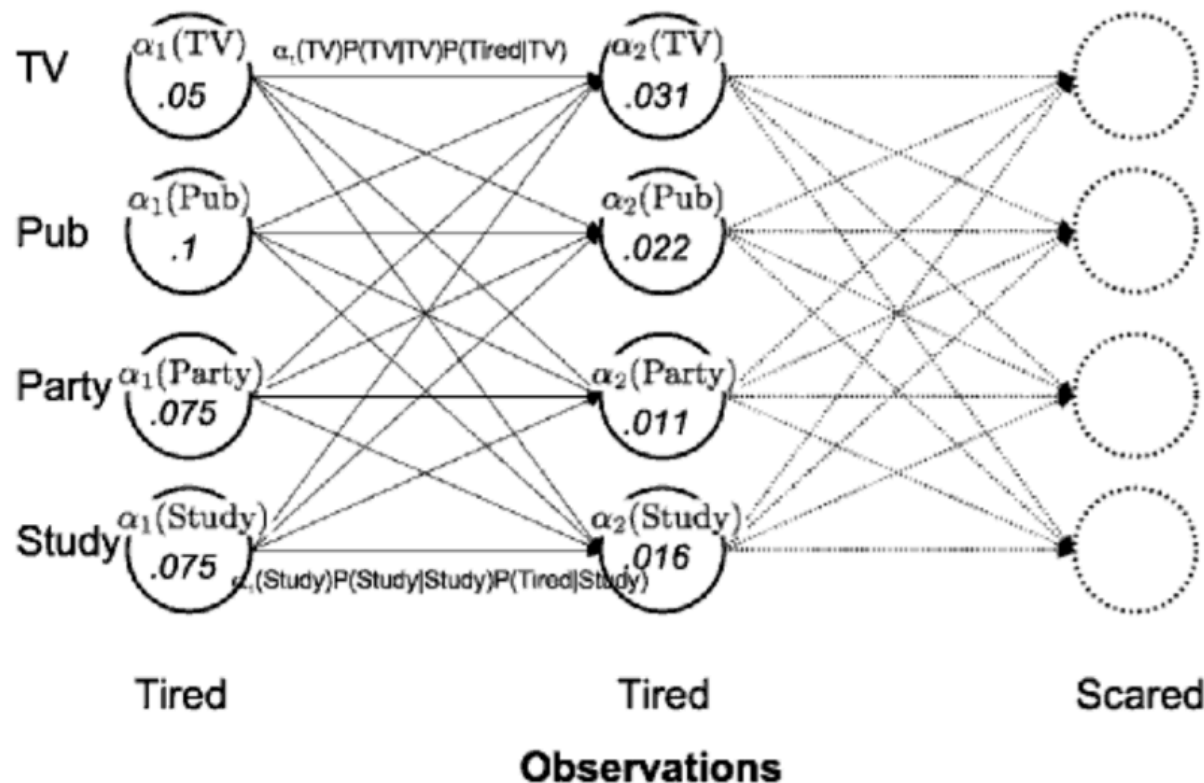
$$\operatorname{argmax}_{y_{1:t-1}} p(y_1, y_2, \dots, y_t, x_1, x_2, \dots, x_t)$$

- Select $\operatorname{argmax}_{y_{1:T}} p(y_1, y_2, \dots, y_T, x_1, x_2, \dots, x_T)$

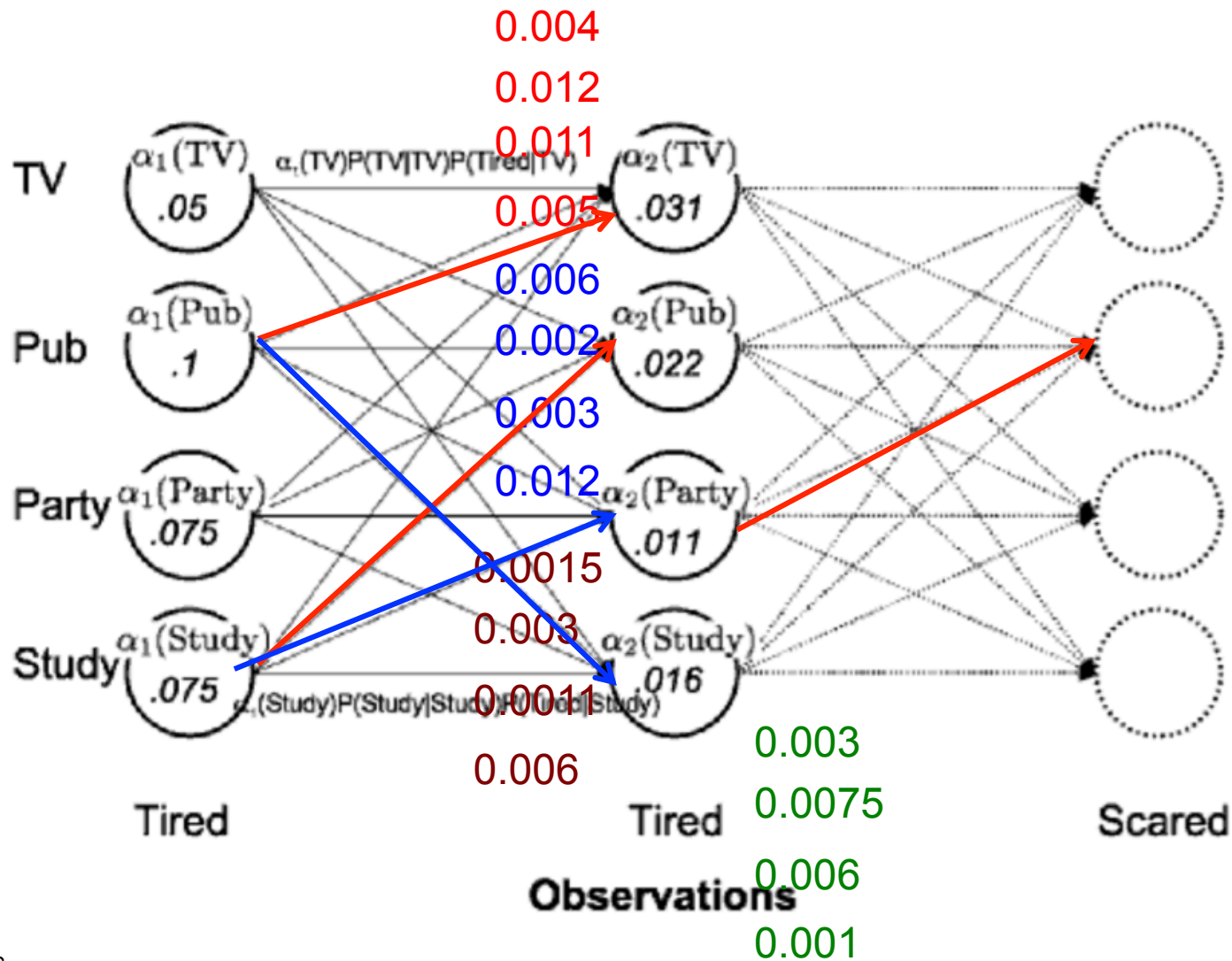


HIDDEN MARKOV MODELS

- Could the results from the forward algorithm be used for Viterbi algorithm?

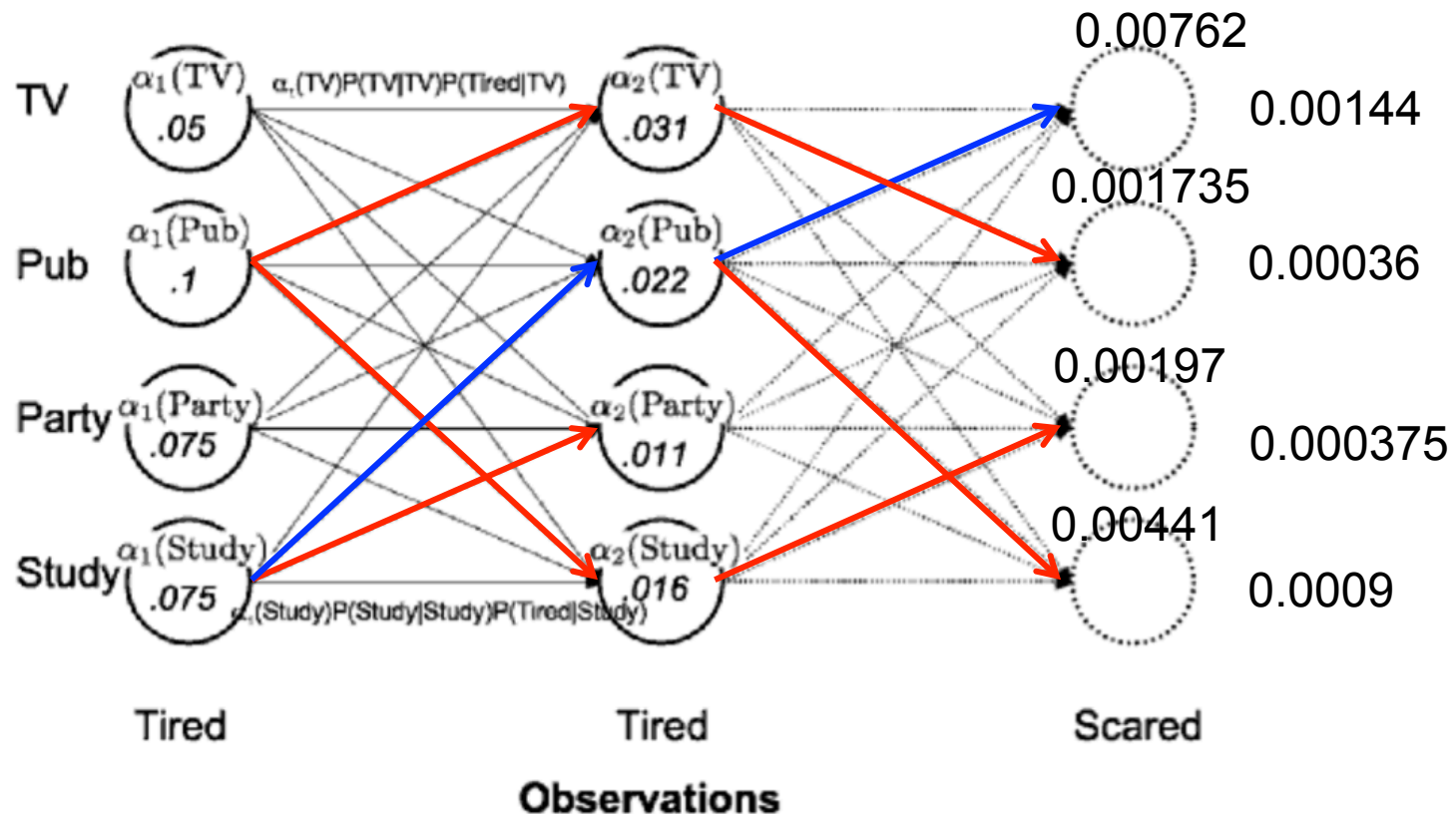


HIDDEN MARKOV MODELS



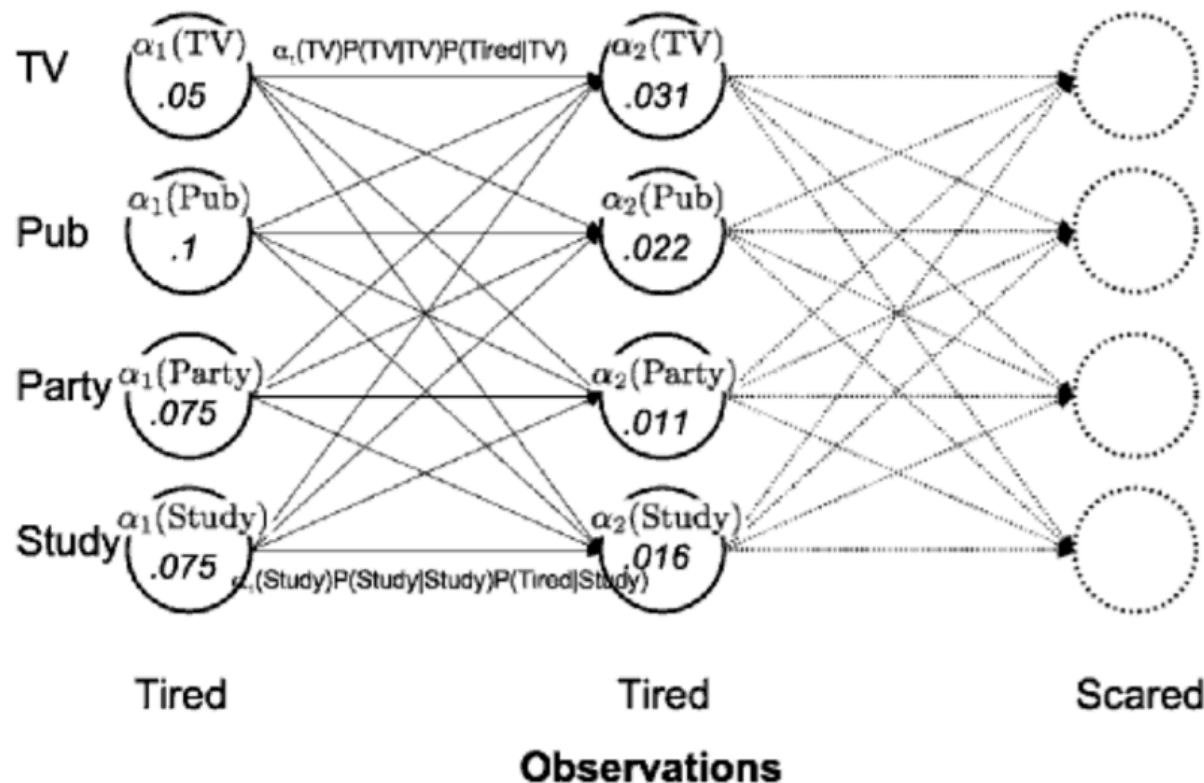
HIDDEN MARKOV MODELS

- Could the results from the forward algorithm be used for Viterbi algorithm?



HIDDEN MARKOV MODELS

- Could the results from the forward algorithm be used for Viterbi algorithm?



READING HOMEWORK 1

- **Marsland, S. (2009) Machine learning: An algorithmic perspective. Chapter 15 (graphical models).**

EXERCISES 1

- **Apply Viterbi algorithm to find the most probable 3-state sequence in the looking-activity example in the lecture.**

HIDDEN MARKOV MODELS

- Where does an HMM come from?

HIDDEN MARKOV MODELS

- **Training HMMs:**
 - Topology is designed beforehand.
 - Parameters to be learned: **emission** and **transition** probabilities.
 - Supervised or unsupervised training.

HIDDEN MARKOV MODELS

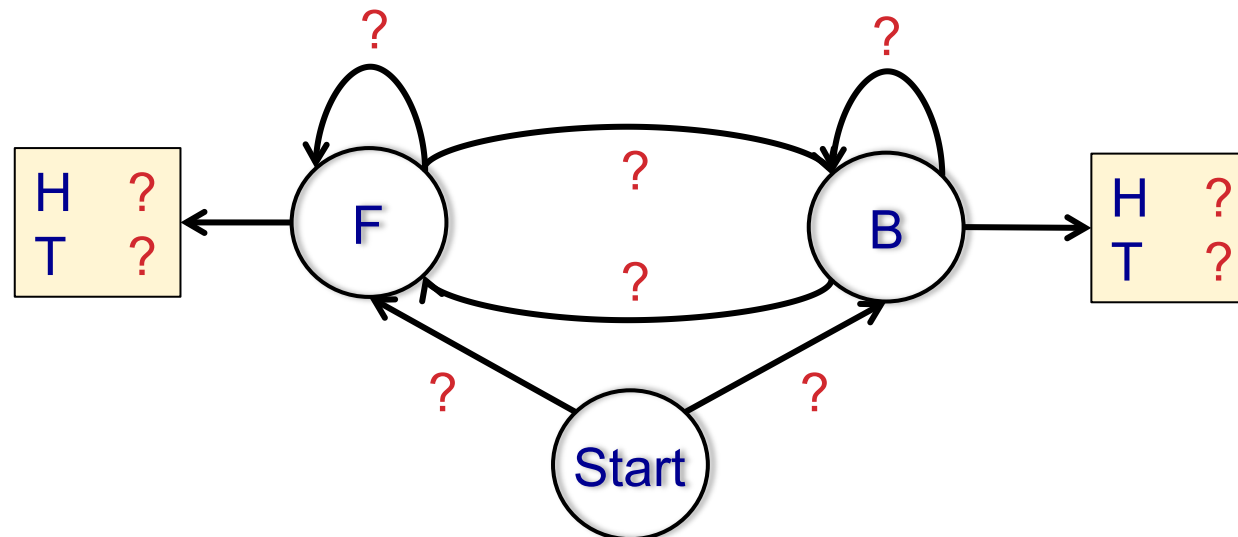
- **Supervised learning:**

- Training data: paired sequences of states and observations $(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_T, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$
- $p(\mathbf{y}_i)$ = num. of sequences starting with \mathbf{y}_i /num. of all sequences.
- $p(\mathbf{y}_j | \mathbf{y}_i)$ = number of $(\mathbf{y}_i, \mathbf{y}_j)$'s/number of all $(\mathbf{y}_i, \mathbf{y})$'s.
- $p(\mathbf{x}_j | \mathbf{y}_i)$ = number of $(\mathbf{y}_i, \mathbf{x}_j)$'s/number of all $(\mathbf{y}_i, \mathbf{x})$'s.

HIDDEN MARKOV MODELS

- Supervised learning example:

| | | | |
|--------|--------|---------|--------|
| FFFBFF | BFFBFF | FFBFFF | FFFFBF |
| HHTHTH | THTHTH | THHTTH | THTTTH |
| BFFFBF | FFFBBF | BFFFFFF | BFBFFF |
| THHTHT | HHTHHT | HHTTHT | HTTTHH |



HIDDEN MARKOV MODELS

- **Unsupervised learning:**
 - Only observation sequences are available.

| | | | |
|-------------------|-------------------|--------------------|-------------------|
| FFFBFF | BFFBFF | FFBFFF | FFFFBF |
| HHTHTH | THTHTH | THHTTH | THTTTH |
| BFFFBF | FFFBBF | BFFFFFF | BFBFFF |
| THHTHT | HHTHHT | HHTTHT | HTTTHH |

HIDDEN MARKOV MODELS

- **Unsupervised learning:**

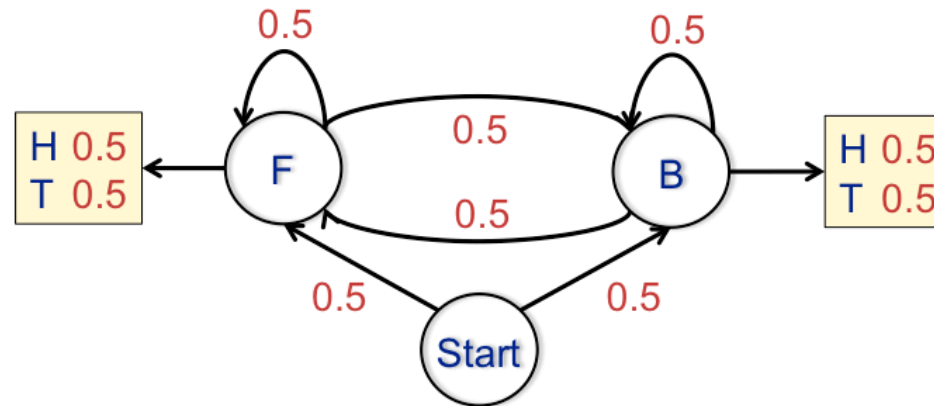
- Only observation sequences are available.

| | | | |
|-------------------|-------------------|--------------------|-------------------|
| FFFBFF | BFFBFF | FFBFFF | FFFFBF |
| HHTHTH | THTHTH | THHTTH | THTTTH |
| BFFFBF | FFFBBF | BFFFFFF | BFBFFF |
| THHTHT | HHTHHT | HHTTHT | HTTTHH |

- Iterative improvement of model parameters.
 - How?

HIDDEN MARKOV MODELS

- **Unsupervised learning:**
 - Initialize estimated parameters.



- For each observation sequence, compute the most probable state sequence, using Viterbi algorithm.
- Update the parameters using supervised learning on obtained paired state-observation sequences.
- Repeat it until convergence.

HIDDEN MARKOV MODELS

- **Application to NER:**

- Example: “Facebook CEO Zuckerberg visited Vietnam”.

ORG NIL PER NIL LOC

HIDDEN MARKOV MODELS

- **Application to NER:**

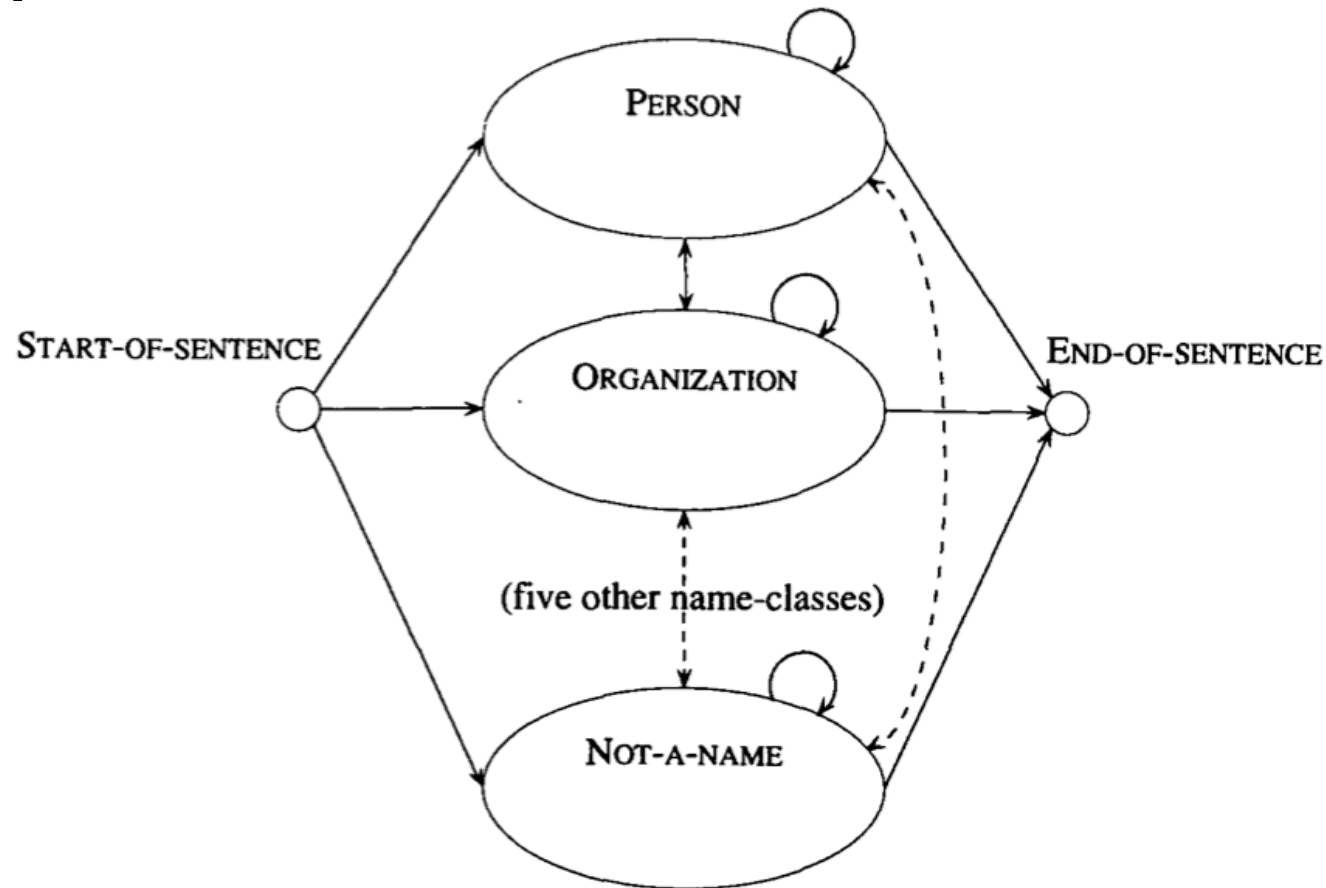
- Example: “Facebook CEO Zuckerberg visited Vietnam”.

 | | | | |
 ORG **NIL** **PER** **NIL** **LOC**

- States = Class labels
 - Observations = Words + Features

HIDDEN MARKOV MODELS

- Application to NER:



HIDDEN MARKOV MODELS

- **Application to NER:**
 - What if a name is a multi-word phrase?
 - Example: "... **John von Neumann** is ..."

HIDDEN MARKOV MODELS

- **Application to NER:**

- What if a name is a multi-word phrase?
- Example: "... **John von Neumann** is ..."


B-PER I-PER I-PER O

- BIO notation:
{B-PER, I-PER, B-ORG, I-ORG, B-LOC, I-LOC, B-MISC, I-MISC, O}

READING HOMEWORK 2

- **Marsland, S. (2009) Machine learning: An algorithmic perspective. Chapter 15 (graphical models).**
- **Bikel, D.M. et al. (1997) Nymble: a high performance learning name-finder.**

EXERCISES 2

- **Write a program to carry out the unsupervised learning example for HMM in the lecture. Discuss on the result, in particular the convergence of the process.**