

# CS540 Introduction to Artificial Intelligence

## Lecture 23

Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

June 13, 2020

# Special Bayesian Network for Sequences

## Motivation

- A sequence of features  $X_1, X_2, \dots$  can be modelled by a Markov Chain but they are not observable.
- A sequence of labels  $Y_1, Y_2, \dots$  depends only on the current hidden features and they are observable.
- This type of Bayesian Network is call a Hidden Markov Model.

# HMM Applications Part 1

## Motivation

- Weather prediction.
- Hidden states:  $X_1, X_2, \dots$  are weather that is not observable by a person staying at home (sunny, cloudy, rainy).
- Observable states:  $Y_1, Y_2, \dots$  are Badger Herald newspaper reports of the condition (dry, dryish, damp, soggy).
- Speech recognition.
  - Hidden states:  $X_1, X_2, \dots$  are words.
  - Observable states:  $Y_1, Y_2, \dots$  are acoustic features.

re k  


---

wreck

a ni s  


---

a nice

b ch  


---

beach

# HMM Applications Part 2

## Motivation

- Stock or bond prediction.
- Hidden states:  $X_1, X_2, \dots$  are information about the company (profitability, risk measures).
- Observable states:  $Y_1, Y_2, \dots$  are  <sup>$Y_{t+1}$</sup>  stock or bond prices.
- Speech synthesis: Chatbox.
- Hidden states:  $X_1, X_2, \dots$  are context or part of speech.
- Observable states:  $Y_1, Y_2, \dots$  are words.

# Other HMM Applications

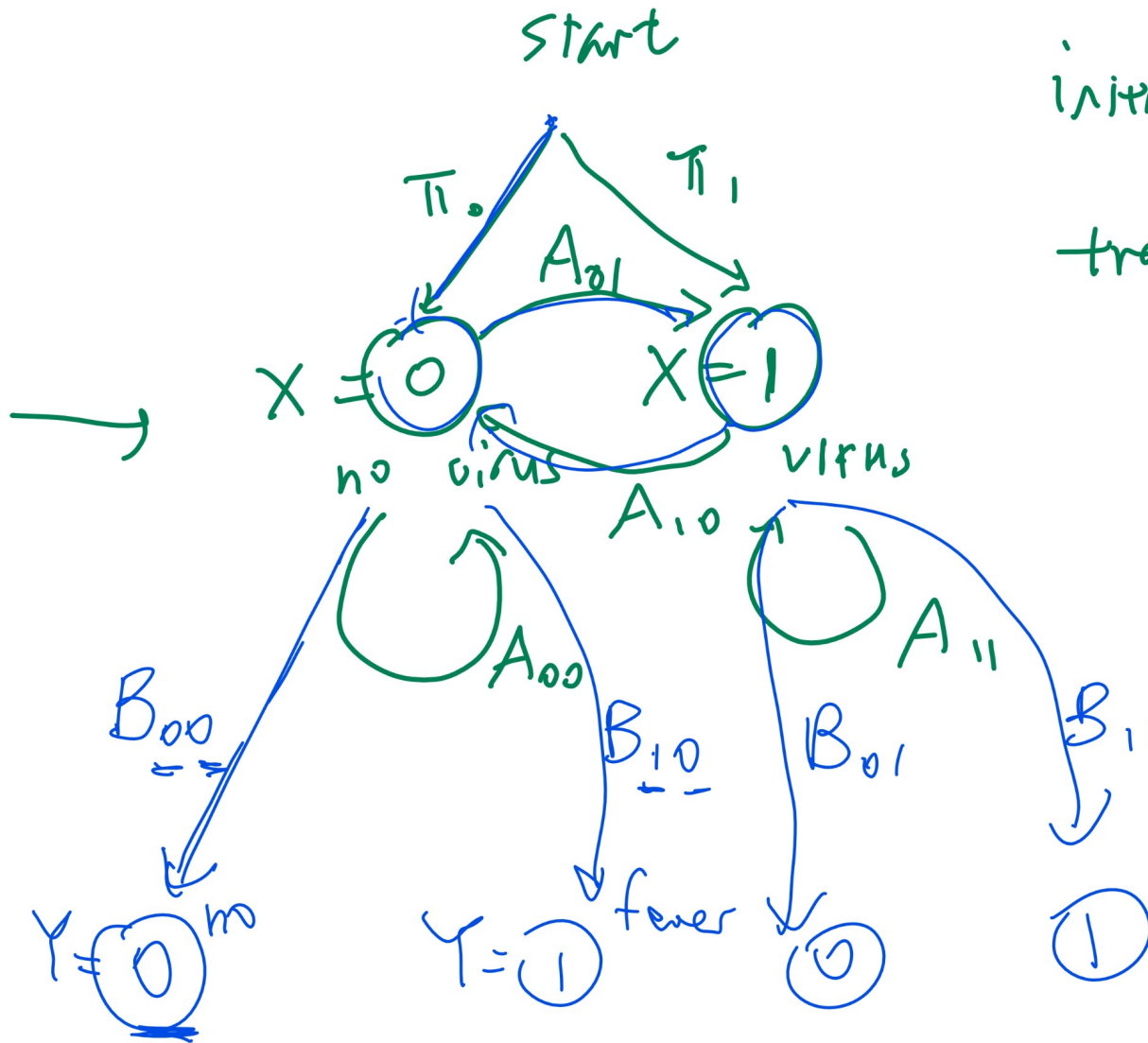
## Motivation

- Machine translation.
- Handwriting recognition.
- Gene prediction.
- Traffic control.



# Hidden Markov Model Diagram

## Motivation



initial prob  $(\pi_0, \pi_1)$

transition matrix  $A = \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix}$

conditional prob

$B_0 = [B_{00}, B_{01}]$

$B_1 = [B_{10}, B_{11}]$

# Transition and Likelihood Matrices

## Motivation

- An initial distribution vector and two state transition matrices are used to represent a hidden Markov model.

- 1 Initial state vector:  $\pi$ .

$$\pi_i = \mathbb{P}\{X_1 = i\}, i \in 1, 2, \dots, |X|$$

- 2 State transition matrix:  $A$ .

$$A_{ij} = \mathbb{P}\{X_t = j | X_{t-1} = i\}, i, j \in 1, 2, \dots, |X|$$

- 3 Observation Likelihood matrix (or output probability distribution):  $B$ .

$$B_{ij} = \mathbb{P}\{Y_t = i | X_t = j\}, i \in 1, 2, \dots, |Y|, j \in 1, 2, \dots, |X|$$

# Markov Property

## Motivation

- The Markov property implies the following conditionally independence property.

$$\mathbb{P}\{x_t | \underbrace{x_{t-1}, x_{t-2}, \dots, x_1}\} = \mathbb{P}\{x_t | x_{t-1}\}$$
$$\mathbb{P}\{\underbrace{y_t | x_t, x_{t-1}, \dots, x_1}\} = \mathbb{P}\{y_t | x_t\}$$



# Evaluation and Training

## Motivation

- There are three main tasks associated with a HMM.

- ① Evaluation problem: finding the probability of an observed sequence given an HMM:  $y_1, y_2, \dots$
- ② Decoding problem: finding the most probable hidden sequence given the observed sequence:  $x_1, x_2, \dots$
- ③ Learning problem: finding the most probable HMM given an observed sequence:  $\pi, A, B, \dots$

Training HMM

EM

CPT

# Expectation Maximization Algorithm

## Description

- Start with a random guess of  $\pi, A, B$ .
  - Compute the forward probabilities: the joint probability of a observed sequence and its hidden state.
  - Compute the backward probabilities: the probability of a observed sequence given its hidden state.
  - Update the model  $\pi, A, B$  using Bayes rule.
  - Repeat until convergence.
  - Sometimes, it is called the Baum-Welch Algorithm.
-

# Evaluation Problem

## Definition

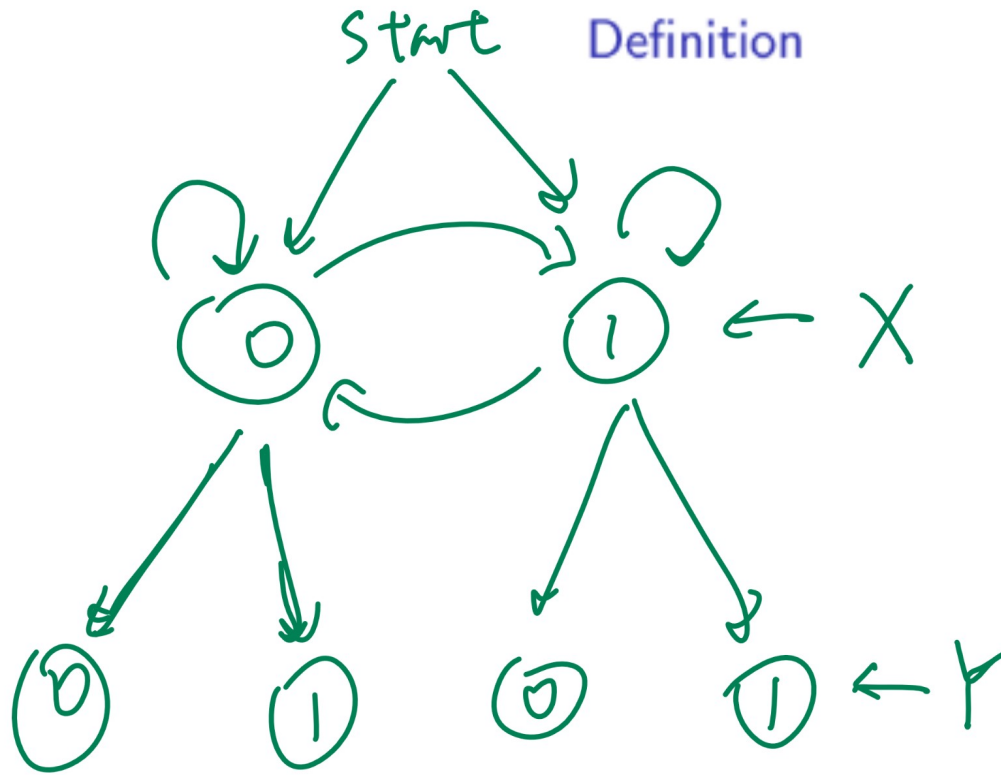
- The task is to find the probability  $\mathbb{P}\{y_1, y_2, \dots, y_T | \pi, A, B\}$ .

$$\begin{aligned}
 & \mathbb{P}\{y_1, y_2, \dots, y_T | \pi, A, B\} \\
 &= \sum_{x_1, x_2, \dots, x_T} \mathbb{P}\{y_1, y_2, \dots, y_T | x_1, x_2, \dots, x_T\} \mathbb{P}\{x_1, x_2, \dots, x_T\} \\
 &= \sum_{x_1, x_2, \dots, x_T} \left( \prod_{t=1}^T B_{y_t x_t} \right) \left( \pi_{x_1} \prod_{t=2}^T A_{x_{t-1} x_t} \right)
 \end{aligned}$$

*Handwritten notes: "0.1 0.1 0.1" above the probability expression, and "0 0 1 1" below it. A green bracket groups the entire equation. Blue underlines are under the product terms in the final equation.*

- This is also called the Forward Algorithm.

# Evaluation Problem Example, Part 1



$$\begin{aligned}
 & P_r \{ Y_1 = \underset{y_1}{0}, Y_2 = \underset{y_2}{1} \mid \pi, A, B \} \\
 &= P_r \{ \underline{Y_1 = 0}, \underline{Y_2 = 1}, \boxed{X_1 = \underset{\cdot}{0}, X_2 = \underset{\cdot}{0}} \mid \pi, A, B \} + \\
 & \quad + \\
 & \quad +
 \end{aligned}$$

# Evaluation Problem Example, Part 2

## Definition

$$\begin{aligned}
 & P_r \{ Y_1 = 0, Y_2 = 1, X_1 = 0, X_2 = 0 \mid \pi, A, B \} \\
 &= \underbrace{P_r \{ Y_1 = 0 \mid X_1 = 0 \}}_{B_{00}} \cdot \underbrace{P_r \{ Y_2 = 1 \mid X_2 = 0 \}}_{B_{10}} \\
 &\quad \underbrace{P_r \{ X_2 = 0 \mid X_1 = 0 \}}_{A_{00}} \cdot \underbrace{P_r \{ X_1 = 0 \}}_{\pi_0}
 \end{aligned}$$

$$\begin{aligned}
 & P_r \{ \underline{Y_1 = a}, \underline{Y_2 = b}, Y_3 = c, \underline{X_1 = d}, \underline{X_2 = e}, X_3 = f \} \\
 &= \underbrace{B_{ad} \cdot B_{be} \cdot B_{cf}}_{\text{emissions}} \cdot \underbrace{\pi_d \cdot A_{de} \cdot A_{ef}}_{\text{transitions}}
 \end{aligned}$$



# Decoding Problem

## Definition

- The task is to find  $x_1, x_2, \dots, x_T$  that maximizes  $\mathbb{P}\{x_1, x_2, \dots, x_T | y_1, y_2, \dots, y_T, \pi, A, B\}$ .
- Direct computation is too expensive.
- Dynamic programming needs to be used to save computation.
- This is called the Viterbi Algorithm.

any all combinations

$$\frac{P_r(y_1, \dots, y_T | x_1, \dots, x_T)}{P_r(y_1, \dots, y_T)}$$

577

# Viterbi Algorithm Value Function

## Definition

- Define the value functions to keep track of the maximum probabilities at each time  $t$  and for each state  $k$ .

$$\begin{aligned}
 V_{1,k} &= \mathbb{P}\{y_1|X_1 = k\} \cdot \mathbb{P}\{X_1 = k\} = \underbrace{\text{Pr}\{y_1, x_1\}} \\
 &= \underbrace{B_{y_1 k} \pi_k} \\
 V_{t,k} &= \max_x \mathbb{P}\{y_t|X_t = k\} \mathbb{P}\{X_t = k|X_{t-1} = x\} V_{1,k} \\
 &= \max_x \underbrace{B_{y_t k} A_{kx}} \underbrace{V_{1,k}} \quad \underbrace{\text{Pr}\{y_t, x_1, \dots, x_t\}}
 \end{aligned}$$



# Viterbi Algorithm Policy Function

## Definition

- Define the policy functions to keep track of the  $x_t$  that maximizes the value function.

$$\text{policy}_{t,k} = \arg \max_x B_{y_t k} A_{kx} V_{1,k}$$

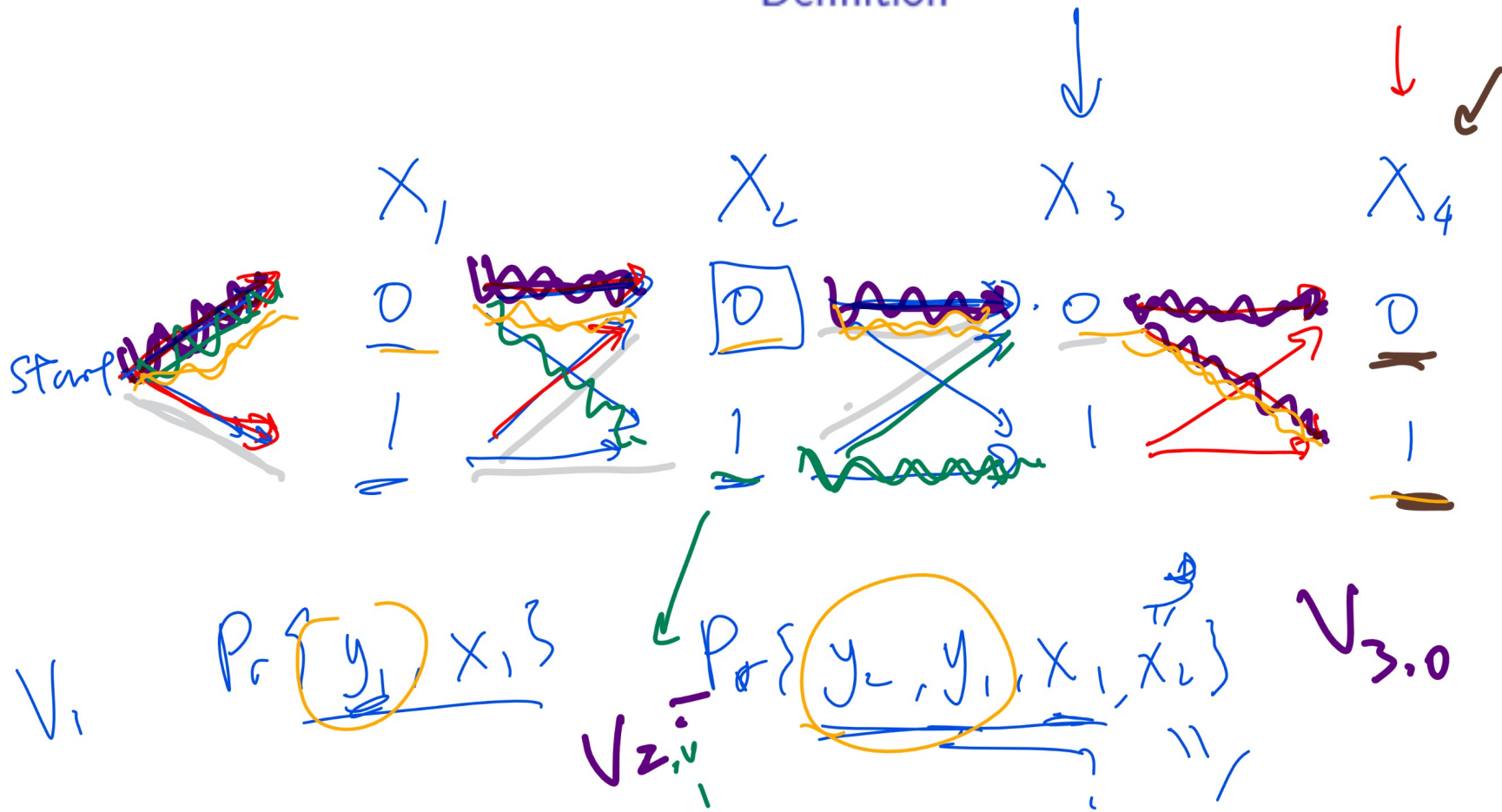
- Given the policy functions, the most probable hidden sequence can be found easily.

$$x_T = \arg \max_x V_{T,x}$$

$$x_t = \text{policy}_{t+1, x_{t+1}}$$

# Dynamic Programming Diagram

Definition



# Viterbi Algorithm Diagram

## Definition

# Expectation Maximization Algorithm (for HMM), Part 1

## Algorithm

- Initialize the hidden Markov model.

CPT

$$\pi \sim D(|X|), A \sim D(|X|, |X|), B \sim D(|Y|, |X|)$$

- Perform the forward pass.

count  $X_t = i, X_{t+1} = j$   
 $X_t = i$

$\alpha_{i,t}$  represents  $\mathbb{P}\{y_1, y_2, \dots, y_t, X_t = i | \pi, A, B\}$

$\alpha_{i,1} = \pi_i B_{y_1,i}$   $\Pr\{y_1, X_1 = i\}$

$\alpha_{i,t+1} = \sum_{j=1}^{|X|} \alpha_{j,t} A_{ji} B_{y_{t+1},i}$   $\Pr\{y_1, \dots, y_t, X_t = j, X_{t+1} = i\}$

$\Pr\{y_1, \dots, y_t, X_1 = i, X_2 = j\}$

The diagram illustrates the forward pass calculation of  $\alpha_{i,t}$ . It shows the recursive formula  $\alpha_{i,t+1} = \sum_{j=1}^{|X|} \alpha_{j,t} A_{ji} B_{y_{t+1},i}$  with handwritten annotations. Arrows indicate the flow of information from the previous time step  $t$  to the current time step  $t+1$ . The terms  $\alpha_{j,t}$ ,  $A_{ji}$ , and  $B_{y_{t+1},i}$  are circled in green. The diagram also shows the initial condition  $\alpha_{i,1} = \pi_i B_{y_1,i}$  and the probability expression  $\Pr\{y_1, X_1 = i\}$ . The sequence of observations  $y_1, \dots, y_t, y_{t+1}, \dots, y_T$  is shown at the bottom, with  $X_t = j$  and  $X_{t+1} = i$  indicated below the corresponding observations.

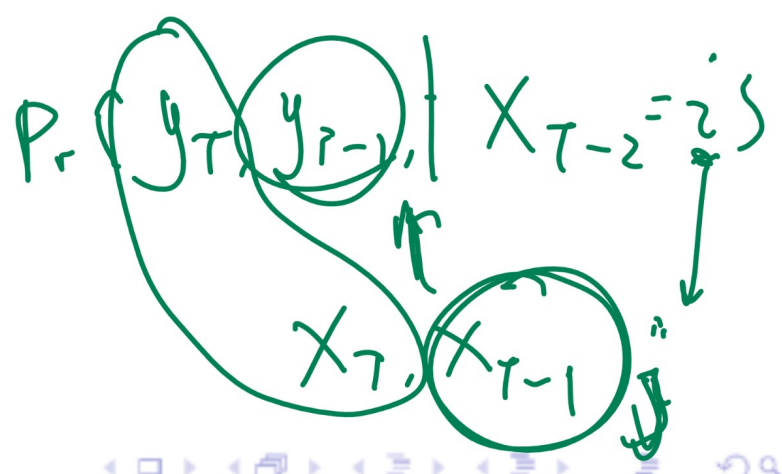
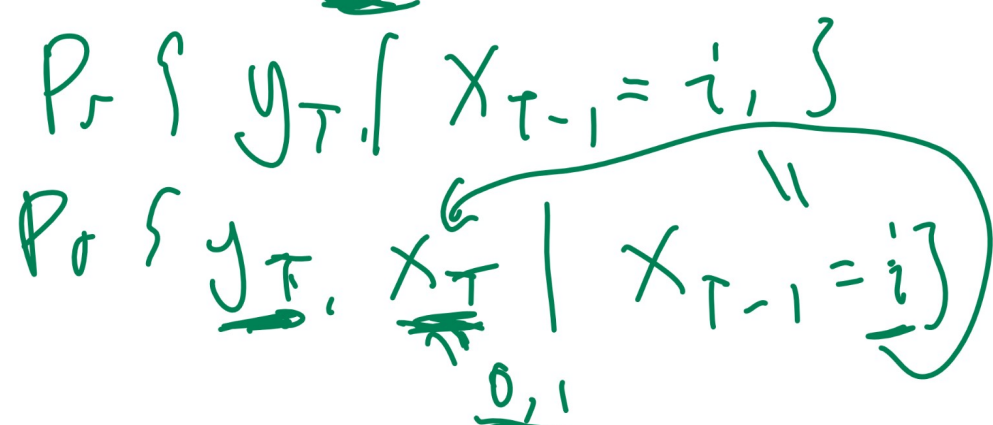
# Expectation Maximization Algorithm (for HMM), Part 2

## Algorithm

- Perform the backward pass.

$\beta_{i,t}$  represents  $\mathbb{P}\{y_{t+1}, y_{t+2}, \dots, y_T | X_t = i, \pi, A, B\}$

$\beta_{i,T} = 1$   $\mathbb{P}\{y_T | X_T = i\}$   
 $\beta_{i,t} = \sum_{j=1}^{|X|} A_{ij} B_{y_{t+1},j} \beta_{j,t+1}$



# Expectation Maximization Algorithm (for HMM), Part 3

## Algorithm

- Define the conditional hidden state probabilities for each training sequence  $n$ .

count  $X_t = i$

$\gamma_{n,i,t}$  = represents  $\mathbb{P}\{X_t = i | y_1, y_2, \dots, y_T, \pi, A, B\}$

$$\gamma_{n,i,t} = \frac{\alpha_{i,t} \beta_{i,t}}{\sum_{j=1}^{|X|} \alpha_{j,t} \beta_{j,t}}$$

$\frac{\mathbb{P}\{X_t = i, y_1, \dots, y_T\}}{\sum_{j=1}^{|X|} \mathbb{P}\{X_t = j, y_1, \dots, y_T\}}$

$\mathbb{P}\{y_1, \dots, y_t, X_t = i\} \mathbb{P}\{y_{t+1}, \dots, y_T | X_t = i\}$

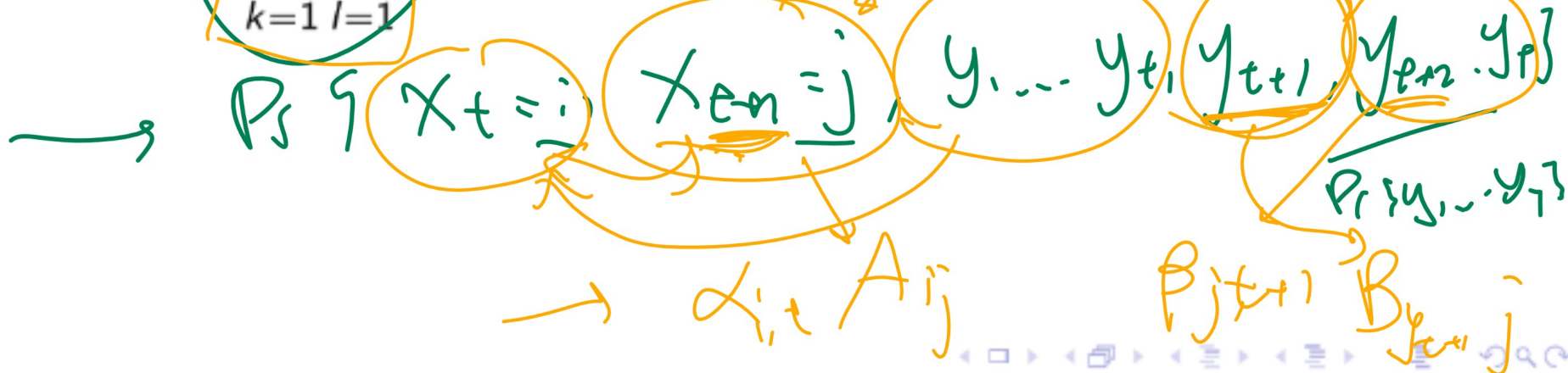
# Expectation Maximization Algorithm (for HMM), Part 4 Algorithm

- Define the conditional hidden state probabilities for each training sequence  $n$ .

count  $X_t = i, X_{t+1} = j$

$\xi_{n,i,j,t}$  represents  $\mathbb{P}\{X_t = i, X_{t+1} = j | y_1, y_2, \dots, y_T, \pi, A, B\}$

$$\xi_{n,i,j,t} = \frac{\alpha_{i,t} A_{ij} \beta_{j,t+1} B_{y_{t+1}j}}{\sum_{k=1}^{|\mathcal{X}|} \sum_{l=1}^{|\mathcal{X}|} \alpha_{k,t} A_{kl} \beta_{l,t+1} B_{y_{t+1}l}}$$



# Expectation Maximization Algorithm (for HMM), Part 5

## Algorithm

- Update the model.

$$\hat{\pi}_j = \frac{\sum_{n=1}^N \gamma_{n,j}(1)}{N}$$

$\# X_1 = i$

---


$$\hat{A}'_{ij} = \frac{\sum_{n=1}^N \sum_{t=1}^{T-1} \xi_{n,i,j,t}}{\sum_{n=1}^N \sum_{t=1}^{T-1} \gamma_{n,i,t}}$$

$\# X_t = i, X_{t+1} = j$

---

$\# X_t = i$



# Expectation Maximization Algorithm (for HMM), Part 6

## Algorithm

- Update the model, continued.

$$B'_{ij} = \frac{\sum_{n=1}^N \sum_{t=1}^T \mathbb{1}_{\{y_{n,t}=j\}} \gamma_{n,i,t}}{\sum_{n=1}^N \sum_{t=1}^T \gamma_{n,i,t}}$$

$(\# X_e=i, Y_e=j)$

$\# X_e=i$

- Repeat until  $\pi, A, B$  converge. ↩