

# CS540 Introduction to Artificial Intelligence

## Lecture 23

Young Wu

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

June 18, 2020

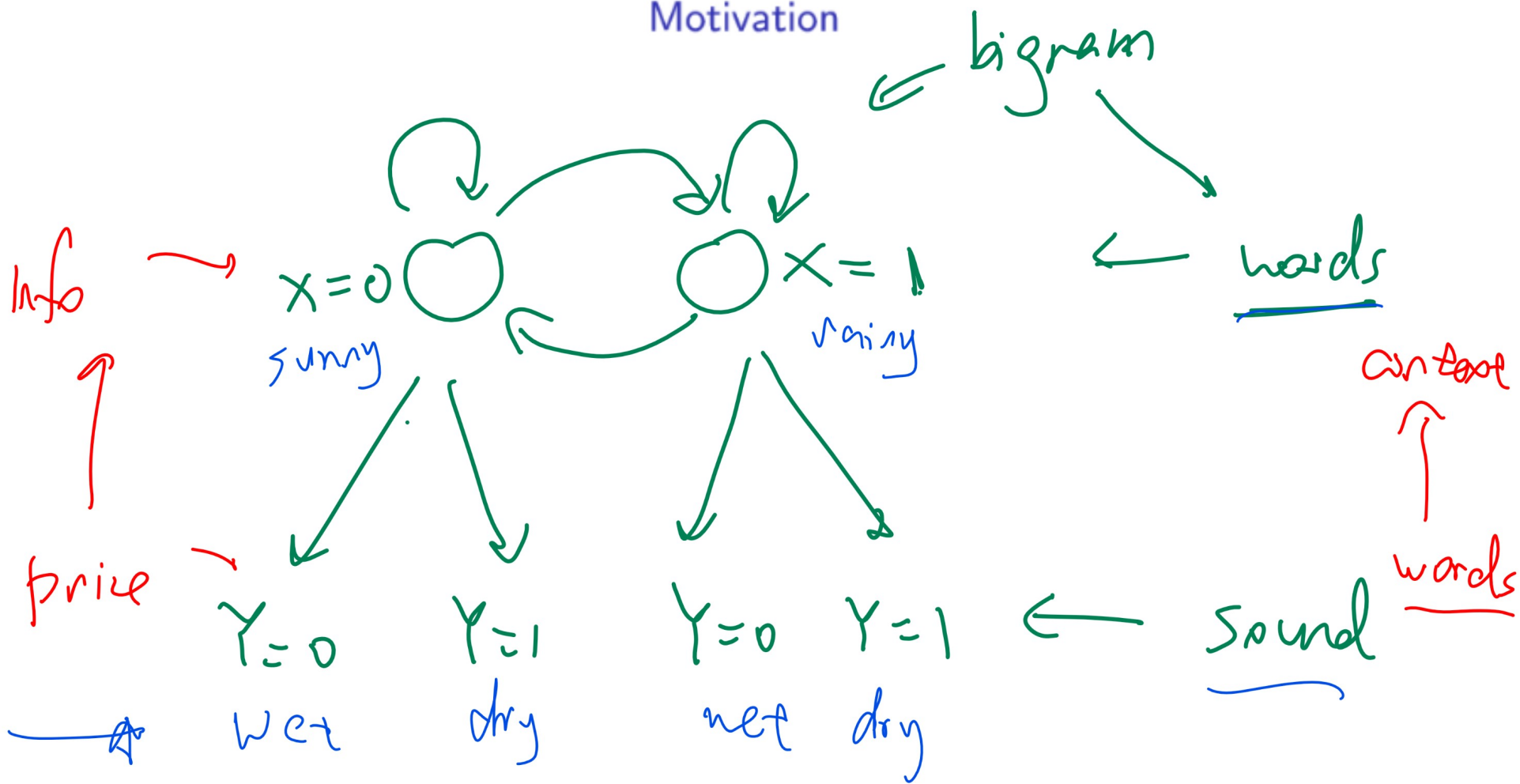
# Special Bayesian Network for Sequences

## Motivation

- A sequence of features  $X_1, X_2, \dots$  can be modeled 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 called a Hidden Markov Model.

# Hidden Markov Model Diagram

Motivation



# Evaluation and Training

## Motivation

- There are three main tasks associated with an 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:  $\pi, A, B, \dots$

# Evaluation Problem

## Definition

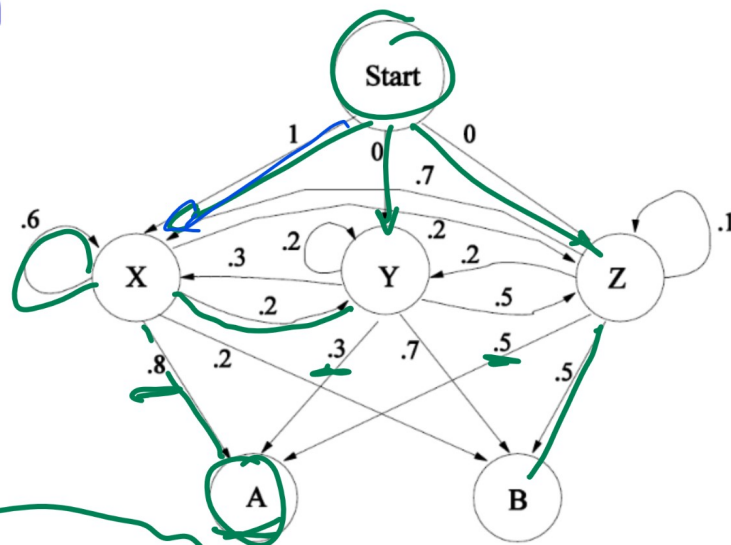
- The task is to find the probability  $\mathbb{P}\{y_1, y_2, \dots, y_T | \pi, A, B\}$ .

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

- This is also called the Forward Algorithm.

# Evaluation Problem Example, Part 1

## Definition



training HMM  
↳ EM algorithm.

2018?

- Fall 2018 Final Q28 and Q29
- Compute  $\mathbb{P}\{X_4 = Y, X_5 = Z | X_3 = X\}$ .
- Compute  $\mathbb{P}\{X_1 = X, X_2 = Z | Y_1 = A, Y_2 = B\}$ .

$$\pi = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \begin{matrix} X_1 = X \\ X_2 = Y \\ X_3 = Z \end{matrix}$$

$$A = \begin{pmatrix} 0.6 & 0.2 & 0.2 \\ 0.3 & 0.2 & 0.5 \\ 0.7 & 0.2 & 0.1 \end{pmatrix}$$

$$B_A = \begin{pmatrix} 0.8 \\ 0.3 \\ 0.5 \end{pmatrix} \leftrightarrow B_B = \begin{pmatrix} 0.2 \\ 0.7 \\ 0.5 \end{pmatrix} \quad B = \begin{pmatrix} 0.8 & 0.2 \\ 0.3 & 0.7 \\ 0.5 & 0.5 \end{pmatrix}$$

## Evaluation Problem Example, Part 2

## Definition

$$\textcircled{1} \quad P_r \{ X_4 = Y, X_5 = Z \mid X_3 = X \}$$

$$= P_r \{ X_4 = Y \mid X_3 = X \} \cdot P_r \{ X_5 = Z \mid X_4 = Y \}$$

$$= 0.2 \cdot 0.5$$

$$\textcircled{2} \quad P_r \{ X_1 = X, X_2 = Z \mid Y_1 = A, Y_2 = B \}$$

$$= P_r \{ Y_1 = A \mid X_1 = X \} \cdot P_r \{ Y_2 = B \mid X_2 = Z \}$$

$$P_r \{ X_1 = X \} \cdot P_r \{ X_2 = Z \mid X_1 = X \}$$

✓

$$P_r \{ Y_1 = A, Y_2 = B \}$$

# Evaluation Problem Example, Part 3



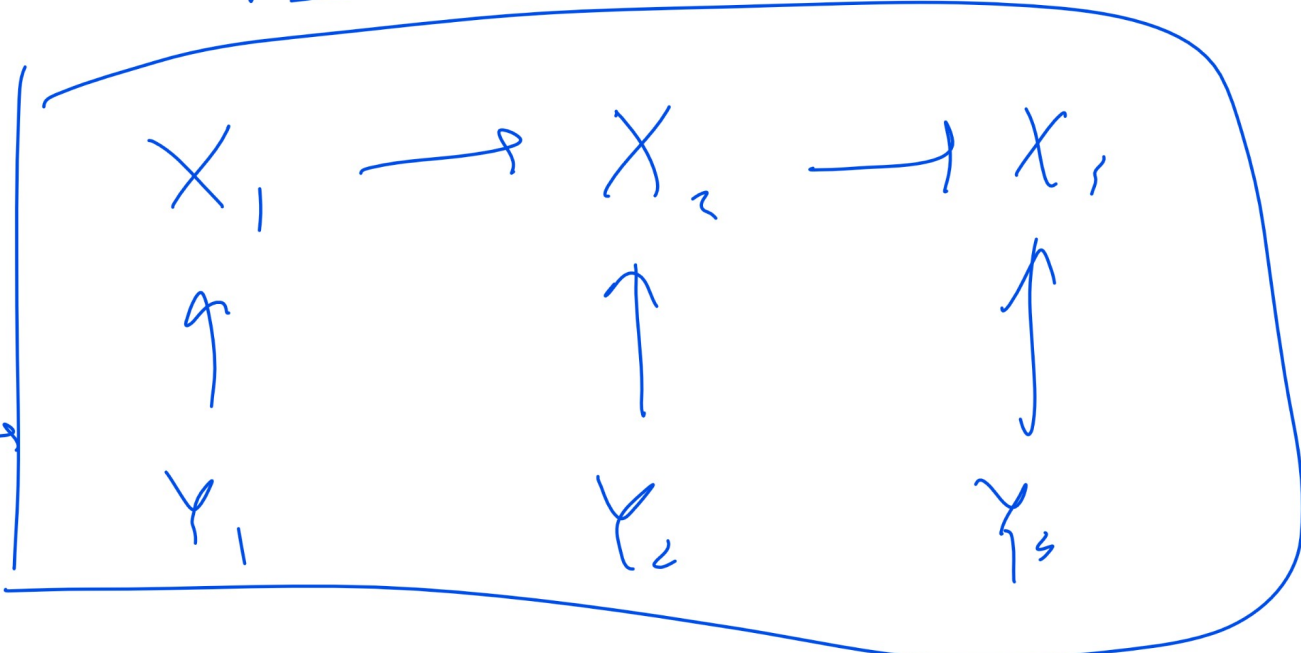
Definition

$$0.8 \cdot 0.5 \cdot 1 - 0.2$$



$$P_0 \left( \begin{matrix} X_1 = X \\ X_2 = Y \\ X_3 = Z \end{matrix} \right), Y_1 = A, Y_2 = B$$

HMM  
Bayesian





# Evaluation Problem Example, Part 4

Definition

$$P_r \{ a, b, c \}$$

$$P_r \{ D \mid a, b, c \}$$

$$= P_r \{ a \} \cdot P_r \{ b \} \cdot P_r \{ c \}$$

unigram count.

$$P_r \{ a, b, c \mid D \} \leftarrow P_r \{ D \}$$

$$P_r \{ a \mid D \} \cdot P_r \{ b \mid D \} \cdot P_r \{ c \mid D \}$$

unigram

# Decoding Problem

## Definition

$$P_r(D|a) = \frac{P_r(a|D) \cdot P_r(D)}{P_r(a|D) P_r(D) + P_r(a|\neg D) \cdot P_r(\neg D)}$$

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

# Viterbi Algorithm Value Function

## Definition

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

$$V_{1,k} = \mathbb{P}\{y_1|X_1 = k\} \cdot \mathbb{P}\{X_1 = k\}$$

$$= B_{y_1 k} \pi_k$$

$$V_{t,k} = \max_x \mathbb{P}\{y_t|X_t = k\} \mathbb{P}\{X_t = k|X_{t-1} = x\} V_{1,k}$$

$$= \max_x B_{y_t k} A_{kx} V_{1,k}$$

# 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.

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

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



# Viterbi Algorithm Diagram

## Definition

# Expectation-Maximization Algorithm (for HMM), Part 1

$\underline{\text{EM}}$  Algorithm

- Initialize the hidden Markov model.

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

- Perform the forward pass.

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

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

$$\alpha_{i,t+1} = \sum_{j=1}^{|X|} \alpha_{j,t} A_{ji} B_{y_{t+1},i}$$

# Expectation-Maximization Algorithm (for HMM), Part 2

## Algorithm

- Perform the backward pass.

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

$$\beta_{i,T} = 1$$

$$\beta_{i,t} = \sum_{j=1}^{|\mathcal{X}|} A_{ij} B_{y_{t+1}j} \beta_{j,t+1}$$



# Expectation-Maximization Algorithm (for HMM), Part 3

## Algorithm

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

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

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

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

## Algorithm

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

$\xi_{n,i,j,t}$  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.

$$\pi'_i = \frac{\sum_{n=1}^N \gamma_{n,i,1}}{N} \quad \frac{C_{x_0=i}}{N}$$
$$A'_{ij} = \frac{\sum_{n=1}^N \sum_{t=1}^{T-1} \xi_{n,i,j,t}}{\sum_{n=1}^N \sum_{t=1}^{T-1} \gamma_{n,i,t}} \quad \frac{C_j}{C_i}$$

# Expectation-Maximization Algorithm (for HMM), Part 6

## Algorithm

- Update the model, continued.

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

$$\frac{C_{Y=j \mid X=i}}{C_{X=i}}$$

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