

CS540 Introduction to Artificial Intelligence

Lecture 12

Young Wu

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

July 8, 2021

Homework

Admin

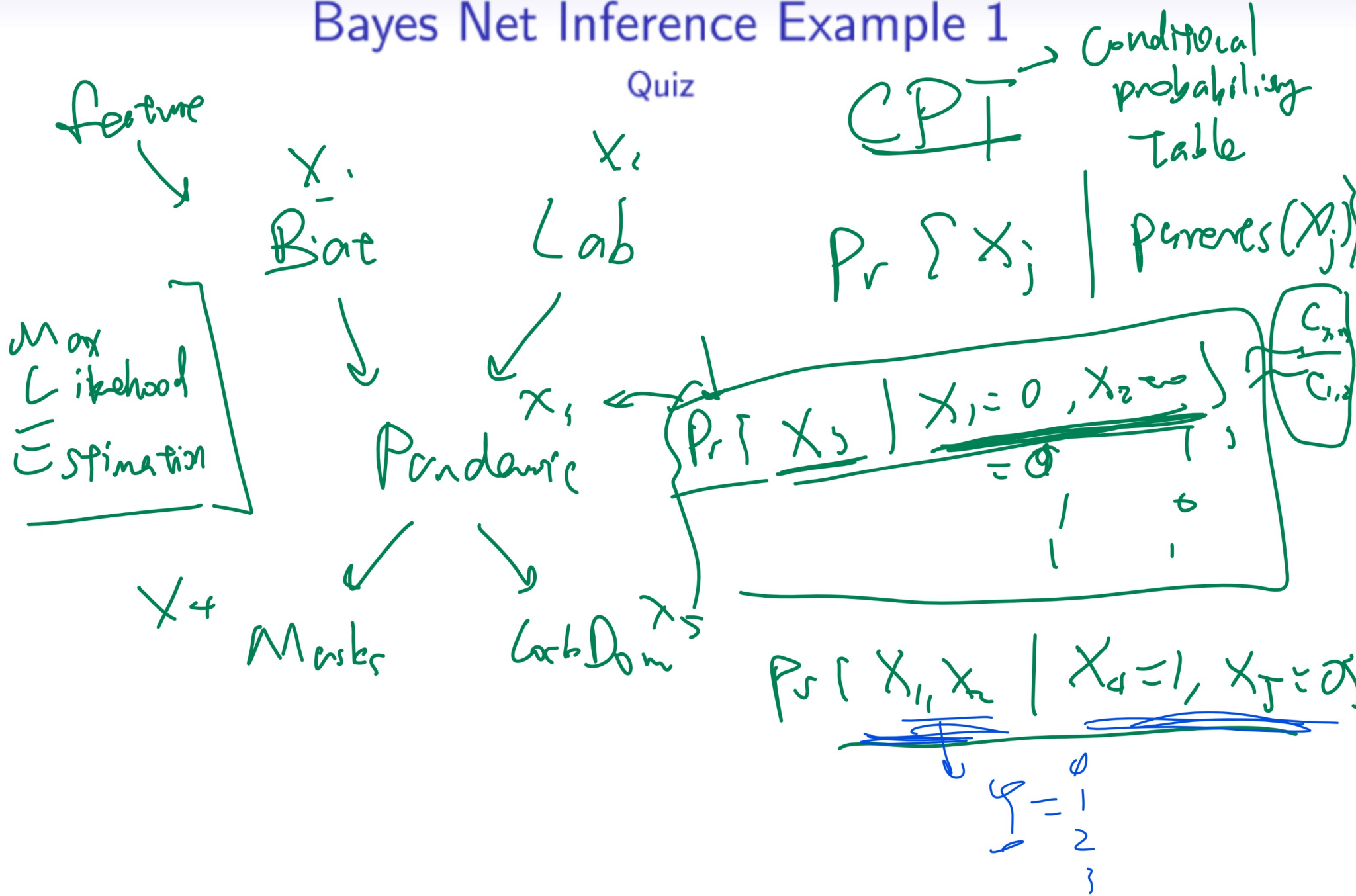
- Please do not submit M8 to M11.
- If you haven't started P1: you should.
- P2 solutions are posted.
- If you use another student's code (or find code on the Internet), you must give attribution at the beginning of your code.
- You are not allowed use another student's output.

Remind Me to Start Recording

Admin

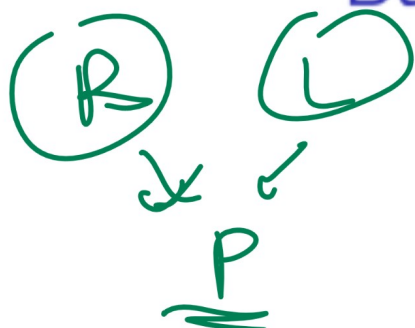
- The messages you send in chat will be recorded: you can change your Zoom name now before I start recording.

Bayes Net Inference Example 1



Bayes Net Inference Example 2

Quiz



- Compute $\hat{\mathbb{P}}\{B = 1 | L = 1\}$?

$$\hat{\mathbb{P}}\{B = 1\} = 0.001, \hat{\mathbb{P}}\{L = 1\} = 0.001$$

$$\hat{\mathbb{P}}\{P = 1 | B = 1, L = 1\} = 0.95, \hat{\mathbb{P}}\{P = 1 | B = 1, L = 0\} = 0.29$$

$$\hat{\mathbb{P}}\{P = 1 | B = 0, L = 1\} = 0.94, \hat{\mathbb{P}}\{P = 1 | B = 0, L = 0\} = 0.00$$

- A: 0, B: 0.001, C: 0.0094, D: 0.0095, E: 1

Handwritten notes in green and red ink:

- $P_r\{B=1\}$ and $P_r\{L=1\}$ (crossed out)
- $P_r\{B=1, L=1\}$
- $P_r\{L=1\}$ (crossed out)
- $P_r\{B=1\}$ (crossed out)

Bayes Net Inference Example 3



Quiz

• Compute $\hat{P}\{B = 1, L = 1 | P = 1\}$?

$P, \{B=1, L=1, P=1\}$

$P, \{P=1\}$

CPT

$\hat{P}\{B = 1\} = 0.001, \hat{P}\{L = 1\} = 0.001$

$\hat{P}\{P = 1 | B = 1, L = 1\} = 0.95, \hat{P}\{P = 1 | B = 1, L = 0\} = 0.29$

$\hat{P}\{P = 1 | B = 0, L = 1\} = 0.94, \hat{P}\{P = 1 | B = 0, L = 0\} = 0.00$

$P, \{B=1\} \cdot P, \{L=1\} \cdot P, \{P=1 | B=1, L=1\}$

• A: $\frac{0.001 \cdot 0.001}{0.001}$, B: $\frac{0.001 \cdot 0.001 \cdot 0.95}{0.001}$

• C: $\frac{0.001 \cdot 0.95 + 0.999 \cdot (0.94 + 0.29)}{0.001 \cdot 0.001}$

• D: $\frac{0.001 \cdot 0.001}{0.001 \cdot 0.95 + 0.999 \cdot (0.94 + 0.29)}$

• E: $\frac{0.001 \cdot 0.95 \cdot 0.001}{0.001 \cdot 0.95 + 0.999 \cdot (0.94 + 0.29)}$

$0.001 \cdot 0.001 \cdot 0.95$

$$P_r\{P=1|B=0,L=0\} \cdot P_r\{B=0\} \cdot P_r\{L=0\}$$

Bayes Net Inference Example 3 Computation

$$P_r\{P=1\}$$

Quiz

$P=1$	$B=0$	$L=0$	↑
$P=1$	$B=0$	$L=1$	↑
$P=1$	$B=1$	$L=0$	↑
$P=1$	$B=1$	$L=1$	↑

$$0 \cdot 0.999 \cdot 0.999$$

$$0.94 \cdot 0.999 \cdot 0.001$$

$$0.29 \cdot 0.001 \cdot 0.999$$

$$0.95 \cdot 0.001 \cdot 0.001$$

- Compute $\hat{P}\{B=1, L=1|P=1\}$?

$$\hat{P}\{B=1\} = 0.001, \hat{P}\{L=1\} = 0.001$$

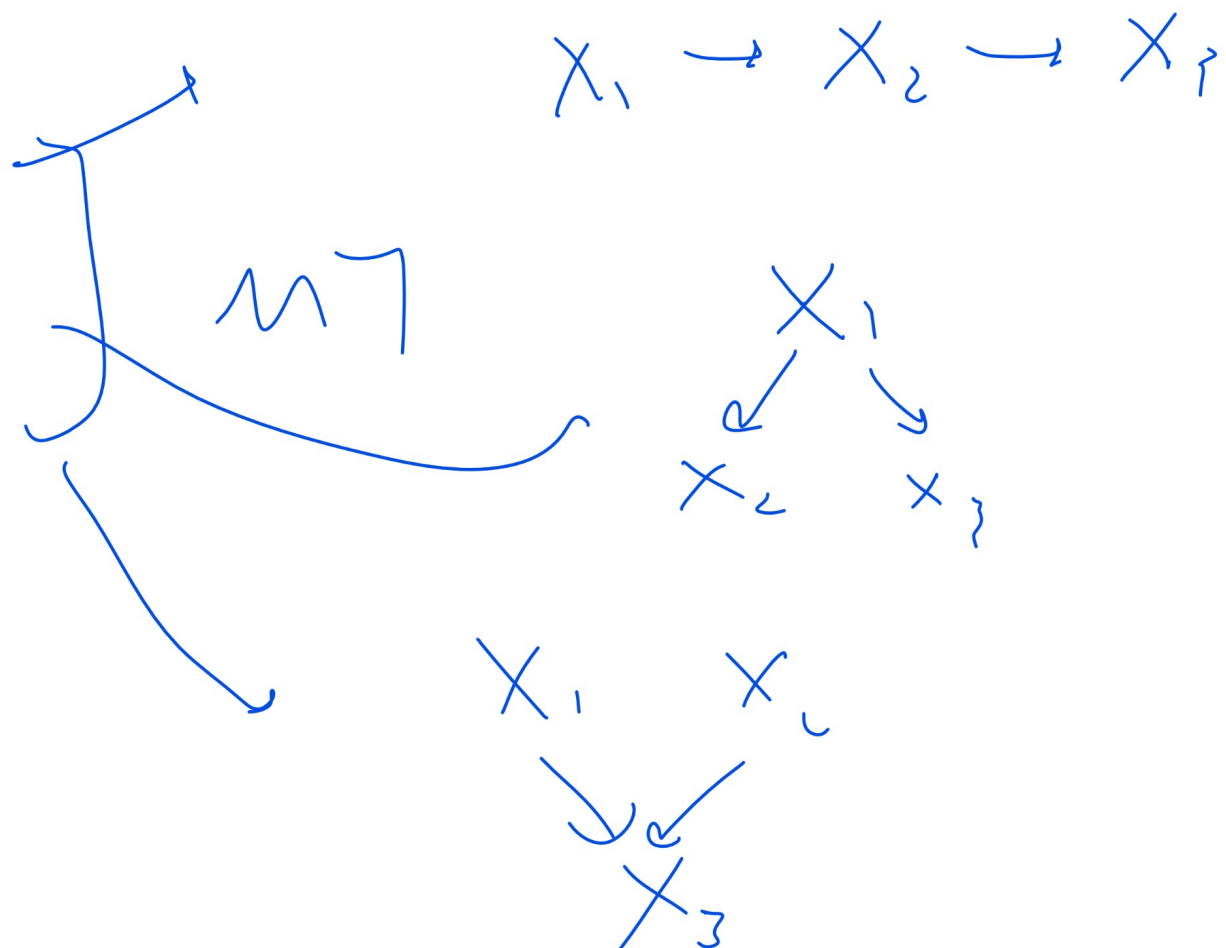
$$\hat{P}\{P=1|B=1, L=1\} = 0.95, \hat{P}\{P=1|B=1, L=0\} = 0.29$$

$$\hat{P}\{P=1|B=0, L=1\} = 0.94, \hat{P}\{P=1|B=0, L=0\} = 0.00$$

Types of Bayes Net Components

Discussion

- Causal Chain
- Common Cause
- Common Effect



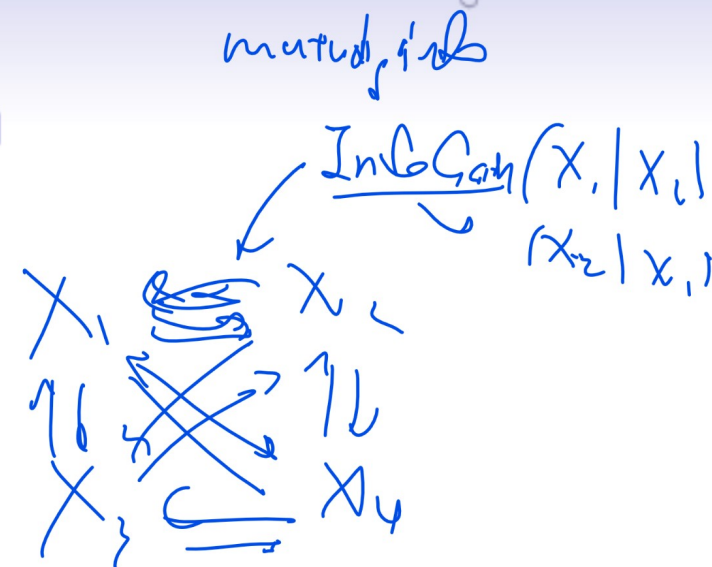
Network Structure

Discussion

- Selecting from all possible structures (DAGs) is too difficult.
- Usually, a Bayesian network is learned with a tree structure.
- Choose the tree that maximizes the likelihood of the training data.

Chow Liu Algorithm

Discussion



- Add an edge between features X_j and $X_{j'}$ with edge weight equal to the information gain of X_j given $X_{j'}$ for all pairs j, j' .
- Find the maximum spanning tree given these edges. The spanning tree is used as the structure of the Bayesian network.

Classification Problem

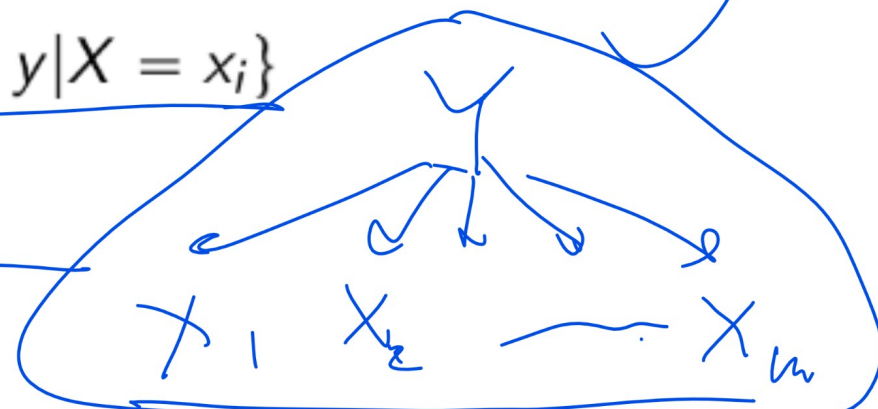
Discussion

- Bayesian networks do not have a clear separation of the label Y and the features X_1, X_2, \dots, X_m .
- The Bayesian network with a tree structure and Y as the root and X_1, X_2, \dots, X_m as the leaves is called the Naive Bayes classifier.
- Bayes rules is used to compute $\mathbb{P}\{Y = y|X = x\}$, and the prediction \hat{y} is y that maximizes the conditional probability.

$$\hat{y}_i = \arg \max_y \mathbb{P}\{Y = y|X = x_i\}$$

CPT } $P_r(X_1 | Y)$
 $P_r(X_2 | Y)$

from
frequency



Naive Bayes Diagram

Discussion

$$P_r \{ Y | X_1, \dots, X_n \} = \frac{P_r \{ Y, X_1, \dots, X_n \}}{P_r \{ X_1, \dots, X_n \}}$$

$$= \frac{P_r \{ Y \} \cdot P_r \{ X_1 | Y \} \cdot P_r \{ X_2 | Y \}}{\sum_y P_r \{ X_1, \dots, X_n, Y=y \}}$$

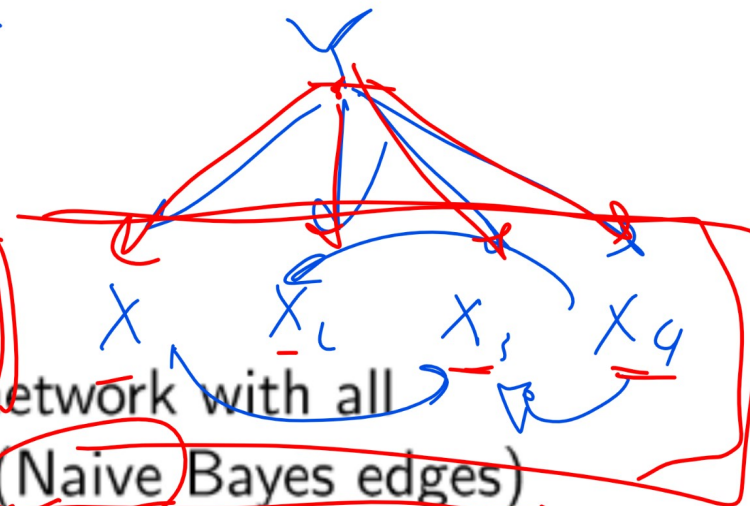
max $P_r \{ Y=y | X_1, \dots, X_n \}$

\rightarrow 0, 1, 2, \dots, K

Tree Augmented Network Algorithm

Discussion

Max Spanning Tree



- It is also possible to create a Bayesian network with all features X_1, X_2, \dots, X_m connected to Y (Naive Bayes edges) and the features themselves form a network, usually a tree (MST edges).
- Information gain is replaced by conditional information gain (conditional on Y) when finding the maximum spanning tree.
- This algorithm is called TAN: Tree Augmented Network.

Tree Augmented Network Algorithm Diagram

Discussion

Special Bayesian Network for Sequences



- 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.

→ recognize speech
→ want a nice beach
 Y_1 Y_2

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: π, A, B, \dots

CPT

$$P_r \{x_1\}$$

$$P_c \{x_2 | x_1\}$$

$$P_r \{y_1 | x_1\}$$

Expectation-Maximization Algorithm

Description

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

Hidden Markov Model Example 1

Definition

- Compute $\mathbb{P}\{X_4 = 1, X_5 = 2 | X_3 = 0\}$.

→ on Friday

Hidden Markov Model Example 1 Computations

Definition

Hidden Markov Model Example 2

Definition

- Compute $\mathbb{P}\{Y_1 = 0, Y_2 = 1\}$.

Hidden Markov Model Example 3

Definition

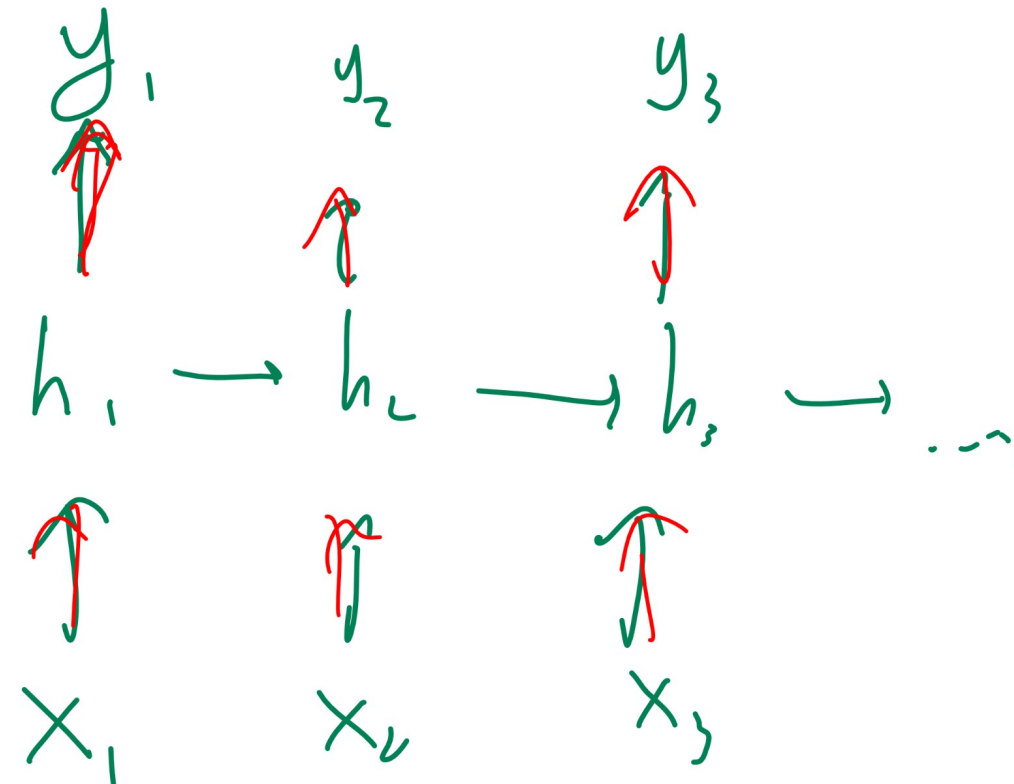
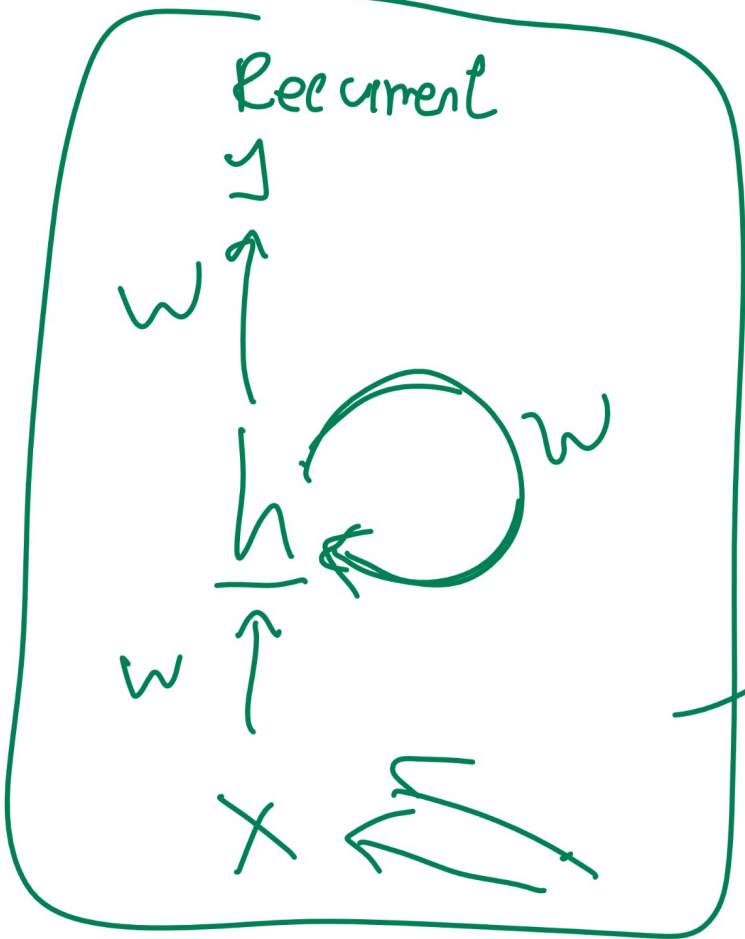
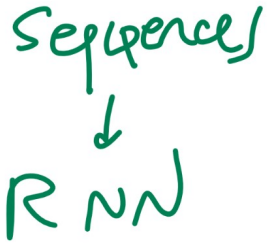
- Compute $\mathbb{P}\{X_1 = 0, X_2 = 2 | Y_1 = 0, Y_2 = 1\}$.

Hidden Markov Model Example 3 Computations

Definition

Dynamic System Diagram

Motivation



Activation Functions

Definition

- The hidden layer activation function can be the tanh activation, and the output layer activation function can be the softmax function.

$$\begin{aligned}
 z_t^{(x)} &= W^{(x)} x_t + W^{(h)} a_{t-1}^{(x)} + b^{(x)} \\
 a_t^{(x)} &= g(z_t^{(x)}), g(\square) = \tanh(\square) \\
 z_t^{(y)} &= W^{(y)} a_t^{(x)} + b^{(y)} \\
 a_t^{(y)} &= g(z_t^{(y)}), g(\square) = \text{softmax}(\square)
 \end{aligned}$$

Handwritten annotations: A large green bracket on the left groups the equations. A green circle highlights the hidden layer activation function $a_t^{(x)}$. A vertical green arrow on the right points from $a_t^{(x)}$ down to $a_t^{(y)}$, with a circled 'G' and an 'X' next to it.

Cost Functions

Definition

- Cross entropy loss is used with softmax activation as usual.

$$\left[\begin{array}{l} C_t = H(y_t, a_t^{(y)}) \\ C = \sum_t C_t \end{array} \right]$$

$$\frac{\partial C}{\partial w^{(x)}}$$

BackPropogation Through Time

Definition

- The gradient descent algorithm for recurrent neural networks is called BackPropogation Through Time (BPTT). The update procedure is the same as standard neural networks using the chain rule.

$$w = w - \alpha \frac{\partial \mathcal{C}}{\partial w}$$
$$b = b - \alpha \frac{\partial \mathcal{C}}{\partial b}$$

Vanishing and Exploding Gradient

Discussion

- If the weights are small, the gradient through many layers will shrink exponentially. This is called the vanishing gradient problem.
- If the weights are large, the gradient through many layers will grow exponentially. This is called the exploding gradient problem.
- Fully connected and convolutional neural networks only have a few hidden layers, so vanishing and exploding gradient is not a problem in training those networks.
- In a recurrent neural network, if the sequences are long, the gradients can easily vanish or explode.

→ LSTM GRU ← Transformer