

# Practice Problems

**ECE/CS 498 DS U/G**

## **Lecture 28: Practice Problems**

Ravi K. Iyer

Dept. of Electrical and Computer Engineering

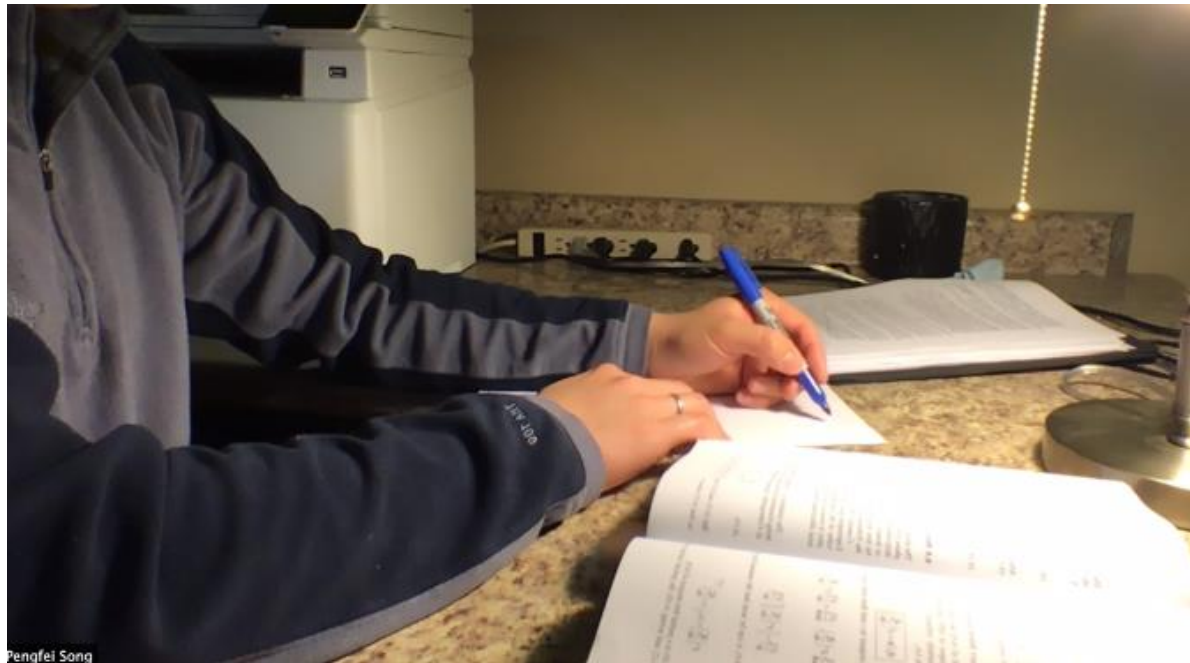
University of Illinois at Urbana Champaign

# Announcements

- **HW 5 due tonight by 11:59 PM on Compass**
- **Final Project:**
  - **Final Presentations will be this Saturday 5/9 from 2-5 PM via Zoom**
    - All graduate students **must attend** the first three presentations
    - Afterwards, there will be two streams
    - Undergraduate students are strongly encouraged to attend
    - We will release presentation signups tonight
  - **Final Report is due Tuesday 5/12 at 11:59 PM via Compass**
    - Up to 8-page IEEE conference style report
    - Must be done in LaTeX
    - We will release report template tonight
- **TA Office Hours:**
  - This week's TA office hours will occur as normal (MW 2-3 PM)
  - Next week, there will be a special extended TA office hours on **Thursday 5/14 from 4-6 PM**

# Final Exam Update

- Due to issues with Proctorio integration, we will now be using Zoom for the final exam
- We will share the Zoom link soon
- You should sign into Zoom at least 5 minutes (that is, 7:55 AM) before the start of the exam (8:00 AM)
- **Make sure your webcam is on.** You need to show your hands and your workspace throughout the exam. For example,



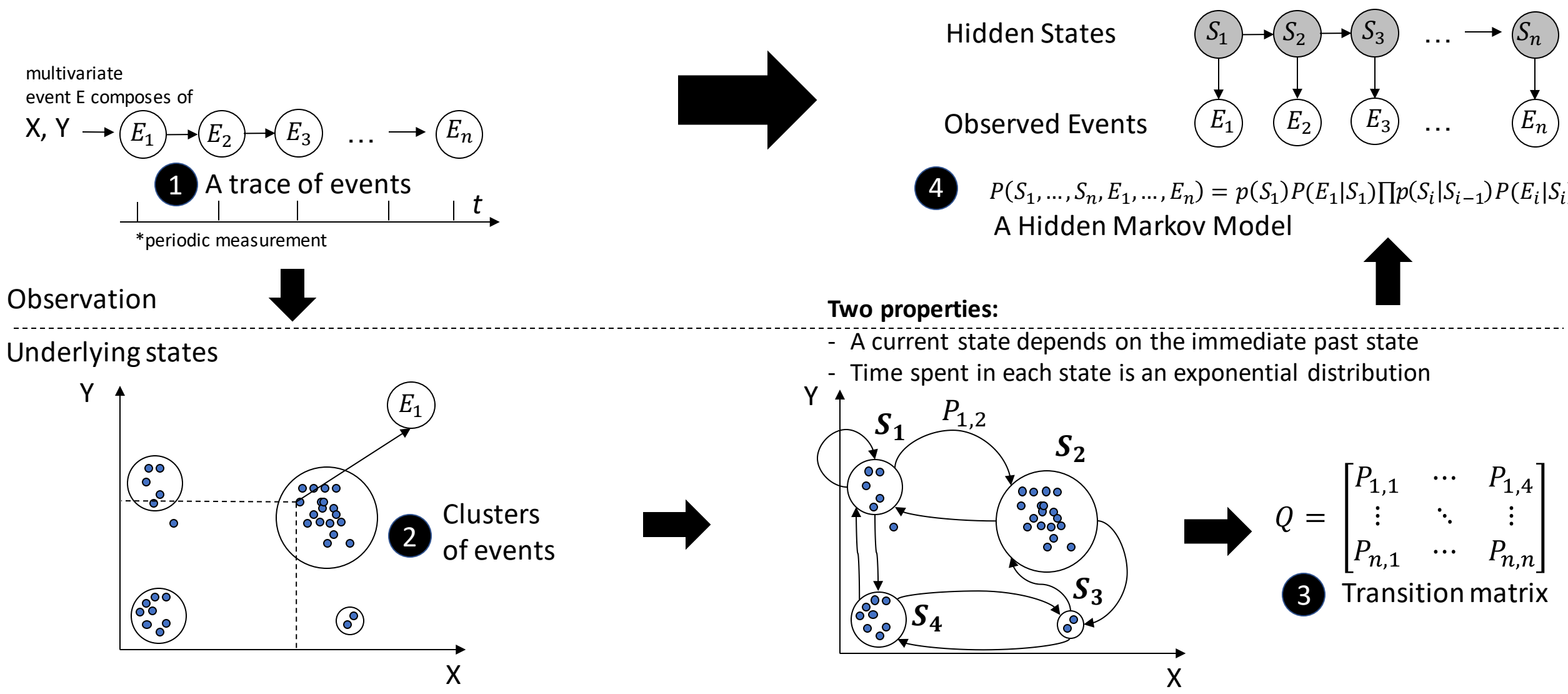
# Final Exam Update

- **Make sure your microphone is on**
- Exam PDF will be released on Compass after you sign an honor agreement
- Exam schedule on **Friday 5/15**
  - 8:00 - 8:10 AM : Read through exam (no writing)
  - 8:10 – 10:40 AM : Complete exam and write out solutions. You will have 2.5 hours
  - 10:40 – 11:10 AM : Scan and submit your solution to Compass (no writing)
    - Late submissions will not be accepted
- If you have questions, raise your hand in Zoom and the instructors will privately message you
  - You will not be allowed to message other students

# Final Exam Update

- **You will be allowed three 8.5x11 cheat sheets on your table, but nothing else**
  - It will be a **closed-book exam** – i.e. no external references (books, Google, electronics, other people, etc.) are allowed
  - During the exam, you should not use your computer except for (i) joining Zoom, (ii) looking through exam PDF, and (iii) scanning/submitting your solutions
  - You are only allowed to use your phone/tablet to scan your exam after you have completed it
- Here is a non-comprehensive list of topics for the exam:
  - **All course material is fair game**
  - **Material from after midterm exam will be emphasized (~75%)**
    - HMMs and Forward-Backward Algorithm
    - Factor Graphs and Belief Propagation
    - SVM
    - Neural Networks
    - Random Forest
    - Cross Validation, etc.
  - **Earlier course material will also be included (~25%)**
    - Basic Probability
    - Naïve Bayes and Bayesian Networks
    - Clustering (K-means, GMM, Hierarchical)
    - Linear/Logistic Regression
    - PCA, etc.

# From a trace of events to a Hidden Markov Model



# Hidden Markov Models

## Model assumptions

An observation depends on its hidden state

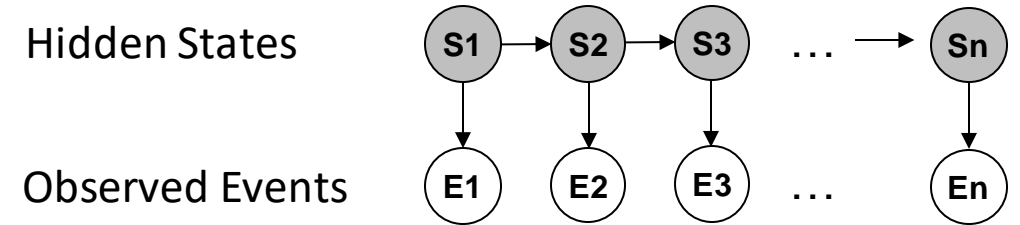
A state variable only depends on the immediate previous state (Markov assumption)

The future observations and the past observations are conditionally independent given the current hidden state

## Advantages:

HMM can model sequential nature of input data (future depends on the past)

HMM has a linear-chain structure that clearly separates system state and observed events.



$$P(S_1, \dots, S_n, E_1, \dots, E_n) = p(S_1)P(E_1|S_1)\prod p(S_i|S_{i-1})P(E_i|S_i)$$

**A Hidden Markov model on observed events and system states**

# Markov Model

- Consider a system which can occupy one of  $N$  discrete *states* or *categories*

$$x_t \in \{1, 2, \dots, N\} \longrightarrow \text{state at time } t$$

- We are interested in *stochastic* systems, in which state evolution is random
- Any *joint* distribution can be factored into a series of *conditional* distributions:

$$p(x_0, x_1, \dots, x_T) = p(x_0) \prod_{t=1}^T p(x_t \mid x_0, \dots, x_{t-1})$$

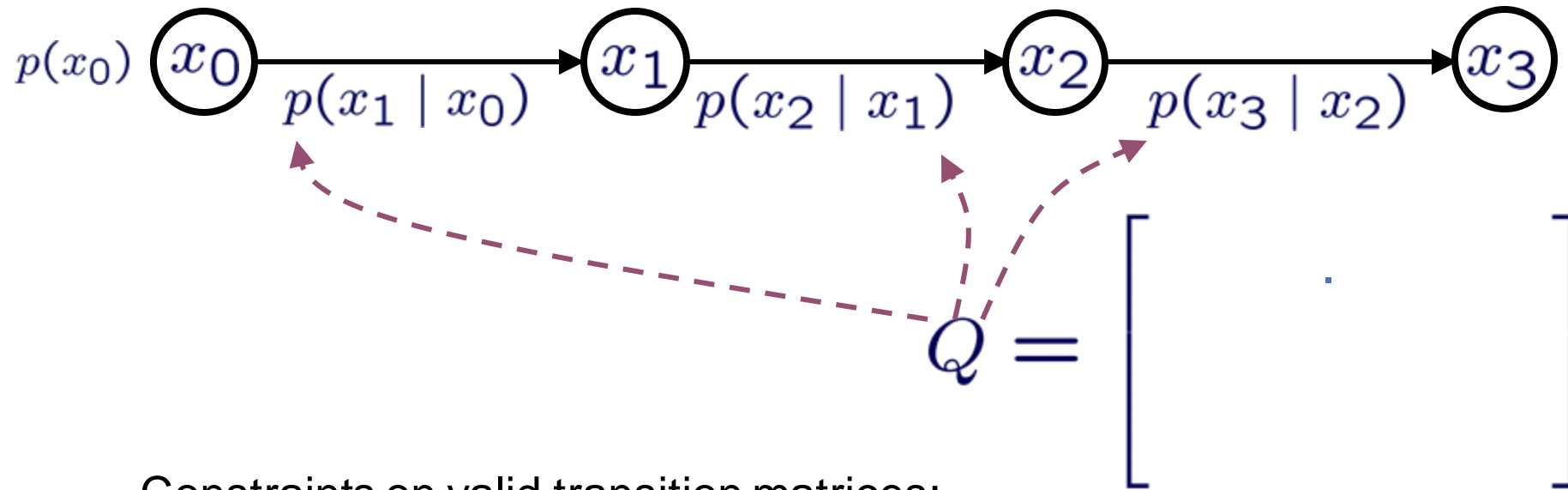
- For a *Markov* process, the next state depends only on the current state:

$$p(x_{t+1} \mid x_0, \dots, x_t) = p(x_{t+1} \mid x_t)$$



# Markov Chains: Graphical Models

$$p(x_0, x_1, \dots, x_T) = p(x_0) \prod_{t=1}^T p(x_t | x_{t-1})$$



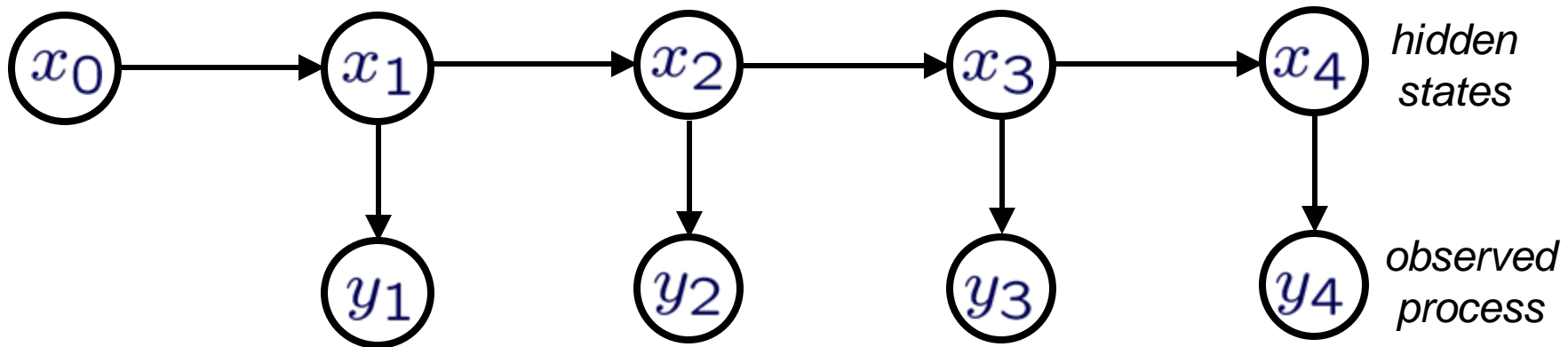
Constraints on valid transition matrices:

$$q_{ij} \geq 0, \quad \sum_{i=1}^N q_{ij} = 1 \quad \text{for all } j$$

$$q_{ij} \triangleq p(x_{t+1} = i | x_t = j)$$

# Hidden Markov Models (Packet Stall Example Cont'd)

- Stall exists due to congestion
- Not directly measurable at runtime (hidden)
- Motivates *hidden Markov models* (*HMM*):



$$p(x_0, x_1, \dots, x_T, y_1, y_2, \dots, y_T) = p(x_0) \prod_{t=1}^T p(x_t | x_{t-1}) p(y_t | x_t)$$

Given  $x_t$ , previous observations do not impact future observations

$$p(y_t, y_{t+1}, \dots, y_T | x_t, y_{t-1}, y_{t-2}, \dots, y_1) = p(y_t, y_{t+1}, \dots, y_T | x_t)$$

# Hidden Markov Models

## Model

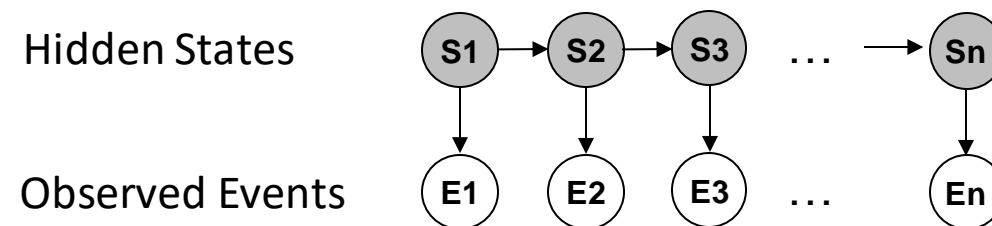
- Set of hidden states  $\mathcal{S} = \{\sigma_1, \dots, \sigma_N\}$
- Set of observable events  $\mathcal{E} = \{\epsilon_1, \dots, \epsilon_M\}$
- Transition probability matrix  $A$
- Observation matrix  $B$
- Initial distribution of hidden states  $\pi$

## Model assumptions

- An observation depends on its hidden state
- A state variable only depends on the immediate previous state (Markov assumption)
- The future observations and the past observations are **conditionally independent** given the current hidden state

## Advantages:

- HMM can model sequential nature of input data (future depends on the past)
- HMM has a linear-chain structure that clearly separates system state and observed events.



## A Hidden Markov model on observed events and system states

$$\begin{aligned} &P(S_1, \dots, S_n, E_1, \dots, E_n) \\ &= P(S_1)P(E_1|S_1) \prod_{i=2}^n P(S_i|S_{i-1})P(E_i|S_i) \end{aligned}$$

# General Inference question

Given the sequence of  $n$  observations  $E_1, E_2, \dots, E_n$ , and the model  $(A, B, \pi)$ , how do we choose a corresponding state sequence  $S_1, S_2, \dots, S_n$  which is optimal in some meaningful sense (i.e., best explains the observations)?

A simpler question: Given the sequence of  $n$  observations  $E_1, E_2, \dots, E_n$ , and the model  $(A, B, \pi)$ , what is the most probable state  $S_t$  at  $t \in \{1, \dots, n\}$ ?

$$\operatorname{argmax}_{j \in \{1, \dots, N\}} P(S_t = \sigma_j | E_1, E_2, \dots, E_n)$$

$$S = \{\sigma_1, \dots, \sigma_N\}$$

# Breaking down the inference question

$$P(S_t | E_1, E_2, \dots, E_n) = \frac{P(E_{t+1}, \dots, E_n | S_t) P(S_t | E_1, \dots, E_t)}{P(E_{t+1}, \dots, E_n | E_1, \dots, E_t)}$$

$P(S_t | E_1, \dots, E_t)$ :

Probability of hidden state at time  $t$  given observation up to time  $t$  (**Forwards algorithm**)

$P(E_{t+1}, \dots, E_n | S_t)$ :

Probability of the future observed sequence given the hidden state at time  $t$  (**Backwards algorithm**)

$P(E_{t+1}, \dots, E_n | E_1, \dots, E_t)$ :

Does not depend on the hidden state (will not affect the maximization because it is just a scaling factor)

# Forwards algorithm: General Expression

Define:  $\alpha_t(j) = P(S_t = \sigma_j | E_1, E_2, \dots, E_t)$  and  $Z_t = P(E_t | E_1, \dots, E_{t-1})$

In general,

$$\alpha_t(j) = \frac{1}{Z_t} P(E_t | S_t = \sigma_j) \sum_{i=1}^N P(S_t = \sigma_j | S_{t-1} = \sigma_i) \alpha_{t-1}(i) \quad Z_t = \sum_{j=1}^N b_t \odot (A^T \alpha_{t-1})$$

Transition probability  $a_{ij}$

Above equation can be written as a matrix for all  $j$ ,

$$\begin{bmatrix} \alpha_t(1) \\ \vdots \\ \alpha_t(j) \\ \vdots \\ \alpha_t(N) \end{bmatrix} \propto \begin{bmatrix} P(E_t | S_t = \sigma_1) \\ \vdots \\ P(E_t | S_t = \sigma_j) \\ \vdots \\ P(E_t | S_t = \sigma_N) \end{bmatrix} \odot \begin{bmatrix} a_{11} & \dots & \dots & \dots & a_{N1} \\ \vdots & \ddots & \dots & \dots & \vdots \\ a_{1j} & \dots & a_{ij} & \dots & a_{Nj} \\ \vdots & \dots & \dots & \ddots & \dots \\ a_{1N} & \dots & \dots & \dots & a_{NN} \end{bmatrix} \begin{bmatrix} \alpha_{t-1}(1) \\ \vdots \\ \alpha_{t-1}(i) \\ \vdots \\ \alpha_{t-1}(N) \end{bmatrix}$$

⊙ Represents elementwise product (Hadamard product)

$$\alpha_t \propto b_t \odot (A^T \alpha_{t-1})$$

$b_t$  is the column of the observation matrix B corresponding to  $E_t$

# Forwards Algorithm: Paleontological Temperature

For observations  $T, L, D, T, L$

$P(S_2|E_1 = T, E_2 = L)$  is,

$$\begin{bmatrix} \alpha_2(H) \\ \alpha_2(C) \end{bmatrix} \propto \begin{bmatrix} 0.5 \\ 0.1 \end{bmatrix} \odot \left( \begin{bmatrix} 0.7 & 0.4 \\ 0.3 & 0.6 \end{bmatrix} \begin{bmatrix} \alpha_1(H) \\ \alpha_1(C) \end{bmatrix} \right)$$

	H	C
H	0.7	0.3
C	0.4	0.6

Transition probability matrix

	T	D	L
H	0.1	0.4	0.5
C	0.7	0.2	0.1

Observation matrix

Similarly,  $P(S_3|E_1 = T, E_2 = L, E_3 = D)$  is,

$$\begin{bmatrix} \alpha_3(H) \\ \alpha_3(C) \end{bmatrix} \propto \begin{bmatrix} 0.4 \\ 0.2 \end{bmatrix} \odot \left( \begin{bmatrix} 0.7 & 0.4 \\ 0.3 & 0.6 \end{bmatrix} \begin{bmatrix} \alpha_2(H) \\ \alpha_2(C) \end{bmatrix} \right)$$

# Forwards Algorithm

1. Input:  $(A, B, \pi)$  and observed sequence  $E_1, \dots, E_n$
2.  $[\alpha_1, Z_1] = \text{normalize}(b_1 \pi)$
3. **for**  $t = 2:n$  **do**  
     $[\alpha_t, Z_t] = \text{normalize}(b_t (A^T \alpha_{t-1}))$
4. return  $\alpha_1, \dots, \alpha_n$  and  $\log(P(E_1, \dots, E_n)) = \sum_t \log(Z_t)$

Note:

Subroutine:  $[v, Z] = \text{normalize}(u)$ :  $Z = \sum_j u_j$ ;  $v_j = u_j/Z$ ;



# Breaking down the inference question

$$P(S_t | E_1, E_2, \dots, E_n) = \frac{P(E_{t+1}, \dots, E_n | S_t) P(S_t | E_1, \dots, E_t)}{P(E_{t+1}, \dots, E_n | E_1, \dots, E_t)}$$

$P(S_t | E_1, \dots, E_t)$ :

Probability of hidden state at time  $t$  given observation up to time  $t$  (**Forwards algorithm**)

$P(E_{t+1}, \dots, E_n | S_t)$ :

Probability of the future observed sequence given the hidden state at time  $t$  (**Backwards algorithm**)

$P(E_{t+1}, \dots, E_n | E_1, \dots, E_t)$ :

Does not depend on the hidden state (will not affect the maximization because it is just a scaling factor)

# Backwards Algorithm

1. Input:  $(A, B, \pi)$  and observed sequence  $E_1, \dots, E_n$
2.  $\beta_n = 1$  ; // initialize  $\beta_n(j)$  to 1 for all states  $\sigma_j$
3. **for**  $t = n - 1$  : **do**  
     $\beta_{t-1} = A(b_t \ \beta_t)$
4. return  $\beta_1, \dots, \beta_n$

# Breaking down the inference question

$$P(S_t | E_1, E_2, \dots, E_n) = \frac{P(E_{t+1}, \dots, E_n | S_t) P(S_t | E_1, \dots, E_t)}{P(E_{t+1}, \dots, E_n | E_1, \dots, E_t)}$$

$P(S_t | E_1, \dots, E_t)$ :

Probability of hidden state at time  $t$  given observation up to time  $t$  (**Forwards algorithm**)

$P(E_{t+1}, \dots, E_n | S_t)$ :

Probability of the future observed sequence given the hidden state at time  $t$  (**Backwards algorithm**)

$P(E_{t+1}, \dots, E_n | E_1, \dots, E_t)$ :

Does not depend on the hidden state (will not affect the maximization because it is just a scaling factor)

# Inference: Most likely state

- Forwards-backwards algorithm gives  $P(S_t = \sigma_j | E_1, \dots, E_n)$  for all  $j$
- Find the **individually most likely state** at time  $t$  given all observations

$$S_t^* = \operatorname{argmax}_{j \in \{1, \dots, N\}} \gamma_t(j)$$

# Optimality of inference

- In the inference problem we attempting to uncover the hidden part of HMM, i.e., find the “correct” state sequence
- It is impossible to find the “correct” state sequence (solution)
- Use optimality criterion to find the “best” possible solution
- **Several reasonable criteria** exist and is a strong function of the intended application
  - **Most likely state given observations**
    - Application in finding average statistics, expected number of correct states
    - Solved using **Forwards-Backwards algorithm**
  - **Single best sequence that maximises probability of observed events**
    - Application in continuous speech recognition
    - Solved using **Viterbi algorithm**

# Hidden Markov Model – Online Battle Simulator Game

You are playing an online battle simulator named *WT*. To balance the gameplay, the game has a "balancer" that decides the difficulty of the game before the start of each round based on the previous round's difficulty. There are two difficulty levels *hard* (H) and *easy* (E), which are hidden from you. For the first round you play, the probability of the round being H is 0.25 and the probability of the round being E is 0.75. The transition probabilities between the states is given as A.

$$A = \begin{matrix} & \begin{matrix} H & E \end{matrix} \\ \begin{matrix} H \\ E \end{matrix} & \begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix} \end{matrix}$$

$$\pi = \begin{matrix} & \begin{matrix} H & E \end{matrix} \\ \begin{bmatrix} 0.25 & 0.75 \end{bmatrix} \end{matrix}$$

Each game round can end in either a win (W), draw (D) or loss(L). The observation matrix is given below.

$$B = \begin{matrix} & \begin{matrix} W & D & L \end{matrix} \\ \begin{matrix} H \\ E \end{matrix} & \begin{bmatrix} 1/6 & 2/6 & 3/6 \\ 4/6 & 1/6 & 1/6 \end{bmatrix} \end{matrix}$$

Suppose you played for eight consecutive rounds and have the following result: **W, D, D, L, W, L, D, L**. You would like to infer the difficulty for each round you have played by using an HMM.

# HMM Solution (1)

## Forward algorithm

1. Input:  $(A, B, \pi)$  and observed sequence  $E_1, \dots, E_n$
2.  $[\alpha_1, Z_1] = \text{normalize}(b_1 \odot \pi)$
3. **for**  $t = 2:n$  **do**  
     $[\alpha_t, Z_t] = \text{normalize}(b_t \odot (A^T \alpha_{t-1}))$
4. return  $\alpha_1, \dots, \alpha_n$  and  $\log(P(E_1, \dots, E_n)) = \sum_t \log(Z_t)$
5. Subroutine:  $[v, Z] = \text{normalize}(u)$ :  $Z = \sum_j u_j$ ;  
     $v_j = u_j/Z$ ;

NOTE:  $\odot$  represents elementwise product (Hadamard product)

## Backward algorithm

1. Input:  $(A, B, \pi)$  and observed sequence  $E_1, \dots, E_n$
2.  $\beta_n = 1$  ; // initialize  $\beta_n(j)$  to 1 for all states  $\sigma_j$
3. **for**  $t = n - 1:1$  **do**  
     $\beta_{t-1} = A(b_t \odot \beta_t)$
4. return  $\beta_1, \dots, \beta_n$

# HMM Solution (2)

What is the most likely state at time step  $t=2$  given the evidence?

$$S_2^* = \operatorname{argmax}_{j \in \{H,E\}} \gamma_2(j) = \operatorname{argmax}_{j \in \{H,E\}} \alpha_2(j) * \beta_2(j)$$

$$\alpha_1 \propto b_1 \odot \pi = \begin{bmatrix} 1/6 \\ 4/6 \end{bmatrix} \odot \begin{bmatrix} 1/4 \\ 3/4 \end{bmatrix} = \begin{bmatrix} 1/24 \\ 1/2 \end{bmatrix}$$

$$\operatorname{Normalize}(\alpha_1) = \begin{bmatrix} 1/13 \\ 12/13 \end{bmatrix}$$

$$\alpha_2 \propto b_2 \odot (A^T \alpha_1) = \begin{bmatrix} 2/6 \\ 1/6 \end{bmatrix} \odot \left( \begin{bmatrix} 0.8 & 0.5 \\ 0.2 & 0.5 \end{bmatrix} \begin{bmatrix} 1/13 \\ 12/13 \end{bmatrix} \right) = \begin{bmatrix} 34/195 \\ 31/390 \end{bmatrix}$$

$$\operatorname{Normalize}(\alpha_2) = \begin{bmatrix} 68/99 \\ 31/99 \end{bmatrix}$$

Computing  $\beta_2, \gamma_2$  is left as an exercise.

$$\pi = \begin{matrix} & H & E \\ \begin{matrix} H \\ E \end{matrix} & \begin{bmatrix} 0.25 & 0.75 \end{bmatrix} \end{matrix}$$

$$A = \begin{matrix} & H & E \\ \begin{matrix} H \\ E \end{matrix} & \begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix} \end{matrix}$$

$$B = \begin{matrix} & W & D & L \\ \begin{matrix} H \\ E \end{matrix} & \begin{bmatrix} 1/6 & 2/6 & 3/6 \\ 4/6 & 1/6 & 1/6 \end{bmatrix} \end{matrix}$$

Observations

W	D	D	L	W	L	D	L
---	---	---	---	---	---	---	---

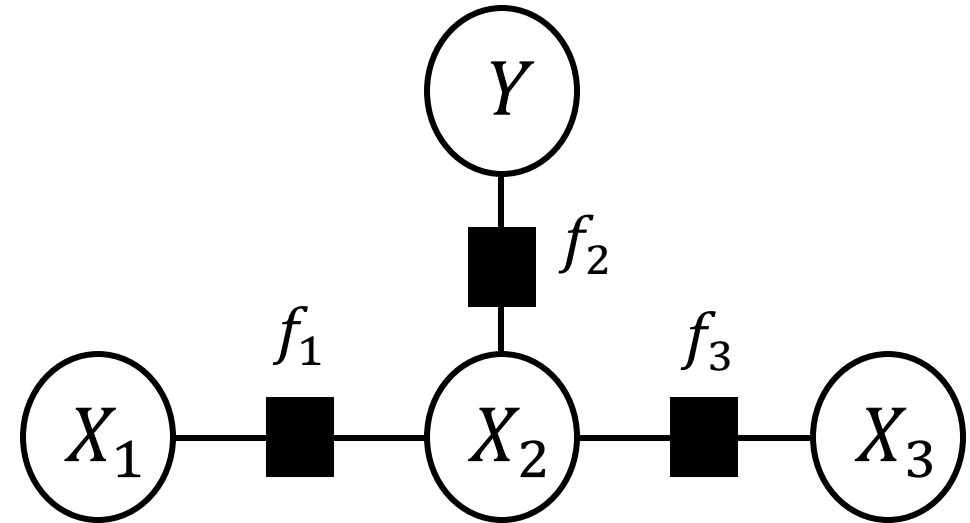


# Factor Graphs Belief Propagation

Suppose  $X_1$  is observed to be 0.

Using belief propagation with the factor graph to the right, calculate

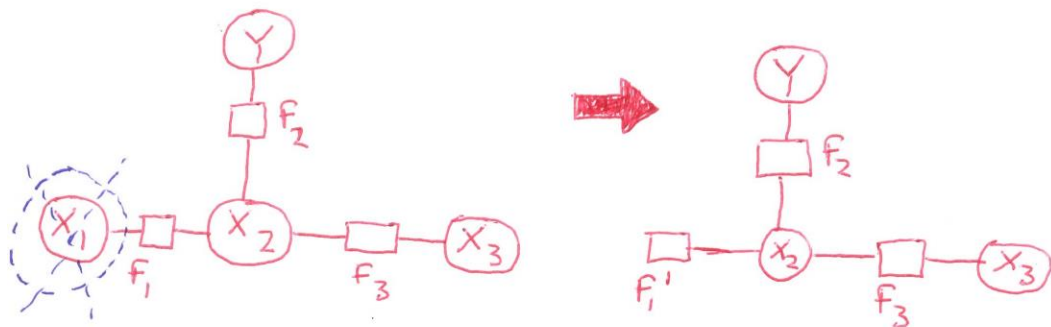
- $P(X_2 | X_1 = 0)$
- $P(Y | X_1 = 0)$
- $P(X_3 | X_1 = 0)$



$X_1$	$X_2$	$f_1(X_1, X_2)$	$X_2$	$Y$	$f_2(X_2, Y)$	$X_2$	$X_3$	$f_3(X_2, X_3)$
0	0	0.9	0	0	0.1	0	0	1
0	1	1	0	1	1	0	1	0.1
1	0	1	1	0	1	1	0	0.1
1	1	0.9	1	1	0.1	1	1	1

Since  $X_1$  has been observed to be 0, we can make the following adjustments:

- Graph structure:



- Factor table for  $f_1$

$X_1$	$X_2$	$f_1(X_1, X_2)$
0	0	0.9
0	1	1
1	0	+
+	+	0.9

$X_2$	$f_1'(X_2)$
0	0.9
1	1

- Factor tables for  $f_2$  and  $f_3$  remain the same.

Let's find  $P(X_2 | X_1=0)$

$$\mu_{f_1' \rightarrow X_2}(X_2) = f_1'(X_2) = \begin{bmatrix} 0.9 \\ 1 \end{bmatrix} \begin{matrix} (X_2=0) \\ (X_2=1) \end{matrix}$$

$$\begin{aligned} \mu_{F_2 \rightarrow X_2}(X_2) &= \sum_Y f_2(X_2, Y) \mu_{Y \rightarrow F_2}(Y) \\ &= \begin{bmatrix} 0.1 \\ 1 \end{bmatrix} \begin{matrix} (X_2=0, Y=0) \\ (X_2=1, Y=0) \end{matrix} \times 1 \\ &\quad + \begin{bmatrix} 1 \\ 0.1 \end{bmatrix} \begin{matrix} (X_2=0, Y=1) \\ (X_2=1, Y=1) \end{matrix} \times 1 = \begin{bmatrix} 1.1 \\ 1.1 \end{bmatrix} \begin{matrix} (X_2=0) \\ (X_2=1) \end{matrix} \end{aligned}$$

$$\begin{aligned} \mu_{F_3 \rightarrow X_2}(X_2) &= \sum_{X_3} F_3(X_2, X_3) \mu_{X_3 \rightarrow F_3}(X_3) \\ &= \begin{bmatrix} 1 \\ 0.1 \end{bmatrix} \begin{matrix} (X_2=0, X_3=0) \\ (X_2=1, X_3=0) \end{matrix} \times 1 + \begin{bmatrix} 0.1 \\ 1 \end{bmatrix} \begin{matrix} (X_2=0, X_3=1) \\ (X_2=1, X_3=1) \end{matrix} \times 1 \\ &= \begin{bmatrix} 1.1 \\ 1.1 \end{bmatrix} \begin{matrix} (X_2=0) \\ (X_2=1) \end{matrix} \end{aligned}$$

$$P(X_2 | X_1=0) = \frac{1}{Z} (\mu_{f_1' \rightarrow X_2} \times \mu_{F_2 \rightarrow X_2} \times \mu_{F_3 \rightarrow X_2})$$

$$= \frac{1}{Z} \left( \begin{bmatrix} 0.9 \\ 1 \end{bmatrix} \odot \begin{bmatrix} 1.1 \\ 1.1 \end{bmatrix} \odot \begin{bmatrix} 1.1 \\ 1.1 \end{bmatrix} \right)$$

$$= \frac{1}{Z} \begin{bmatrix} 1.089 \\ 1.21 \end{bmatrix}$$

$$= \frac{1}{1.089 + 1.21} \begin{bmatrix} 1.089 \\ 1.21 \end{bmatrix}$$

$$= \begin{bmatrix} 0.47 \\ 0.53 \end{bmatrix} \begin{matrix} (P(x_2=0|x_1=0)) \\ (P(x_2=1|x_1=0)) \end{matrix}$$

Next, let's find  $P(Y|x_1=0)$

$$\mu_{f_1 \rightarrow x_2}(x_2) = \begin{bmatrix} 0.9 \\ 1 \end{bmatrix} \begin{matrix} (x_2=0) \\ (x_2=1) \end{matrix}$$

$$\mu_{f_3 \rightarrow x_2}(x_2) = \begin{bmatrix} 1.1 \\ 1.1 \end{bmatrix} \begin{matrix} (x_2=0) \\ (x_2=1) \end{matrix} \quad \leftarrow \text{This was calculated earlier}$$

$$\begin{aligned} \mu_{x_2 \rightarrow f_2}(x_2) &= \mu_{f_1 \rightarrow x_2}(x_2) \times \mu_{f_3 \rightarrow x_2}(x_2) \\ &= \begin{bmatrix} 0.9 \\ 1 \end{bmatrix} \odot \begin{bmatrix} 1.1 \\ 1.1 \end{bmatrix} = \begin{bmatrix} 0.99 \\ 1.1 \end{bmatrix} \begin{matrix} (x_2=0) \\ (x_2=1) \end{matrix} \end{aligned}$$

$$\mu_{f_2 \rightarrow Y}(Y) = \sum_{x_2} f_2(x_2, Y) \mu_{x_2 \rightarrow f_2}(x_2)$$

$$= \begin{bmatrix} 0.1 \\ 1 \end{bmatrix} \begin{matrix} (x_2=0, Y=0) \\ (x_2=0, Y=1) \end{matrix} \times 0.99$$

$$+ \begin{bmatrix} 1 \\ 0.1 \end{bmatrix} \begin{matrix} (x_2=1, Y=0) \\ (x_2=1, Y=1) \end{matrix} \times 1.1$$

$$= \begin{bmatrix} 1.199 \\ 1.1 \end{bmatrix} \begin{matrix} (Y=0) \\ (Y=1) \end{matrix}$$

$$P(Y|x_1=0) = \frac{1}{Z} \mu_{f_2 \rightarrow Y}(Y)$$

$$= \frac{1}{Z} \begin{bmatrix} 1.199 \\ 1.1 \end{bmatrix}$$

$$= \frac{1}{1.199 + 1.1} \begin{bmatrix} 1.199 \\ 1.1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.52 \\ 0.48 \end{bmatrix} \begin{matrix} (P(Y=0|x_1=0)) \\ (P(Y=1|x_1=0)) \end{matrix}$$

Finally, let's calculate  $\Psi(x_3 | x_1=0)$

$$\mu_{F_1 \rightarrow x_2}(x_2) = \begin{bmatrix} 0.9 \\ 1 \end{bmatrix} \begin{matrix} (x_2=0) \\ (x_2=1) \end{matrix}$$

$$\mu_{F_2 \rightarrow x_2}(x_2) = \begin{bmatrix} 1.1 \\ 1.1 \end{bmatrix} \begin{matrix} (x_2=0) \\ (x_2=1) \end{matrix} \leftarrow \text{This was already calculated previously}$$

$$\begin{aligned} \mu_{x_2 \rightarrow F_3}(x_2) &= \mu_{F_1 \rightarrow x_2} \times \mu_{F_2 \rightarrow x_2} \\ &= \begin{bmatrix} 0.9 \\ 1 \end{bmatrix} \odot \begin{bmatrix} 1.1 \\ 1.1 \end{bmatrix} = \begin{bmatrix} 0.99 \\ 1.1 \end{bmatrix} \begin{matrix} (x_2=0) \\ (x_2=1) \end{matrix} \end{aligned}$$

$$\begin{aligned} \mu_{F_3 \rightarrow x_3}(x_3) &= \sum_{x_2} F_3(x_2, x_3) \times \mu_{x_2 \rightarrow F_3}(x_2) \\ &= \begin{bmatrix} 1 \\ 0.1 \end{bmatrix} \begin{matrix} (x_2=0, x_3=0) \\ (x_2=0, x_3=1) \end{matrix} \times 0.99 \\ &\quad + \begin{bmatrix} 0.1 \\ 1 \end{bmatrix} \begin{matrix} (x_2=1, x_3=0) \\ (x_2=1, x_3=1) \end{matrix} \times 1.1 \\ &= \begin{bmatrix} 1.1 \\ 1.199 \end{bmatrix} \begin{matrix} (x_3=0) \\ (x_3=1) \end{matrix} \end{aligned}$$

$$\text{Thus, } \Psi(x_3 | x_1=0) = \frac{1}{Z} \mu_{F_3 \rightarrow x_3}(x_3)$$

$$= \frac{1}{Z} \begin{bmatrix} 1.1 \\ 1.199 \end{bmatrix}$$

$$= \frac{1}{1.1 + 1.199} \begin{bmatrix} 1.1 \\ 1.199 \end{bmatrix}$$

$$= \begin{bmatrix} 0.48 \\ 0.52 \end{bmatrix} \begin{matrix} (P(x_3=0 | x_1=0)) \\ (P(x_3=1 | x_1=0)) \end{matrix}$$

# Gradient Descent

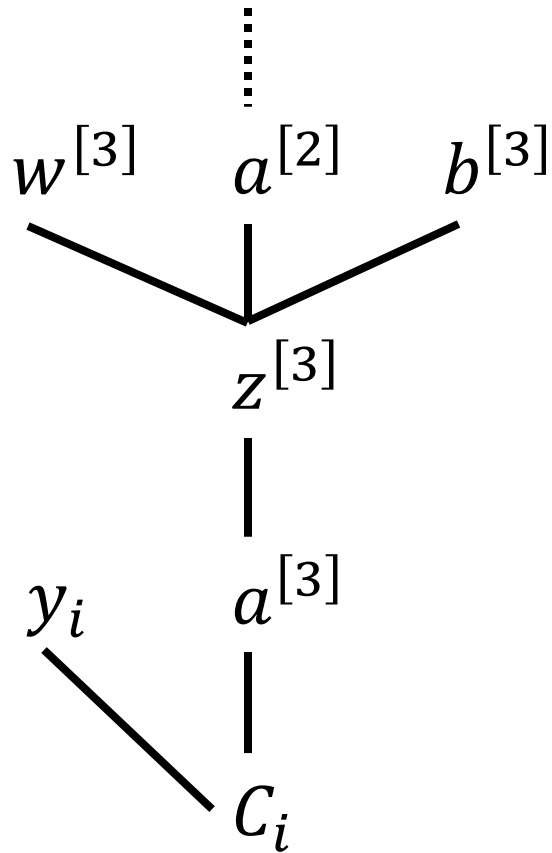
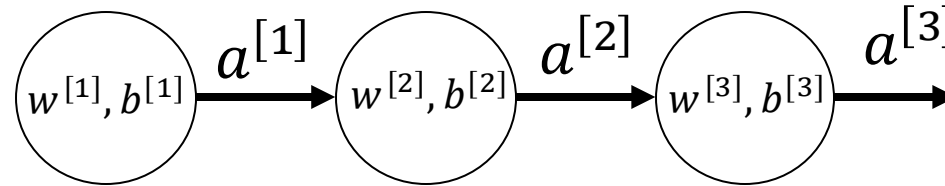
Given  $N$  training data points  $\{(\mathbf{x}^k, y^k)\}$  for  $k = \{1, \dots, N\}$ ,  $\mathbf{x}^k \in R^d$ , and ground truth  $y^k$ , we seek a linear regressor  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  optimizing the loss function  $L(z) = (y - f(x))^4$

1. Find the gradient and gradient descent update equation to find  $\mathbf{w}$ .
2. Suppose you also want to include a penalty term  $\lambda \|\mathbf{w}\|^2$  to the overall loss function. Derive the gradient for gradient descent to update  $\mathbf{w}$ .

# Back Propagation: Chain Rule

Apply *chain rule of derivative* to update parameters:

$$\mathbf{w}^{[1]}, \mathbf{b}^{[1]}, \mathbf{w}^{[2]}, \mathbf{b}^{[2]}$$



- Our goal is to figure out how to update the weights to minimize the cost in the next iteration of gradient descent

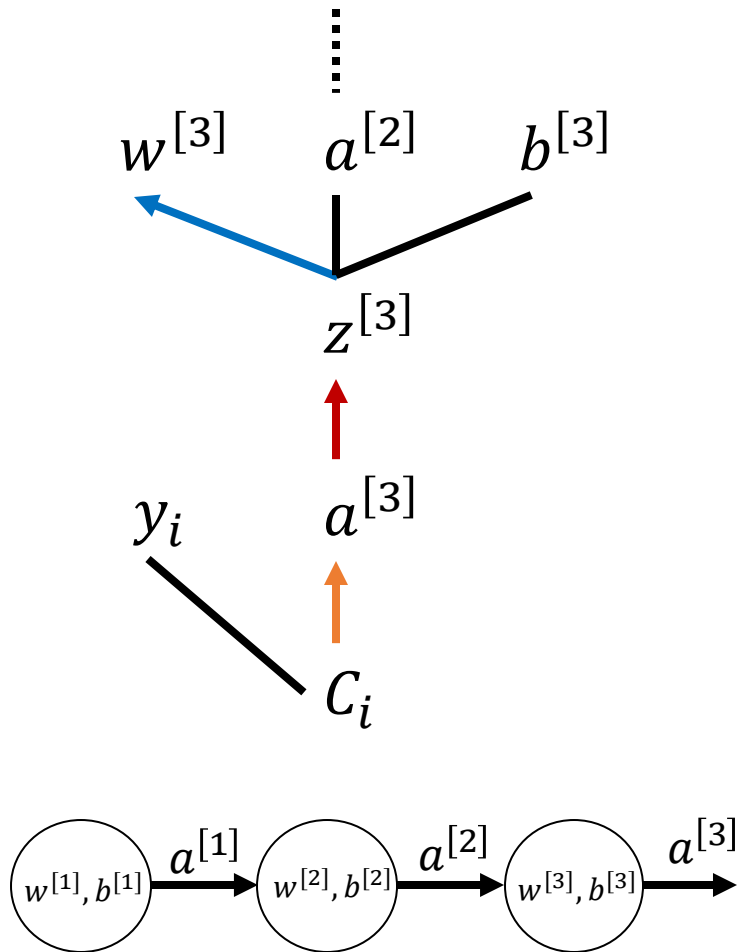
- To begin, we wish to calculate  $\frac{\partial C_i}{\partial w^{[3]}}$

- Recall that

$$\begin{aligned} C_i &= (a^{[3]} - y_i)^2 \\ &= (\sigma(z^{[3]}) - y_i)^2 \\ &= (\sigma(w^{[3]}a^{[2]} + b^{[3]}) - y_i)^2 \end{aligned}$$

- We can visualize this relationship with the dependency tree to the left

# Back Propagation: Chain Rule



- Now, let's compute  $\frac{\partial C_i}{\partial w^{[3]}}$

$$\frac{\partial C_i}{\partial w^{[3]}} = \frac{\partial C_i}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial w^{[3]}}$$

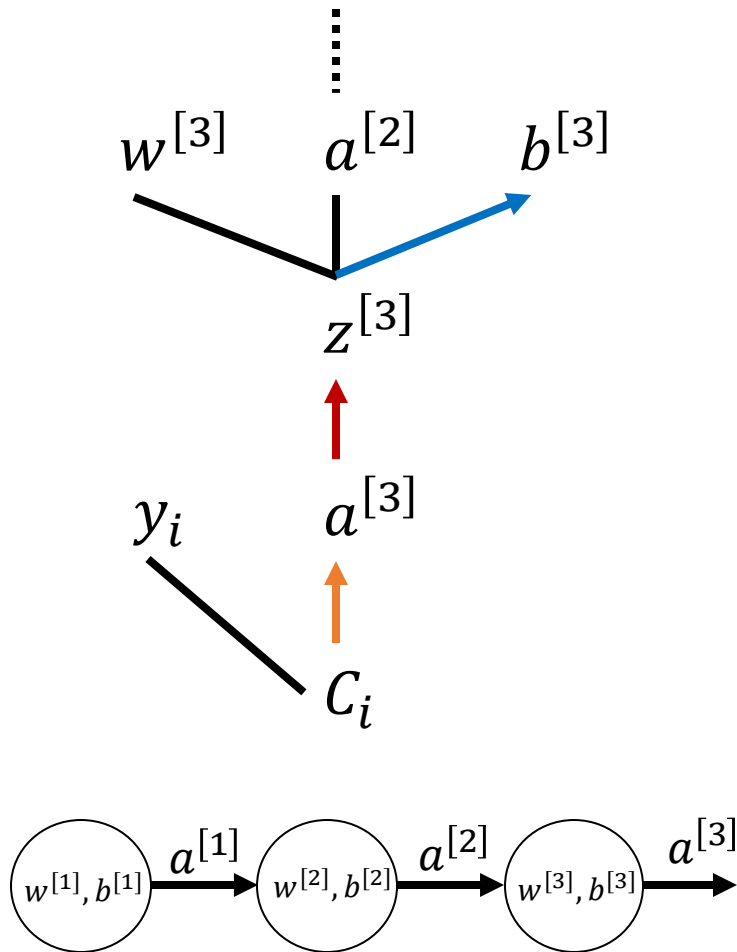
$$\frac{\partial C_i}{\partial a^{[3]}} = \frac{\partial (a^{[3]} - y_i)^2}{\partial a^{[3]}} = 2(a^{[3]} - y_i)$$

$$\frac{\partial a^{[3]}}{\partial z^{[3]}} = \frac{\partial \sigma(z^{[3]})}{\partial z^{[3]}} = \sigma(z^{[3]}) (1 - \sigma(z^{[3]}))$$

$$\frac{\partial z^{[3]}}{\partial w^{[3]}} = \frac{\partial (w^{[3]} a^{[2]} + b^{[3]})}{\partial w^{[3]}} = a^{[2]}$$

$$\Rightarrow \frac{\partial C_i}{\partial w^{[3]}} = 2(a^{[3]} - y_i) \sigma(z^{[3]}) (1 - \sigma(z^{[3]})) a^{[2]}$$

# Back Propagation: Chain Rule



- Similarly, we can compute  $\frac{\partial C_i}{\partial b^{[3]}}$

$$\frac{\partial C_i}{\partial b^{[3]}} = \frac{\partial C_i}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial b^{[3]}}$$

$$\frac{\partial C_i}{\partial a^{[3]}} = \frac{\partial (a^{[3]} - y_i)^2}{\partial a^{[3]}} = 2(a^{[3]} - y_i)$$

$$\frac{\partial a^{[3]}}{\partial z^{[3]}} = \frac{\partial \sigma(z^{[3]})}{\partial z^{[3]}} = \sigma(z^{[3]}) (1 - \sigma(z^{[3]}))$$

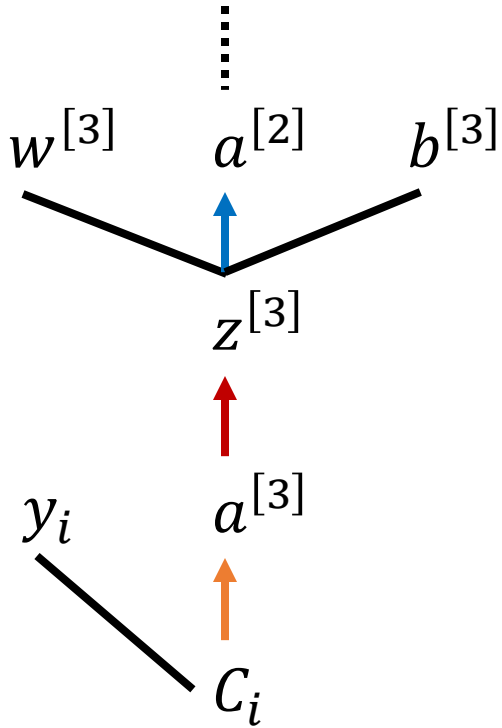
$$\frac{\partial z^{[3]}}{\partial b^{[3]}} = \frac{\partial (w^{[3]} a^{[2]} + b^{[3]})}{\partial b^{[3]}} = 1$$

$$\Rightarrow \frac{\partial C_i}{\partial b^{[3]}} = 2(a^{[3]} - y_i) \sigma(z^{[3]}) (1 - \sigma(z^{[3]}))$$



# Back Propagation: Chain Rule

- Finally, we compute  $\frac{\partial C_i}{\partial a^{[2]}}$



$$\frac{\partial C_i}{\partial a^{[2]}} = \frac{\partial C_i}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial a^{[2]}}$$

$$\frac{\partial C_i}{\partial a^{[3]}} = \frac{\partial (a^{[3]} - y_i)^2}{\partial a^{[3]}} = 2(a^{[3]} - y_i)$$

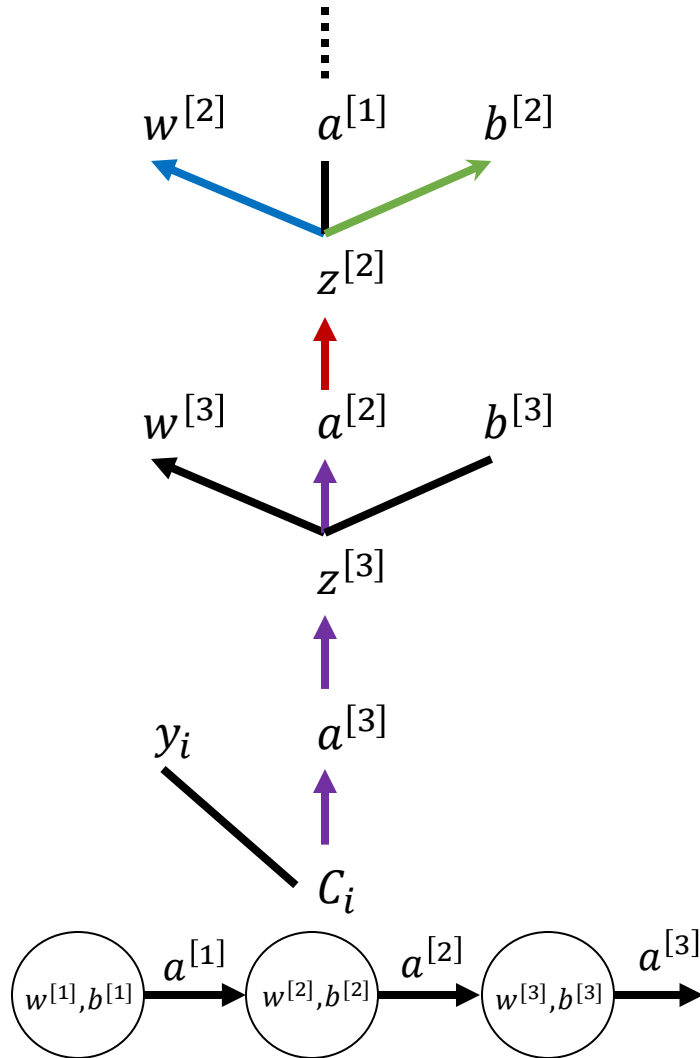
$$\frac{\partial a^{[3]}}{\partial z^{[3]}} = \frac{\partial \sigma(z^{[3]})}{\partial z^{[3]}} = \sigma(z^{[3]})(1 - \sigma(z^{[3]}))$$

$$\frac{\partial z^{[3]}}{\partial a^{[2]}} = \frac{\partial (w^{[3]}a^{[2]} + b^{[3]})}{\partial a^{[2]}} = w^{[3]}$$

$$\begin{array}{c} \text{---} \\ \vdots \\ w^{[3]} \quad a^{[2]} \quad b^{[3]} \\ \swarrow \quad \uparrow \quad \searrow \\ z^{[3]} \\ \uparrow \\ y_i \quad a^{[3]} \\ \swarrow \quad \uparrow \\ C_i \end{array}$$

$$\begin{array}{c} \text{---} \\ \vdots \\ w^{[1], b^{[1]}} \xrightarrow{a^{[1]}} w^{[2], b^{[2]}} \xrightarrow{a^{[2]}} w^{[3], b^{[3]}} \xrightarrow{a^{[3]}} \end{array} \Rightarrow \frac{\partial C_i}{\partial a^{[2]}} = 2(a^{[3]} - y_i) \sigma(z^{[3]}) (1 - \sigma(z^{[3]})) w^{[3]}$$

# Back Propagation: Chain Rule



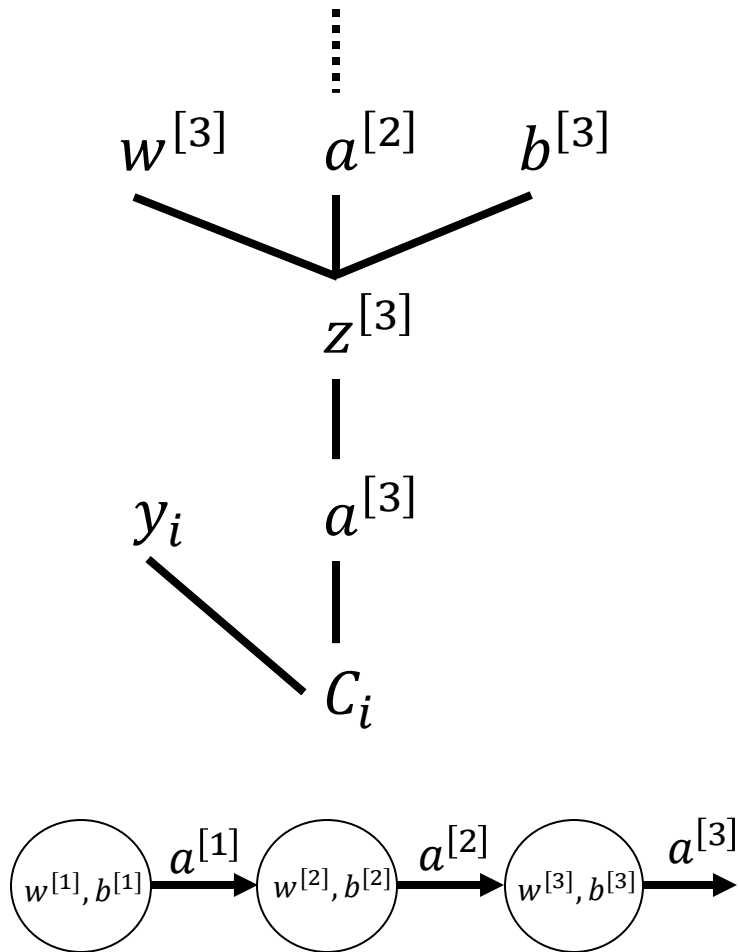
- We can now *propagate* the gradient calculations we have done *backwards* in order to find the gradients for weights/biases in layer 2

Has been computed already in previous slide

$$\begin{aligned} \frac{\partial C_i}{\partial w^{[2]}} &= \frac{\partial C_i}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial w^{[2]}} \\ &= \frac{\partial C_i}{\partial a^{[2]}} \sigma(z^{[2]}) (1 - \sigma(z^{[2]})) a^{[1]} \end{aligned}$$

$$\begin{aligned} \frac{\partial C_i}{\partial b^{[2]}} &= \frac{\partial C_i}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial b^{[2]}} \\ &= \frac{\partial C_i}{\partial a^{[2]}} \sigma(z^{[2]}) (1 - \sigma(z^{[2]})) \end{aligned}$$

# Back Propagation: Full Cost Function



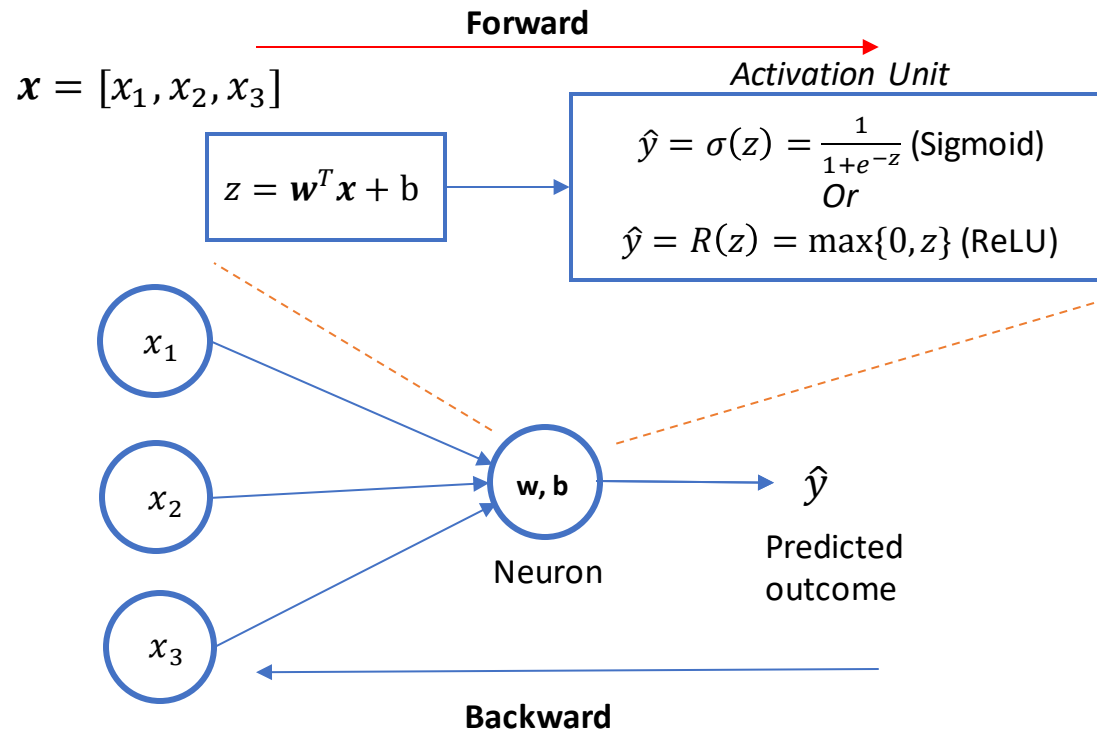
- In order to the partial derivative of the full cost function  $C$  with respect to a weight (say,  $w^{[3]}$ ), we need to average the gradients of cost with respect to that weight for all  $n$  training samples

$$\frac{\partial C}{\partial w^{[3]}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial C_i}{\partial w^{[3]}}$$

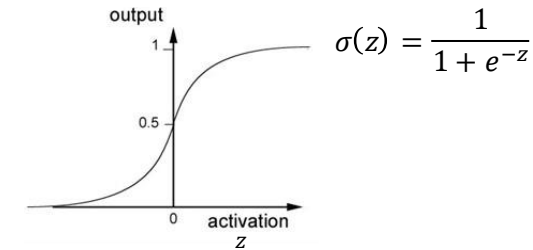
# Apply *chain rule of derivative* to update perceptron Model

parameters:  $\mathbf{b}^{[1]}, \mathbf{w}^{[2]}, \mathbf{b}^{[2]}$

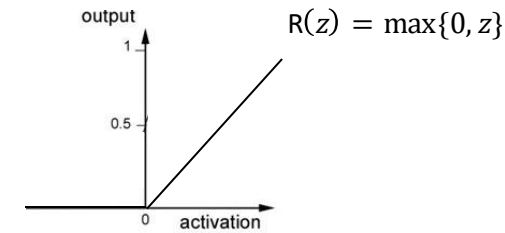
- The core of the neural network is perceptron model



Sigmoid Function



ReLU Function



**Update Rule (Backward):**

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla J(\mathbf{w}_t)$$

$\eta$  : Learning rate

**Loss**

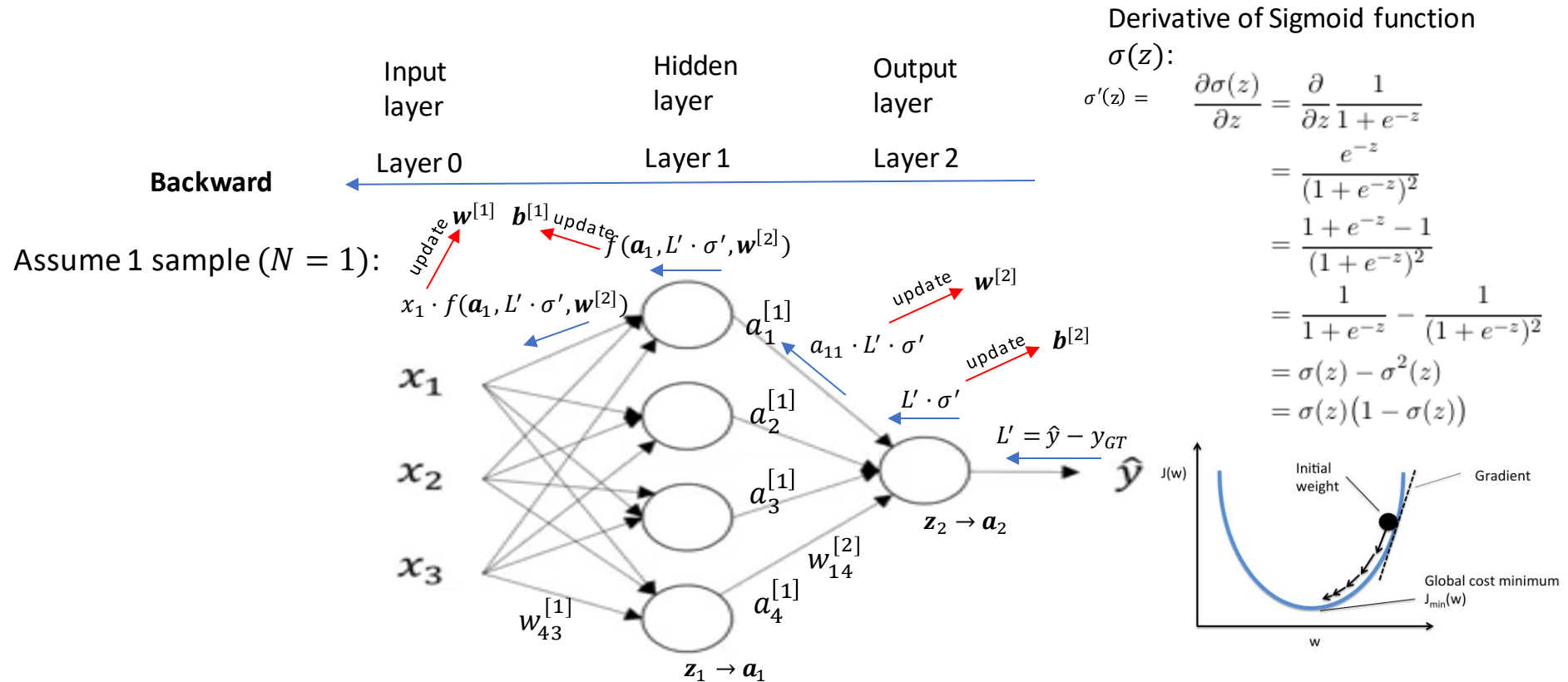
$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{w} \cdot \mathbf{x}^{(i)}, y^{(i)})$$

$N$ : number of samples     $\mathbf{x}^{(i)}$ : feature of  $i^{th}$  sample

**Computing Gradient**

$$\nabla J(\mathbf{w}_0) = \left( \frac{\partial J(\mathbf{w})}{\partial w_0}, \frac{\partial J(\mathbf{w})}{\partial w_1}, \dots, \frac{\partial J(\mathbf{w})}{\partial w_n} \right)_{\mathbf{w}_0}$$

# Backpropagation: Gradient Descent



Purpose of backpropagation:

Apply *chain rule of derivative* to update parameters:  $w^{[1]}$ ,  $b^{[1]}$ ,  $w^{[2]}$ ,  $b^{[2]}$

# Gradient Descent

1. Find the gradient and gradient descent update equation to find  $\mathbf{w}$ .

$$L = \sum_{i=1}^N (y^i - \mathbf{w}^T \mathbf{x}^i)^4$$

$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \quad \nabla_{\mathbf{w}} L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial w_d} \end{bmatrix}$$

$$\frac{\partial L}{\partial w_1} = \sum_{i=1}^N 4 (y^i - \mathbf{w}^T \mathbf{x}^i)^3 (-x_1^i)$$

$$\frac{\partial L}{\partial w_2} = \sum_{i=1}^N 4 (y^i - \mathbf{w}^T \mathbf{x}^i)^3 (-x_2^i)$$

$$\vdots$$

$$\frac{\partial L}{\partial w_d} = \sum_{i=1}^N 4 (y^i - \mathbf{w}^T \mathbf{x}^i)^3 (-x_d^i)$$

$$\nabla_{\mathbf{w}} L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial w_d} \end{bmatrix} = \sum_{i=1}^N 4 (y^i - \mathbf{w}^T \mathbf{x}^i)^3 \begin{bmatrix} -x_1^i \\ -x_2^i \\ \vdots \\ -x_d^i \end{bmatrix} = - \sum_{i=1}^N 4 (y^i - \mathbf{w}^T \mathbf{x}^i)^3 \mathbf{x}^i$$

$$\begin{aligned} \mathbf{w}^{t+1} &= \mathbf{w}^t - \eta \nabla_{\mathbf{w}} L \\ &= \mathbf{w}^t + \eta \sum_{i=1}^N 4 (y^i - \mathbf{w}^t \mathbf{x}^i)^3 \mathbf{x}^i \end{aligned}$$

(Note  $\mathbf{x}^i$  is a  $k$ -dimensional vector)

2. Suppose you also want to include a penalty term  $\lambda \|\mathbf{w}\|^2$  to the overall loss function. Derive the gradient for gradient descent to update  $\mathbf{w}$ .

$$L = \sum_{i=1}^N (y^i - \mathbf{w}^T \mathbf{x}^i)^2 + \lambda \|\mathbf{w}\|^2$$

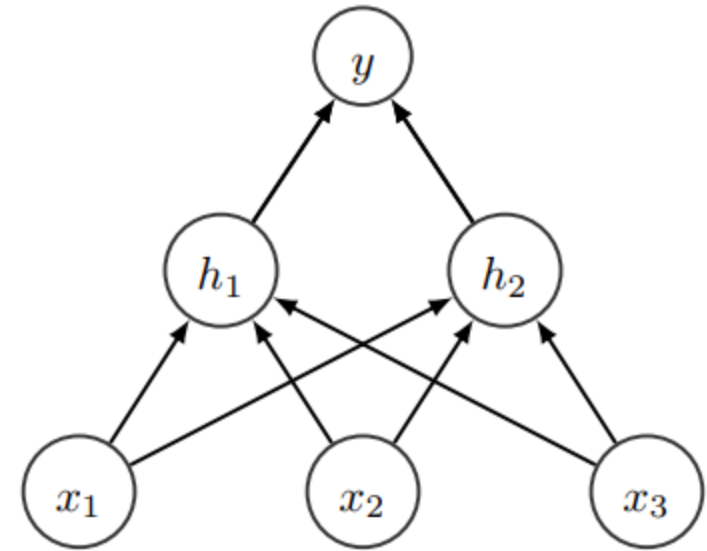
$$\text{gradient: } \nabla_{\mathbf{w}} L = - \sum_{i=1}^N 2(y^i - \mathbf{w}^T \mathbf{x}^i) \mathbf{x}^i + 2\lambda \mathbf{w}$$

$$\begin{aligned} \text{gradient update: } \mathbf{w}^{t+1} &= \mathbf{w}^t - \eta \nabla_{\mathbf{w}} L \\ &= \mathbf{w}^t - \eta \left( - \sum_{i=1}^N 2(y^i - \mathbf{w}^t{}^T \mathbf{x}^i) \mathbf{x}^i + 2\lambda \mathbf{w}^t \right) \\ &= (1 - 2\eta\lambda) \mathbf{w}^t + \eta \sum_{i=1}^N 2(y^i - \mathbf{w}^t{}^T \mathbf{x}^i) \mathbf{x}^i \end{aligned}$$

(Note both  $\mathbf{w}^t$  and  $\mathbf{x}^i$  are  $k$ -dimensional vectors)

# Neural Networks Backpropagation

Consider the neural network given alongside. The hidden units and output layer has ReLU activation function. The loss function is given by  $L(y, y) = \frac{1}{2} (y - t)^2$  where  $t$  is the target value. For simplicity, assume that the bias terms are 0. Weights connecting input to hidden layer and hidden layer to output layer are given by  $W$  and  $V$  respectively.



1. Write the forward equation to map input to output.
2. Compute the output and backpropagation for  $x = [1, 2, 1]$  and  $t = 1$ .

$$W = \begin{bmatrix} 1 & 0 & 1 \\ 1 & -1 & 0 \end{bmatrix}$$

$$V = \begin{bmatrix} 0 & 1 \end{bmatrix}$$



# Neural Network Backprop:

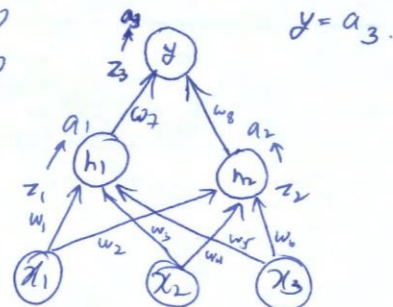
$$1. a_3 = R(z_3) \text{ where } R(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

$$z_3 = a_1 w_7 + a_2 w_8$$

$$a_1 = R(z_1); a_2 = R(z_2)$$

$$z_1 = w_1 x_1 + w_3 x_2 + w_5 x_3$$

$$z_2 = w_2 x_1 + w_4 x_2 + w_6 x_3$$



$$\therefore y = R(w_7 \cdot R(w_1 x_1 + w_3 x_2 + w_5 x_3) + w_8 \cdot R(w_2 x_1 + w_4 x_2 + w_6 x_3))$$

$$2. W = \begin{bmatrix} 1 & 0 & 1 \\ 1 & -1 & 0 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 & w_5 \\ w_2 & w_4 & w_6 \end{bmatrix} \quad V = \begin{bmatrix} 0 & 1 \end{bmatrix} = \begin{bmatrix} w_7 & w_8 \end{bmatrix}$$

substituting  $x = [1, 2, 1] = [x_1, x_2, x_3]$ , we get:

$$z_1 = 2, z_2 = -1 \Rightarrow a_1 = 2, a_2 = 0$$

$$\Rightarrow z_3 = 0 \Rightarrow a_3 = 0 \Rightarrow y = 0.$$

Backpropagation:

$$L = \frac{1}{2} (y - t)^2 = \frac{1}{2} (0 - 1)^2 = \frac{1}{2}.$$

$$\frac{\partial L}{\partial w_7} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_7} = (y - t) \cdot s(z_3) \cdot a_1$$

$$= (0 - 1) \cdot s(0) \cdot 2$$

$$= 0.$$

$$\text{we: } s(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1} = (y - t) \cdot s(z_3) \cdot w_7 \cdot s(z_1) \cdot x_1$$

$$= (0 - 1) \cdot s(0) \cdot 0 \cdot s(2) \cdot 1$$

$$= 0$$

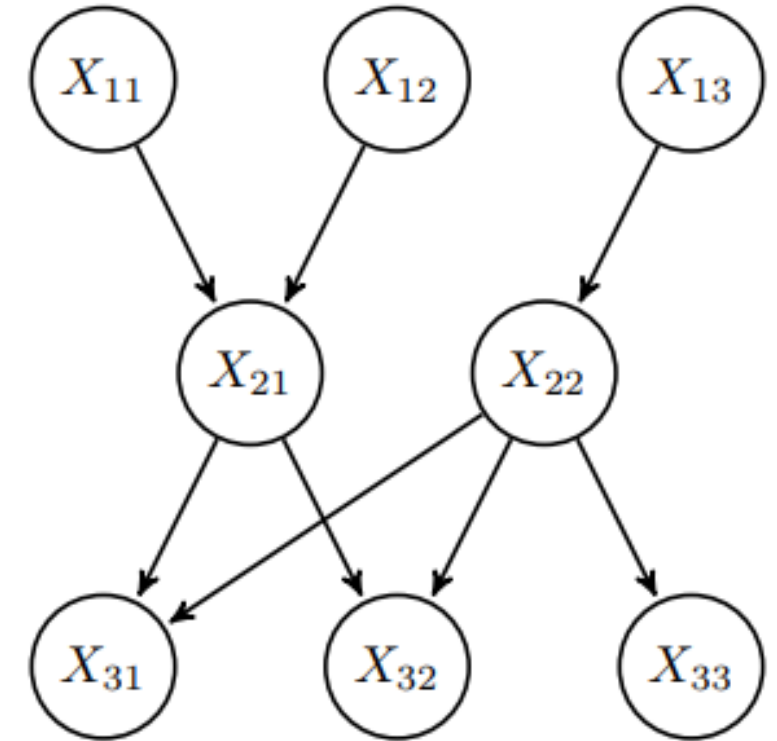
similar computations would give

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial w_4} = \frac{\partial L}{\partial w_5} = \frac{\partial L}{\partial w_6} = \frac{\partial L}{\partial w_8} = 0$$

# Bayesian Network Question

Consider the Bayesian Network alongside with binary variables. Answer the following questions.

1. Is there any variable(s) conditionally independent of  $X_{33}$  given  $X_{11}$  and  $X_{12}$ ? If so, list all.
2. Is there any variable(s) conditionally independent of  $X_{33}$  given  $X_{22}$ ? If so, list all.
3. How many parameters are required to specify the factorized joint distribution?
4. Express  $P(X_{13} = 0, X_{22} = 1, X_{33} = 0)$  in terms of the conditional probabilities from the Bayesian Network.



# Bayesian Network question

1).  $X_{21} \perp\!\!\!\perp X_{33} \mid X_{11}, X_{12}$

by local semantics. Note that for ~~for~~ node  $X_{21}$ ,  $X_{11}$  and  $X_{12}$  are its parents and are observed, while  $X_{33}$  is its non-descendant.

2). For node  $X_{33}$ ,  $X_{22}$  is its parent and is observed. All the nodes in the graph except  $X_{22}$  are its non-descendants. Therefore, by local semantics,  $X_{33}$  ~~is~~ is conditionally independent of all nodes (except  $X_{22}$ ) given  $X_{22}$ .

3). All variables are binary. Number of parameters is:

$$\begin{array}{ccccccc} 1 & + & 1 & + & 1 & + & 4 & + & 2 & + & 4 & + & 4 & + & 2 & = & 19 \\ \uparrow & & \uparrow & & \uparrow & & \uparrow & & \uparrow & & \uparrow & & \uparrow & & \uparrow & & \\ X_{11} & & X_{12} & & X_{13} & & X_{21} & & X_{22} & & X_{31} & & X_{32} & & X_{33} \end{array}$$

factorized distribution:  $P(X_{11})P(X_{12})P(X_{13})P(X_{21} \mid X_{11}, X_{12}) \cdot P(X_{22} \mid X_{13}) \cdot$   
 $P(X_{31} \mid X_{21}, X_{22}) \cdot P(X_{32} \mid X_{21}, X_{22}) \cdot P(X_{33} \mid X_{22})$

$$\begin{aligned} 4) \quad P(X_{13}, X_{22}, X_{33}) &= \sum_{X_{11}, X_{12}, X_{21}, X_{31}, X_{32}} P(X_{11}, X_{12}, X_{13}, X_{21}, X_{22}, X_{31}, X_{32}, X_{33}) \\ &= \sum_{X_{11}, X_{12}, X_{21}, X_{31}, X_{32}} P(X_{11})P(X_{12})P(X_{13})P(X_{21} \mid X_{11}, X_{12})P(X_{22} \mid X_{13})P(X_{31} \mid X_{21}, X_{22}) \\ &\quad P(X_{32} \mid X_{21}, X_{22})P(X_{33} \mid X_{22}) \\ &= P(X_{13})P(X_{22} \mid X_{13})P(X_{33} \mid X_{22}) \cdot \sum_{X_{11}} P(X_{11}) \cdot \sum_{X_{12}} P(X_{12}) \cdot \sum_{X_{21}} P(X_{21} \mid X_{11}, X_{12}) \\ &\quad \sum_{X_{31}} P(X_{31} \mid X_{21}, X_{22}) \cdot \sum_{X_{32}} P(X_{32} \mid X_{21}, X_{22}) \\ &= P(X_{13})P(X_{22} \mid X_{13})P(X_{33} \mid X_{22}) \end{aligned}$$

$\therefore P(X_{13}=0, X_{22}=1, X_{33}=0) = P(X_{13}=0) \cdot P(X_{22}=1 \mid X_{13}=0) \cdot P(X_{33}=0 \mid X_{22}=1)$