

Lecture 18: HMM & Deep Learning

Winter 2018

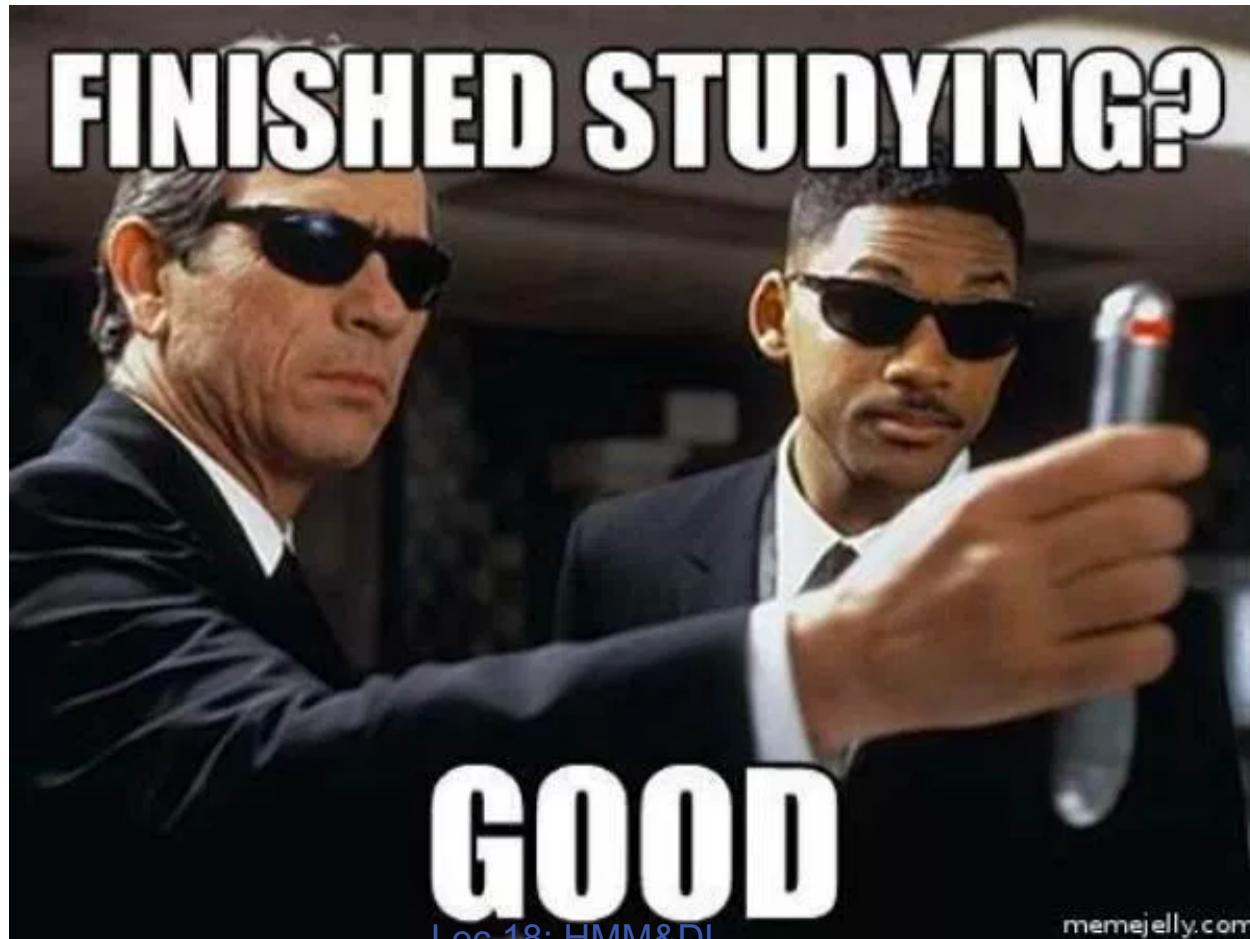
Kai-Wei Chang
CS @ UCLA

kw+cm146@kwchang.net

The instructor gratefully acknowledges Dan Roth, Vivek Srikumar, Sriram Sankararaman, Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

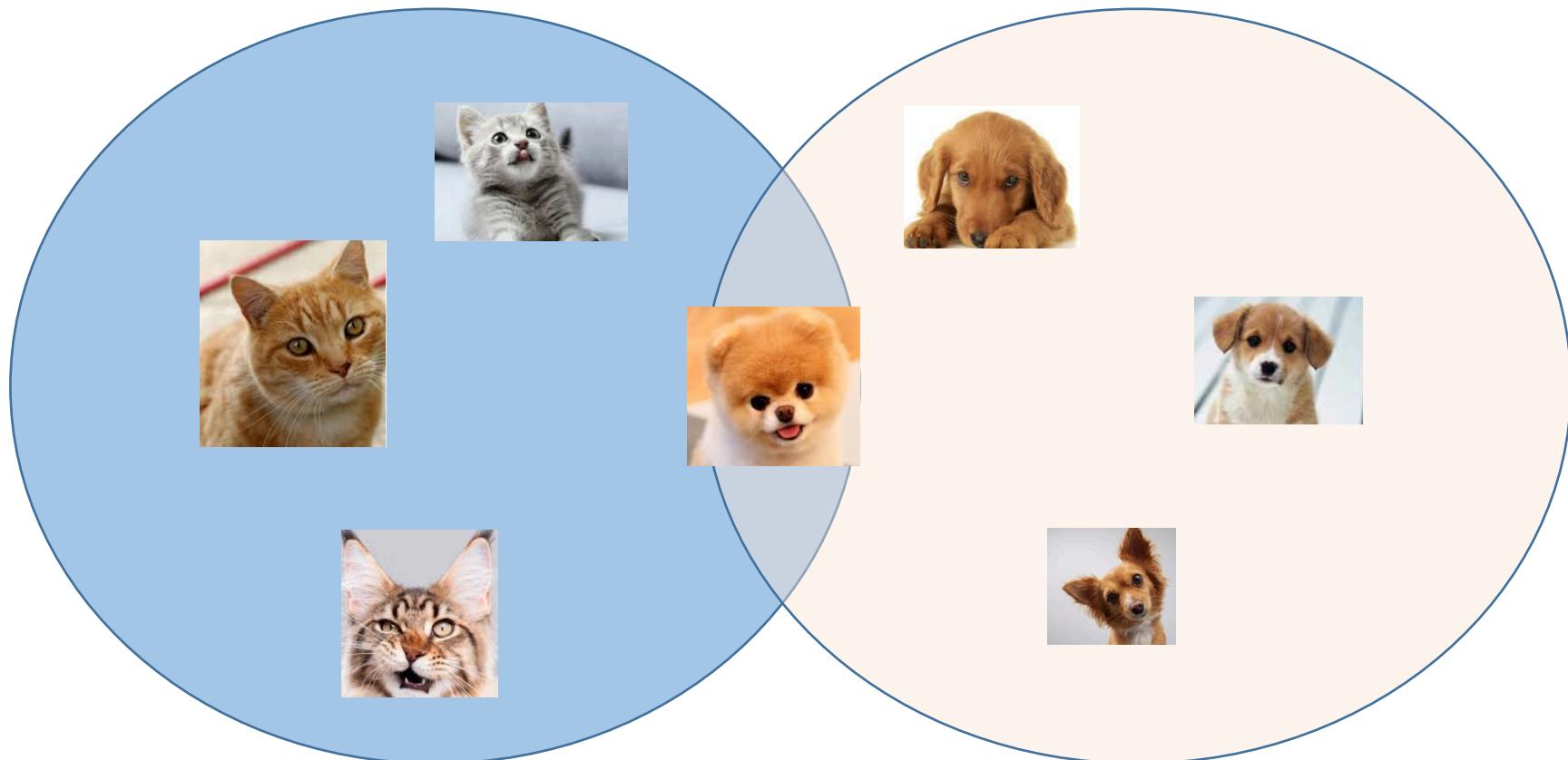
Final Exam 3/22, 11:30am-2:00pm

- ❖ Franz Hall: 1178 , 150min exam



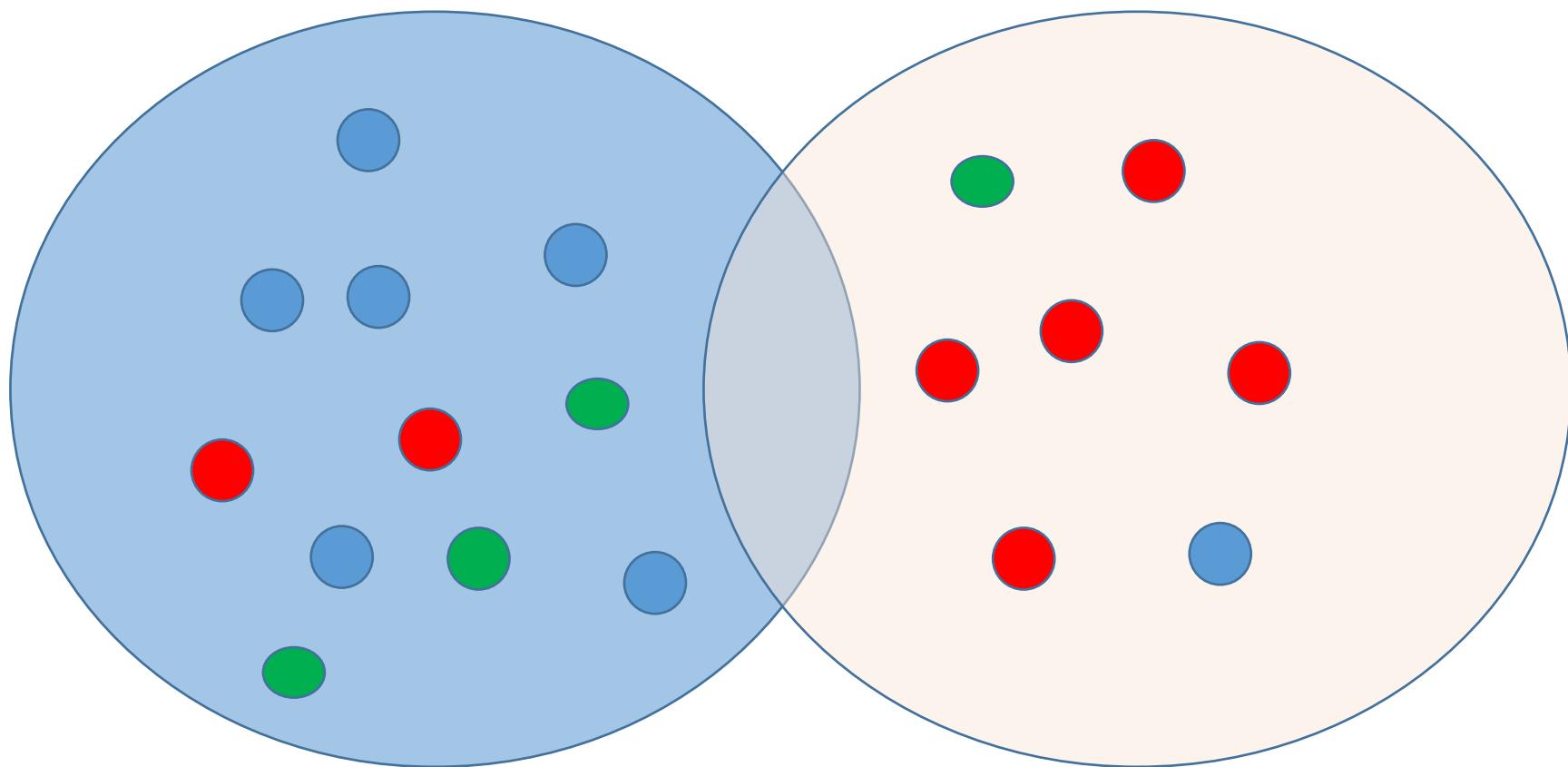
Generative model

- ❖ Idea: build a probabilistic model to measure how likely we would see the observations



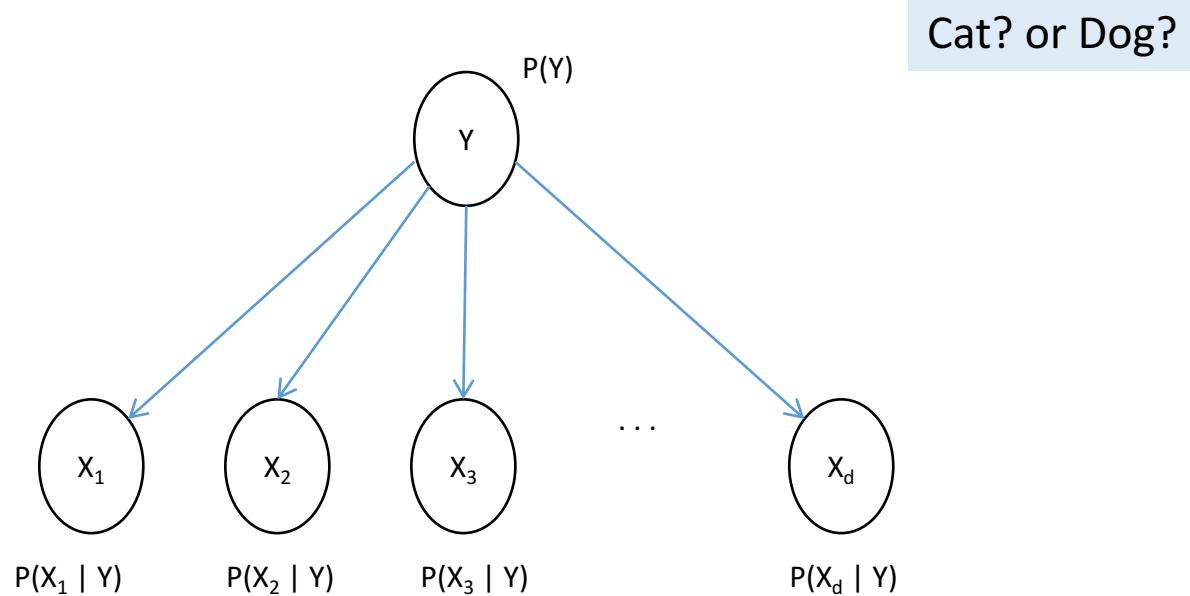
Generative model

- ❖ Idea: build a probabilistic model to measure how likely we would see the observations



Prediction under a generative Model

Given the label, sample the features independently from the conditional distributions



$$\operatorname{argmax}_y P(y) \prod_j P(x_j|y)$$

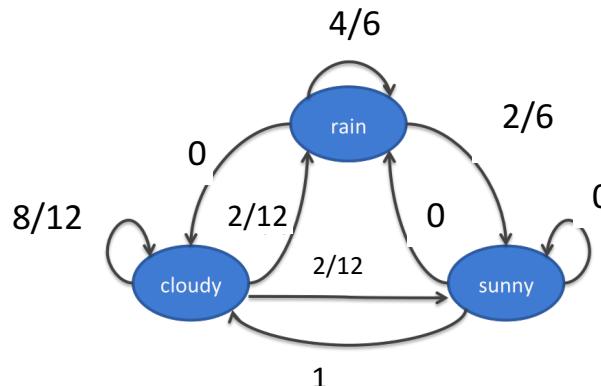
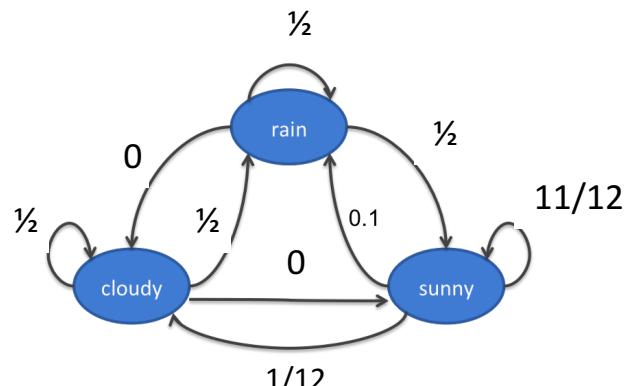
Weather Forecast by Markov Chain

- ❖ Observing a sequence of weather events, predict the weather tomorrow

S: Sunny, C: Cloudy, R: rainy

Los Angeles (winter): S S S S S S S S C C R R S S S S S S S S S S

Seattle (winter): C C C C C C R R C S C C S C R R R R C C C



Discrete Markov Process

- ❖ State at time t is x_t
- ❖ First-order Markov assumption

The state of the system at any time is ***independent*** of the full sequence history given the previous state

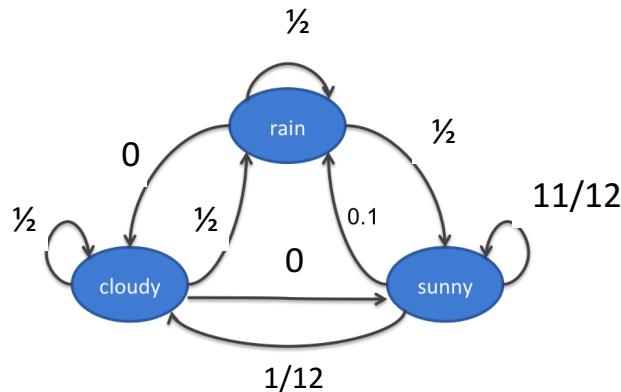
$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_{i-1})$$

- ❖ Defined by two sets of probabilities:
 - ❖ Initial state distribution: $P(x_1 = S_j)$
 - ❖ State transition probabilities: $P(x_i = S_j | x_{i-1} = S_k)$

Example: The weather

- ❖ Three states: rain, cloudy, sunny

State transitions:



- ❖ Observations are Markov chains:

Eg: *cloudy sunny sunny rain*

Probability of the sequence =

$P(\text{cloudy}) P(\text{sunny} | \text{cloudy}) P(\text{sunny} | \text{sunny}) P(\text{rain} | \text{sunny})$

Initial probability

Transition probabilities

Weather Forecast by Markov Chain

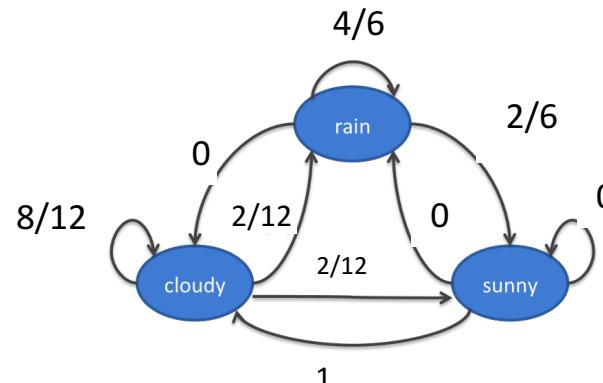
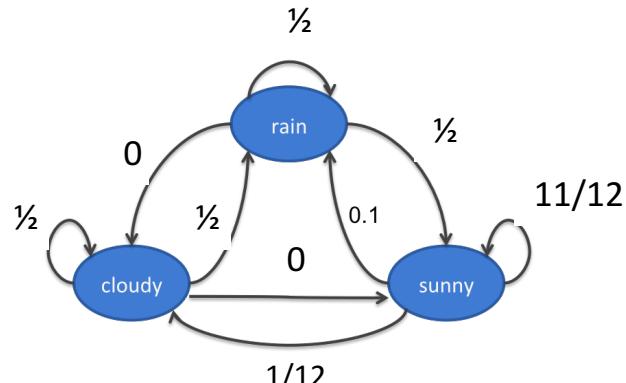
- ❖ Maximum Likelihood estimation:

$$P(y_t = \text{Rain} | y_{t-1} = \text{Sunny}) = \frac{\# S R}{\# S}$$

S: Sunny, C: Cloudy, R: rainy

Los Angeles (winter): S S S S S S S S C C R R S S S S S S S S S

Seattle (winter): C C C C C C R R C S C C S C R R R R C C C



Example: Another language model

It was a bright cold day in April

$$P(\text{It was a bright cold day in April}) =$$

$P(\text{It}) \times$ Probability of a word starting a sentence

$P(\text{was}|\text{It}) \times$ Probability of a word following “It”

$P(\text{a}|\text{was}) \times$ Probability of a word following “was”

$P(\text{bright}|\text{a}) \times$ Probability of a word following “a”

$P(\text{cold}|\text{bright}) \times$

$P(\text{day}|\text{cold}) \times \dots$

If there are K tokens/states, how many parameters do we need? $\mathcal{O}(K^2)$

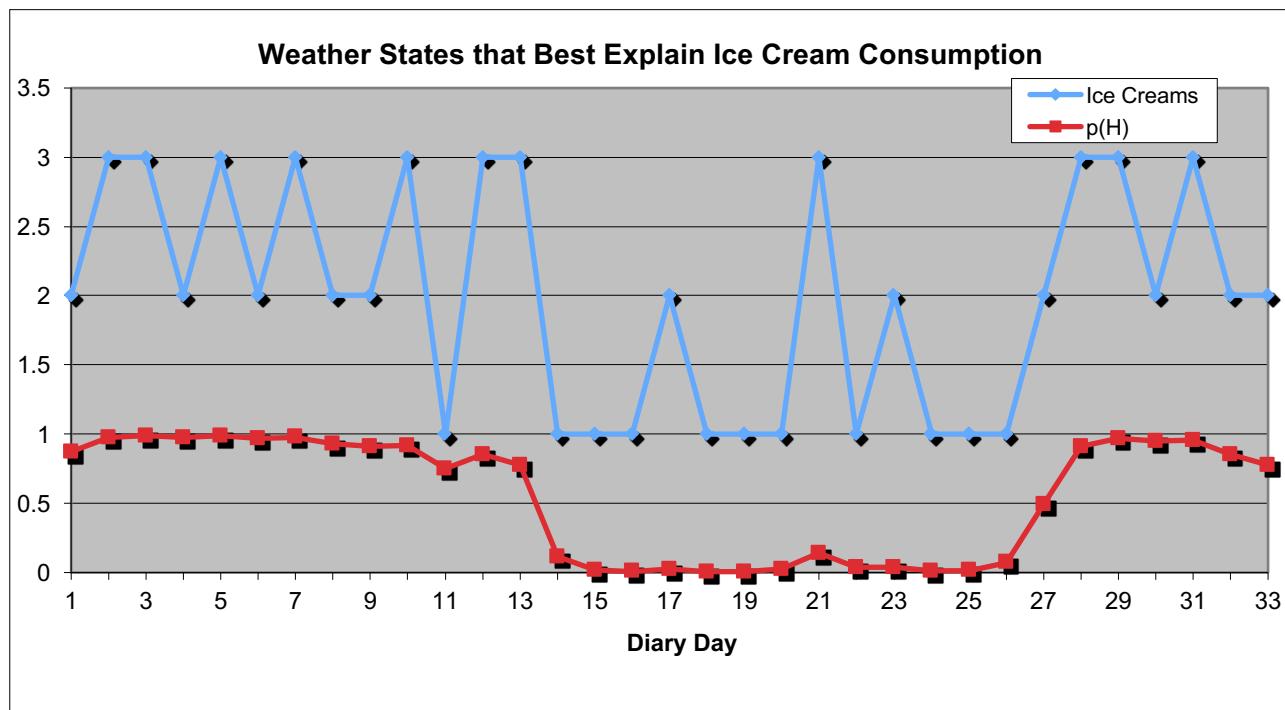
Hidden Markov Model



(C)old day v.s. (H)ot day

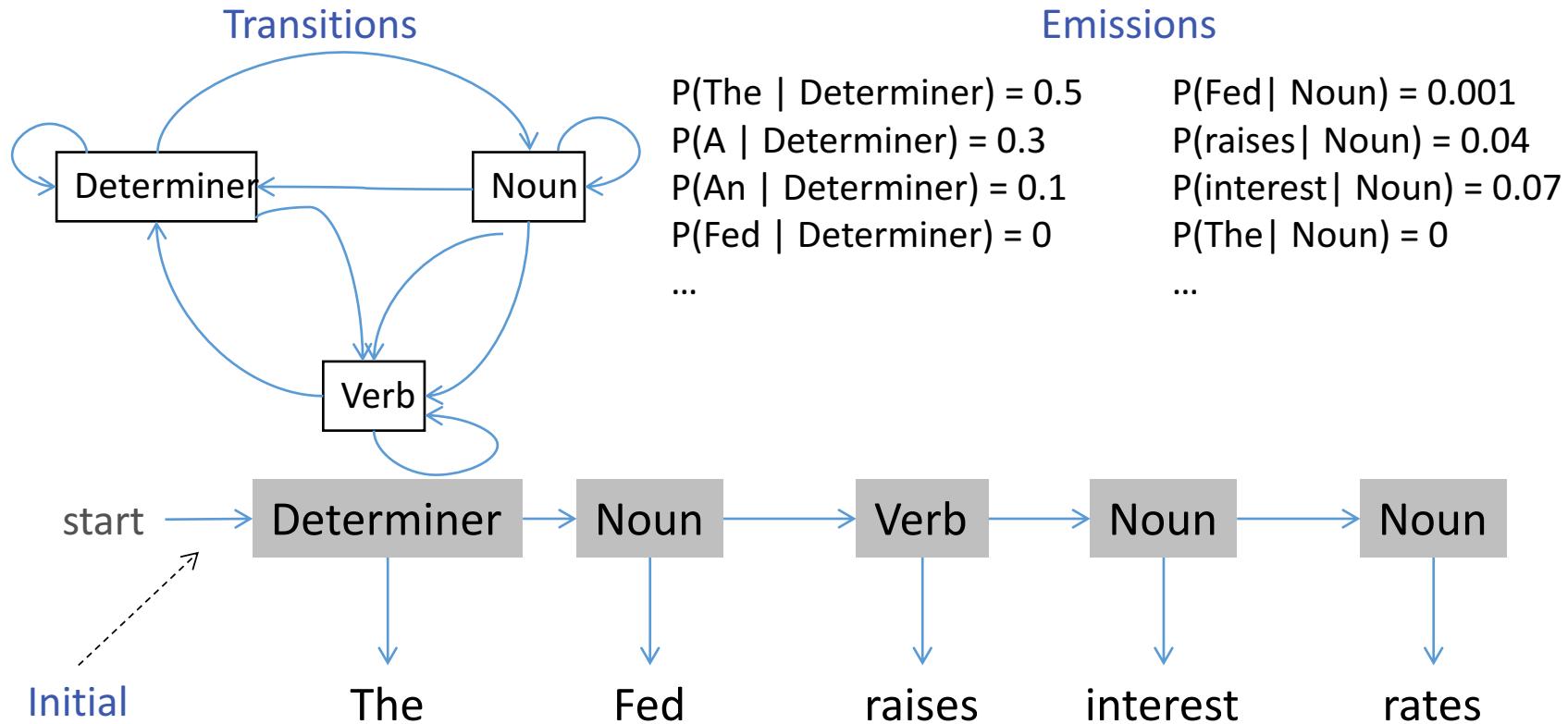
#cones

	$p(\dots C)$	$p(\dots H)$	$p(\dots \text{START})$
(1 ...)	0.7	0.1	
(2 ...)	0.2	0.2	
(3 ...)	0.1	0.7	
(C ...)	0.8	0.1	0.5
(H ...)	0.1	0.8	0.5
...)	0.1	0.1	0



Toy part-of-speech example

The Fed raises interest rates



Joint model over states and observations

- ❖ Notation
 - ❖ Number of states = K , Number of observations = M
 - ❖ π : Initial probability over states (K dimensional vector)
 - ❖ A : Transition probabilities ($K \times K$ matrix)
 - ❖ B : Emission probabilities ($K \times M$ matrix)
- ❖ Probability of states and observations
 - ❖ Denote states by y_1, y_2, \dots and observations by x_1, x_2, \dots

$$\begin{aligned} P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) &= P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i) \\ &= \pi_{y_1} \prod_{i=1}^{n-1} A_{y_i, y_{i+1}} \prod_{i=1}^n B_{y_i, x_i} \end{aligned}$$

Key questions for HMMs

[Rabiner 1999]

- Given an observation sequence, x_1, x_2, \dots, x_n and a model (π, A, B) , how to efficiently calculate the most probable state sequence?

Inference

- How to calculate (π, A, B) from observations?

Learning

Outline

- ❖ *Sequence models*
- ❖ Hidden Markov models
 - ❖ **Inference with HMM**
 - ❖ Supervised Learning for HMM

Most likely state sequence

- ❖ Input:
 - ❖ A hidden Markov model (π, A, B)
 - ❖ An observation sequence $\mathbf{x} = (x_1, x_2, \dots, x_n)$
- ❖ Output: A state sequence $\mathbf{y} = (y_1, y_2, \dots, y_n)$ that corresponds to
 - ❖ Maximum *a posteriori* inference (MAP inference)
$$\arg \max_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}, \pi, A, B)$$
- ❖ Computationally: combinatorial optimization

MAP inference

- ❖ We want

$$\arg \max_{\mathbf{y}} P(\mathbf{y} | \mathbf{x}, \pi, A, B)$$

- ❖ We have defined

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1} | y_i) \prod_{i=1}^n P(x_i | y_i)$$

- ❖ But $P(\mathbf{y} | \mathbf{x}, \pi, A, B) \propto P(\mathbf{x}, \mathbf{y} | \pi, A, B)$
 - ❖ And we don't care about $P(\mathbf{x})$ we are maximizing over \mathbf{y}

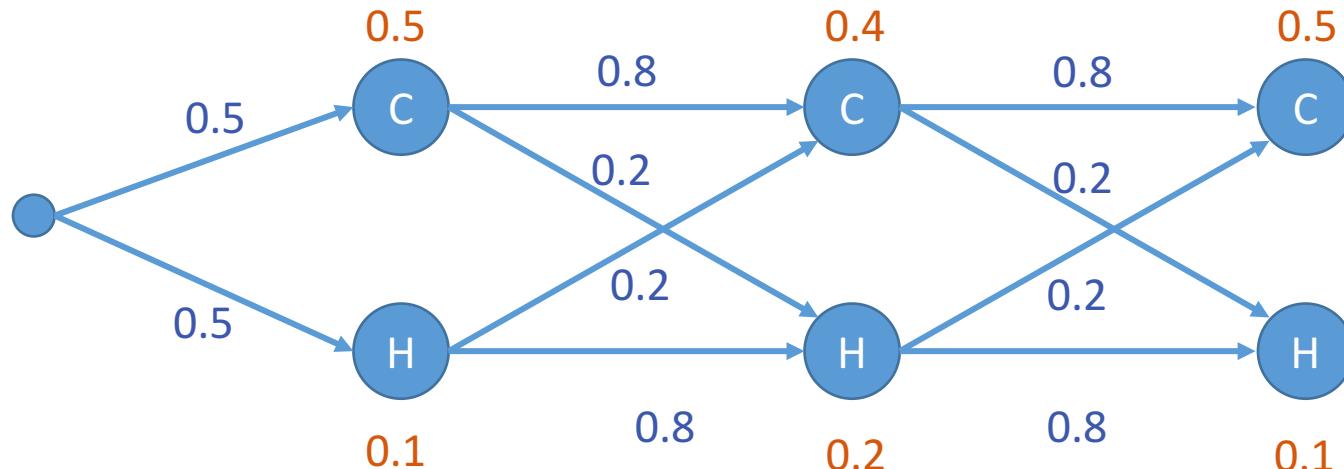
- ❖ So, $\arg \max_{\mathbf{y}} P(\mathbf{y} | \mathbf{x}, \pi, A, B) = \arg \max_{\mathbf{y}} P(\mathbf{y}, \mathbf{x} | \pi, A, B)$

Jason's ice cream

#cones

	$p(\dots C)$	$p(\dots H)$	$p(\dots \text{START})$
(1 ...)	0.5	0.1	
(2 ...)	0.4	0.2	
(3 ...)	0.1	0.7	
(C ...)	0.8	0.2	0.5
(H ...)	0.2	0.8	0.5

❖ Best tag sequence for $P("1,2,1")?$

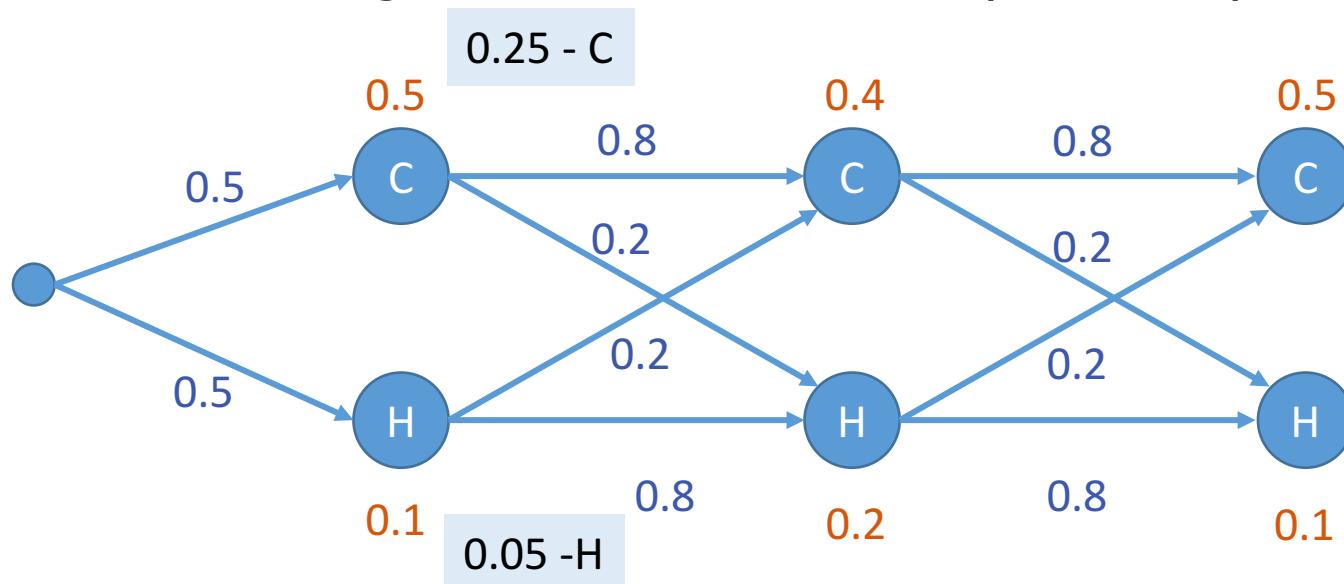


Jason's ice cream

#cones

	$p(\dots C)$	$p(\dots H)$	$p(\dots \text{START})$
(1 ...)	0.5	0.1	
(2 ...)	0.4	0.2	
(3 ...)	0.1	0.7	
(C ...)	0.8	0.2	0.5
(H ...)	0.2	0.8	0.5

❖ Best tag sequence for $P("1,2,1")?$

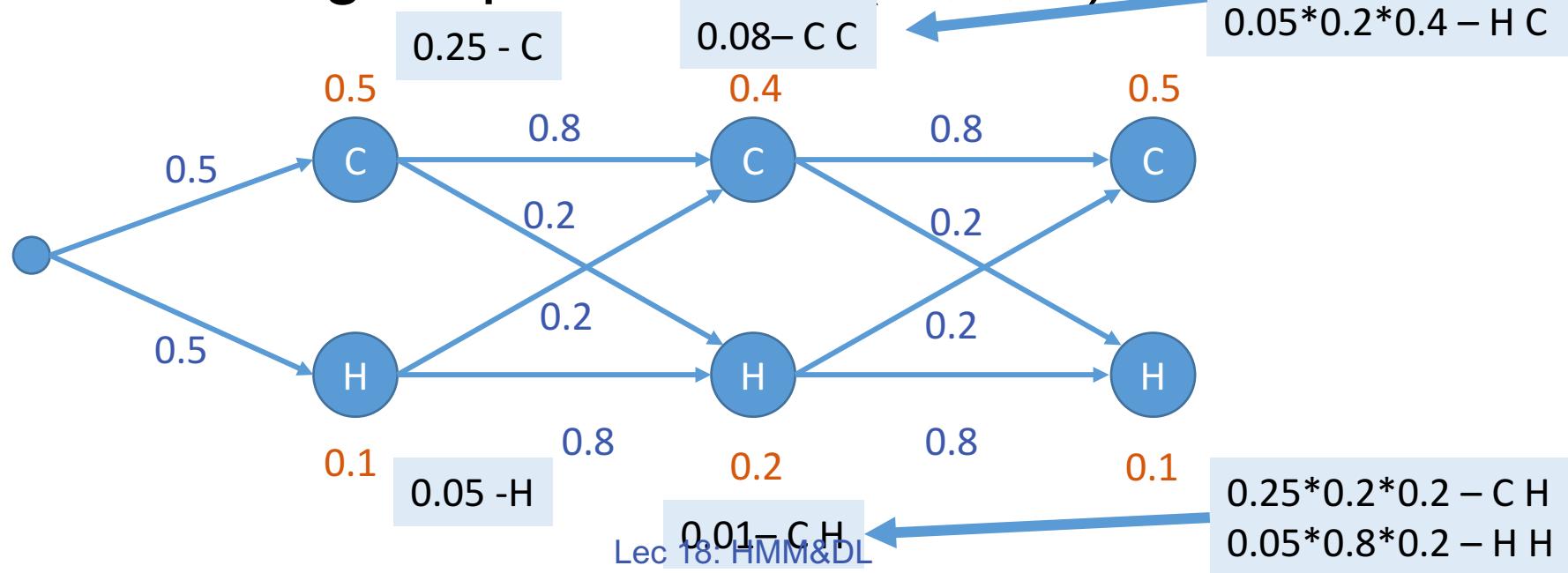


Jason's ice cream

#cones

	$p(\dots C)$	$p(\dots H)$	$p(\dots \text{START})$
(1 ...)	0.5	0.1	
(2 ...)	0.4	0.2	
(3 ...)	0.1	0.7	
(C ...)	0.8	0.2	0.5
(H ...)	0.2	0.8	0.5

❖ Best tag sequence for $P("1,2,1")?$

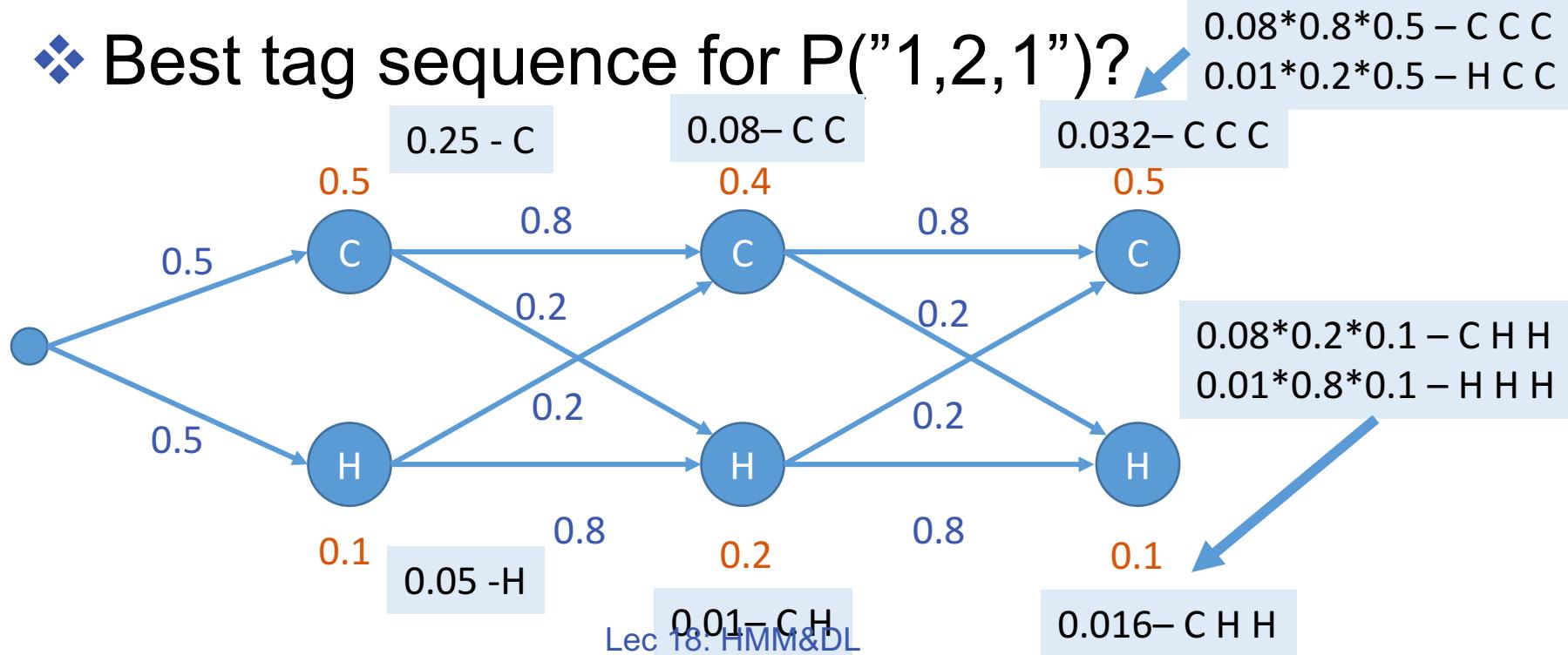


Jason's ice cream

#cones

	$p(\dots C)$	$p(\dots H)$	$p(\dots \text{START})$
(1 ...)	0.5	0.1	
(2 ...)	0.4	0.2	
(3 ...)	0.1	0.7	
(C ...)	0.8	0.2	0.5
(H ...)	0.2	0.8	0.5

❖ Best tag sequence for $P("1,2,1")?$

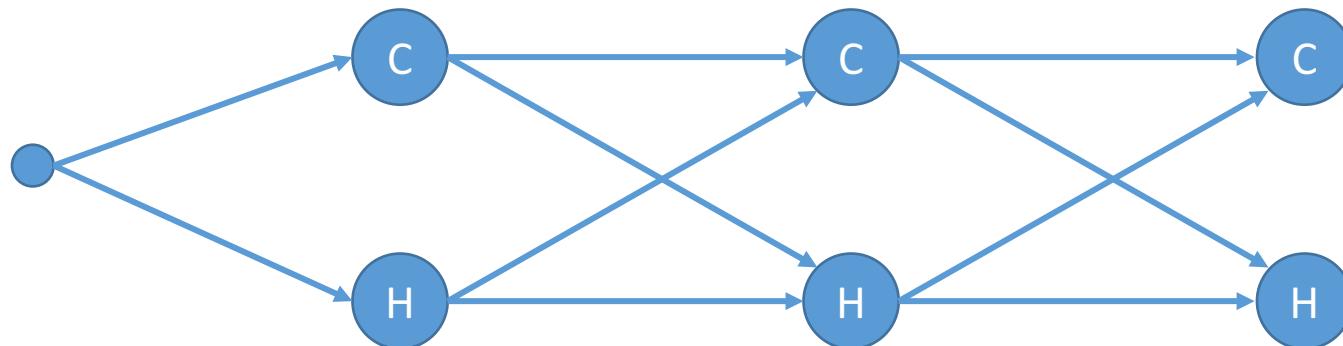


Your turn

#cones

	$p(\dots C)$	$p(\dots H)$	$p(\dots \text{START})$
(1 ...)	0.5	0.1	
(2 ...)	0.4	0.2	
(3 ...)	0.1	0.7	
(C ...)	0.8	0.2	0.5
(H ...)	0.2	0.8	0.5

❖ Best tag sequence for P("3,3,2")?



Viterbi algorithm

Max-product algorithm for first order sequences

π : Initial probabilities
A: Transitions
B: Emissions

1. **Initial:** For each state s , calculate

$$\text{score}_1(s) = P(s)P(x_1|s) = \pi_s B_{x_1,s}$$

1. **Recurrence:** For $i = 2$ to n , for every state s , calculate

$$\begin{aligned}\text{score}_i(s) &= \max_{y_{i-1}} P(s|y_{i-1})P(x_i|s)\text{score}_{i-1}(y_{i-1}) \\ &= \max_{y_{i-1}} A_{y_{i-1},s} B_{s,x_i} \text{score}_{i-1}(y_{i-1})\end{aligned}$$

1. **Final state:** calculate

$$\max_{\mathbf{y}} P(\mathbf{y}, \mathbf{x} | \pi, A, B) = \max_s \text{score}_n(s)$$

This only calculates the max. To get final answer (*argmax*),

- keep track of which state corresponds to the max at each step
- build the answer using these back pointers

General idea

- ❖ Dynamic programming
 - ❖ The best solution for the full problem relies on best solution to sub-problems
 - ❖ Memoize partial computation
- ❖ Examples
 - ❖ Viterbi algorithm
 - ❖ Dijkstra's shortest path algorithm
 - ❖ ...

Complexity of inference

- ❖ Complexity parameters
 - ❖ Input sequence length: n
 - ❖ Number of states: K
- ❖ Memory
 - ❖ Storing the table: nK (scores for all states at each position)
- ❖ Runtime
 - ❖ At each step, go over pairs of states
 - ❖ $O(nK^2)$

Outline

- ❖ *Sequence models*
- ❖ Hidden Markov models
 - ❖ Inference with HMM
 - ❖ Supervised Learning for HMM

Learning HMM parameters

- ❖ Assume we know the number of states in the HMM
- ❖ Two possible scenarios

1. We are given a data set $D = \{\langle x_i, y_i \rangle\}$ of sequences labeled with states

And we have to learn the parameters of the HMM (π, A, B)

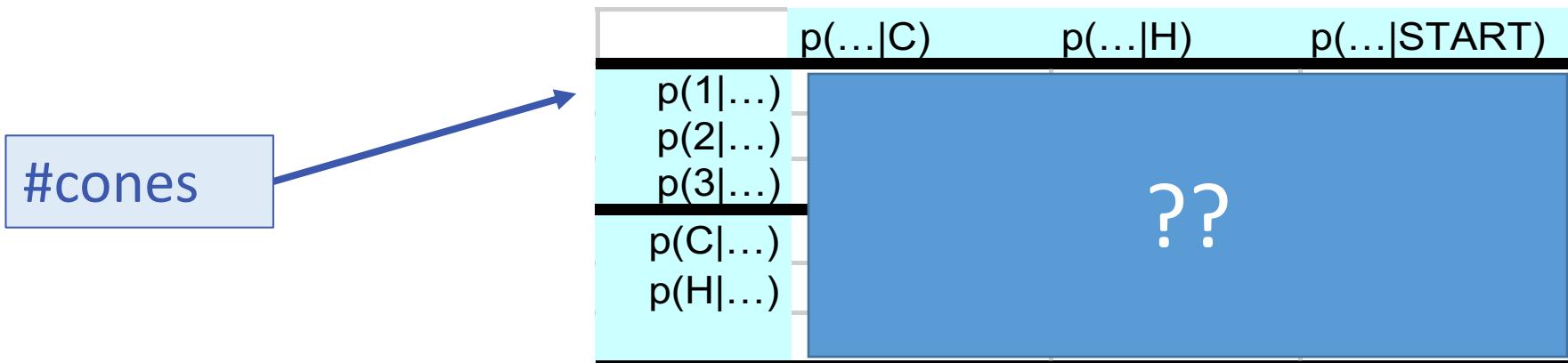
Supervised learning with complete data

2. We are given only a collection of sequences $D = \{x_i\}$

And we have to learn the parameters of the HMM (π, A, B)

Unsupervised learning, with incomplete data

Jason's ice cream



- ❖ Can we figure out (π, A, B) from history records?

#cones	1	2	1	1	2	3	3	3	3	2	3	1	2
Hot/Cold	C	C	C	C	H	H	H	H	C	C	H	H	H

#cones	3	2	2	2	2	3	1	1	1	2	3	2	2
Hot/Cold	H	H	H	C	C	H	H	C	C	H	H	H	C

Supervised learning of HMM

- ❖ We are given a dataset $D = \{\langle \mathbf{x}_i, \mathbf{y}_i \rangle\}$
 - ❖ Each \mathbf{x}_i is a sequence of observations and \mathbf{y}_i is a sequence of states that correspond to \mathbf{x}_i
- Goal: Learn initial, transition, emission distributions (π, A, B)
- ❖ How do we learn the parameters of the probability distribution?
 - ❖ The maximum likelihood principle

Where have we seen this before?

$$(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B) = \max_{\pi, A, B} \prod_i P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

And we know how to write this in terms of the parameters of the HMM

Supervised learning details

$$(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B) = \max_{\pi, A, B} \prod_i P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

π, A, B can be estimated separately just by counting

- ❖ Makes learning simple and fast

[Exercise: Derive the following using derivatives of the log likelihood.]

$$\pi_s = \frac{\text{count}(\text{start} \rightarrow s)}{n}$$

Initial probabilities

Number of instances where the first state is s

Number of examples

$$A_{s',s} = \frac{\text{count}(s \rightarrow s')}{\text{count}(s)}$$

Transition probabilities

$$B_{s,x} = \frac{\text{count} \begin{pmatrix} s \\ \downarrow \\ x \end{pmatrix}}{\text{count}(s)}$$

Emission probabilities

Priors and smoothing

- ❖ Maximum likelihood estimation works best with lots of annotated data
 - ❖ Never the case
- ❖ Priors inject information about the probability distributions
- ❖ Effectively additive smoothing
 - ❖ Add small constants to the counts

Question we will not answer in this lecture

- ❖ Can we figure out (π, A, B) from just observation

#cones

	$p(\dots C)$	$p(\dots H)$	$p(\dots \text{START})$
$p(1 \dots)$			
$p(2 \dots)$			
$p(3 \dots)$			
$p(C \dots)$??
$p(H \dots)$			

#cones	1	2	1	1	2	3	3	3	3	2	3	1	2
Hot/Cold							??						

#cones	3	2	2	2	2	3	1	1	1	2	3	2	2
Hot/Cold							??						

Hidden Markov Models summary

- ❖ Predicting sequences
 - ❖ As many output states as observations
- ❖ Markov assumption helps decompose the score
- ❖ Several algorithmic questions
 - ❖ Most likely state (decoding)
 - ❖ Learning parameters
 - ❖ Supervised, Unsupervised

Neural Network

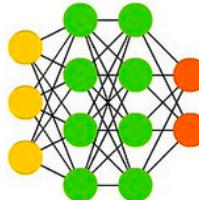
A mostly complete chart of

Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool

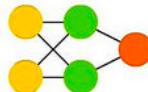
Deep Feed Forward (DFF)



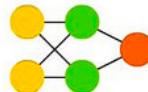
Perceptron (P)



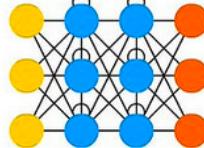
Feed Forward (FF)



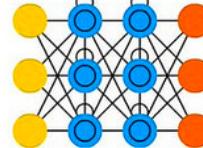
Radial Basis Network (RBF)



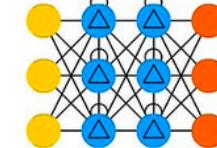
Recurrent Neural Network (RNN)



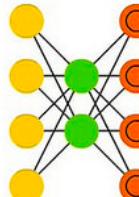
Long / Short Term Memory (LSTM)



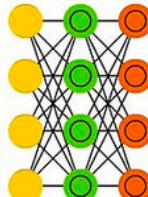
Gated Recurrent Unit (GRU)



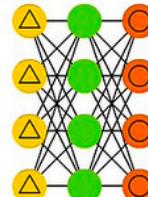
Auto Encoder (AE)



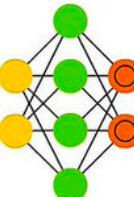
Variational AE (VAE)



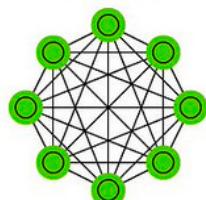
Denoising AE (DAE)



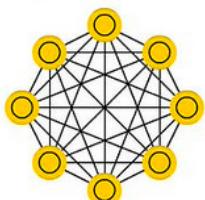
Sparse AE (SAE)



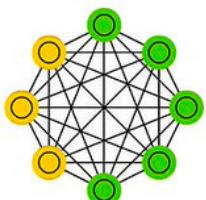
Markov Chain (MC)



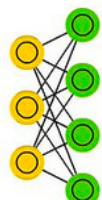
Hopfield Network (HN)



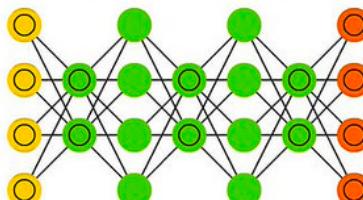
Boltzmann Machine (BM)



Restricted BM (RBM)



Deep Belief Network (DBN)

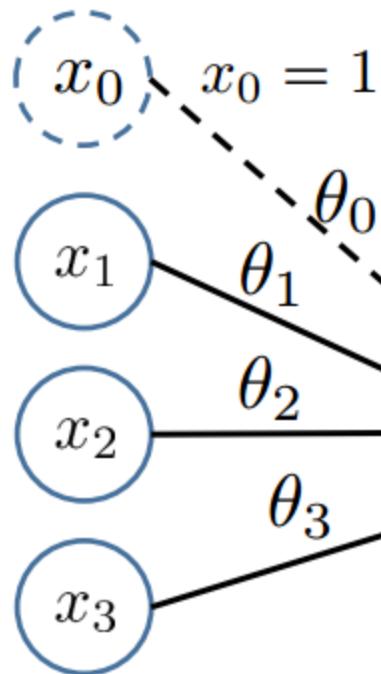


This lecture

- ❖ Neural Network
- ❖ Recurrent Neural Network

Neural Network

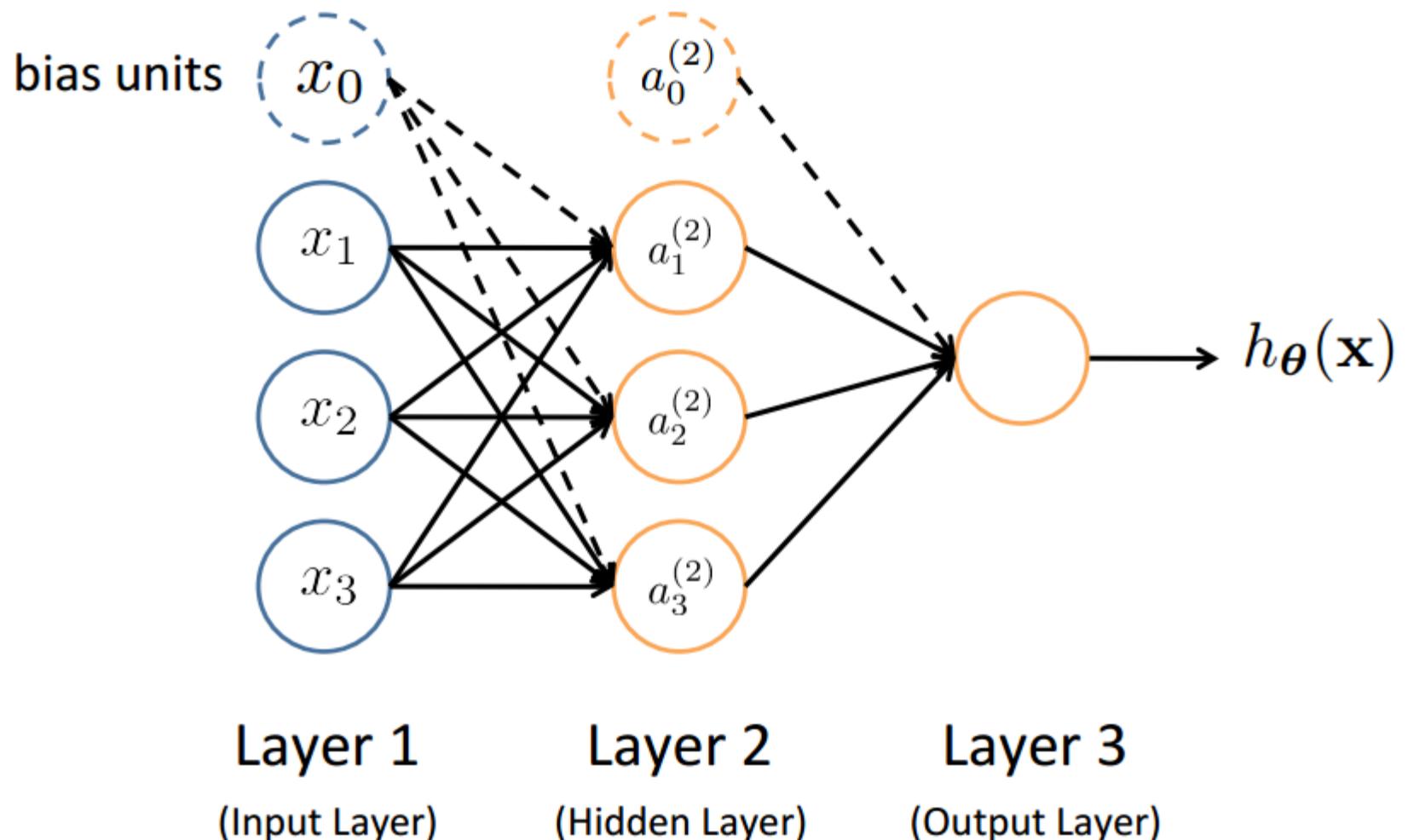
“bias unit”



$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$
$$h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x})$$
$$= \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

Sigmoid (logistic) activation function: $g(z) = \frac{1}{1 + e^{-z}}$

Neural Network (feed forward)



Feed-Forward Process

- ❖ Input layer units are features (e.g., words)
 - ❖ Usually, one-hot vector or word embedding
- ❖ Working forward through the network, the **input function** is applied to compute the input value
 - ❖ E.g., weighted sum of the input
- ❖ The **activation function** transforms this input function into a final value
 - ❖ Typically a **nonlinear** function (e.g, **sigmoid**)

Activation functions

Also called transfer functions

$$\text{output} = \text{activation}(\mathbf{w}^T \mathbf{x} + b)$$

Name of the neuron	Activation function: $\text{activation}(z)$
Linear unit	z
Threshold/sign unit	$\text{sign}(z)$
Sigmoid unit	$\frac{1}{1 + \exp(-z)}$
Rectified linear unit (ReLU)	$\max(0, z)$
Tanh unit	$\tanh(z)$

Many more activation functions exist (sinusoid, sinc, gaussian, polynomial...)

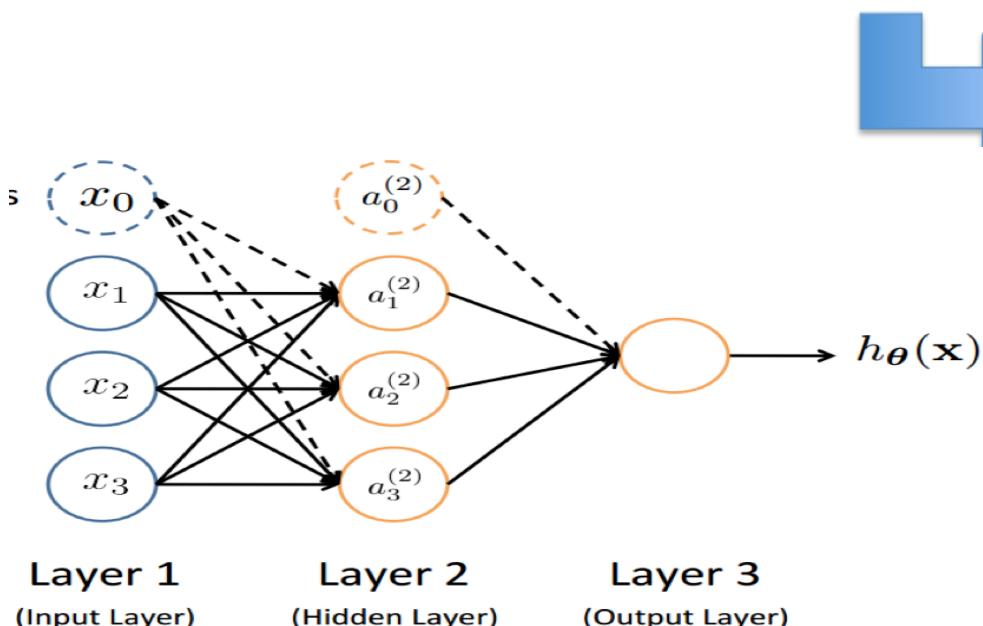
Vector Representation

$$a_1^{(2)} = g \left(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3 \right) = g \left(z_1^{(2)} \right)$$

$$a_2^{(2)} = g \left(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3 \right) = g \left(z_2^{(2)} \right)$$

$$a_3^{(2)} = g \left(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3 \right) = g \left(z_3^{(2)} \right)$$

$$h_{\Theta}(\mathbf{x}) = g \left(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)} \right) = g \left(z_1^{(3)} \right)$$



Feed-Forward Steps:

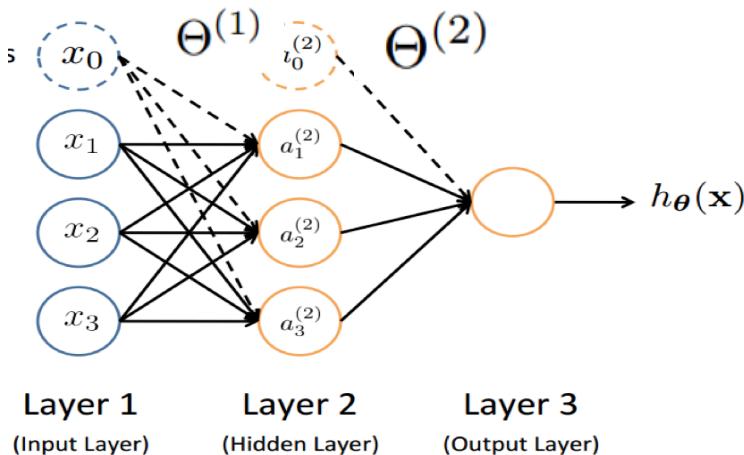
$$\mathbf{z}^{(2)} = \Theta^{(1)} \mathbf{x}$$

$$\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$$

$$\text{Add } a_0^{(2)} = 1$$

$$\mathbf{z}^{(3)} = \Theta^{(2)} \mathbf{a}^{(2)}$$

$$h_{\Theta}(\mathbf{x}) = \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$$



$a_i^{(j)}$ = “activation” of unit i in layer j
 $\Theta^{(j)}$ = weight matrix controlling function
 mapping from layer j to layer $j + 1$

$$a_1^{(2)} = g(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3)$$

$$h_\Theta(x) = a_1^{(3)} = g(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)})$$

If network has s_j units in layer j and s_{j+1} units in layer $j+1$,
 then $\Theta^{(j)}$ has dimension $s_{j+1} \times (s_j + 1)$.

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4} \quad \Theta^{(2)} \in \mathbb{R}^{1 \times 4}$$

Can extend to multi-class



Pedestrian



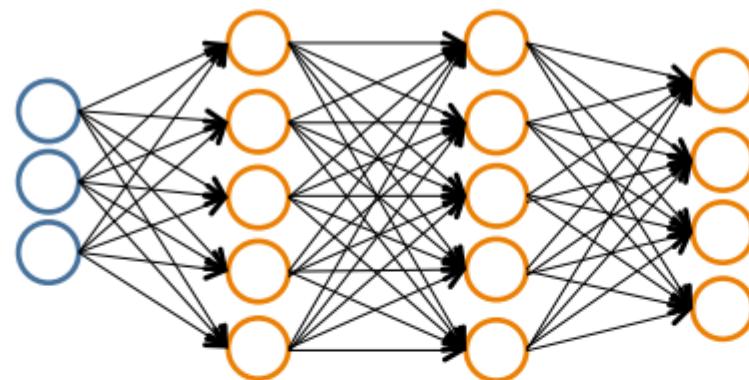
Car



Motorcycle



Truck



$$h_{\Theta}(\mathbf{x}) \in \mathbb{R}^K$$

We want:

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

when pedestrian

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

when car

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

when motorcycle

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

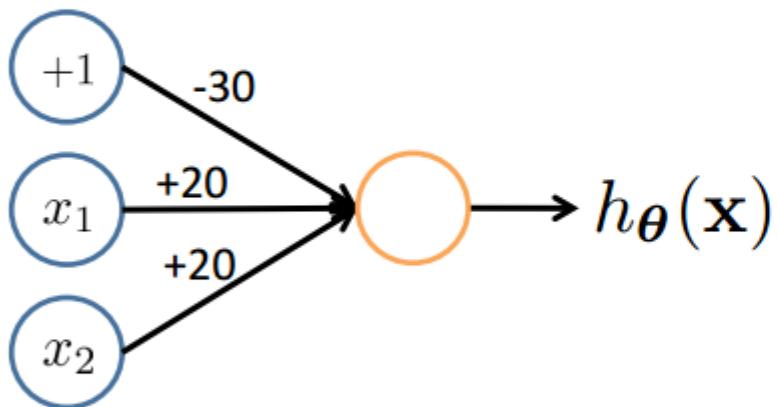
when truck

Why staged predictions?

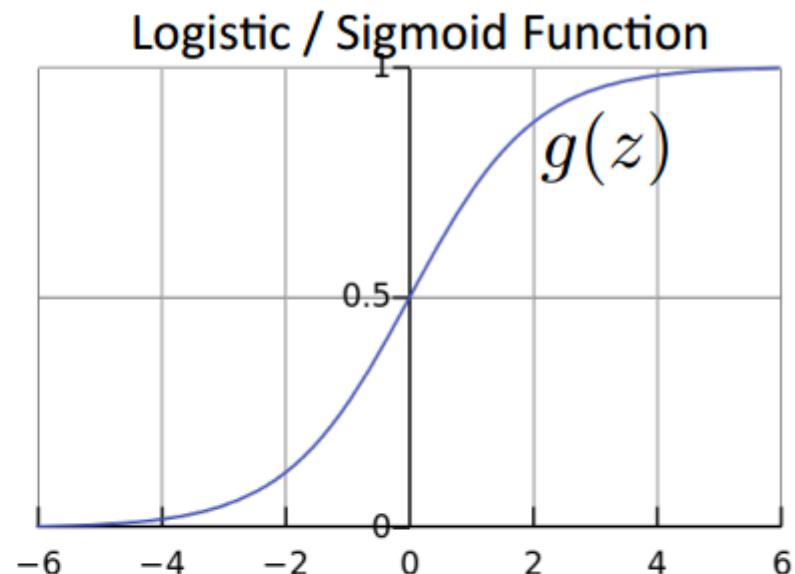
Simple example: AND

$$x_1, x_2 \in \{0, 1\}$$

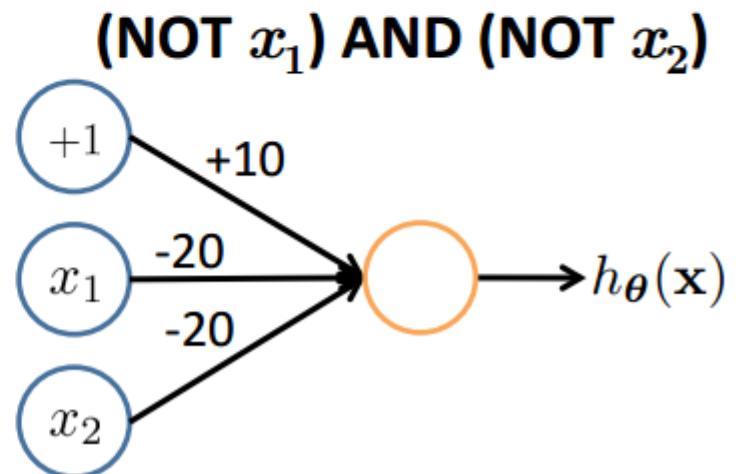
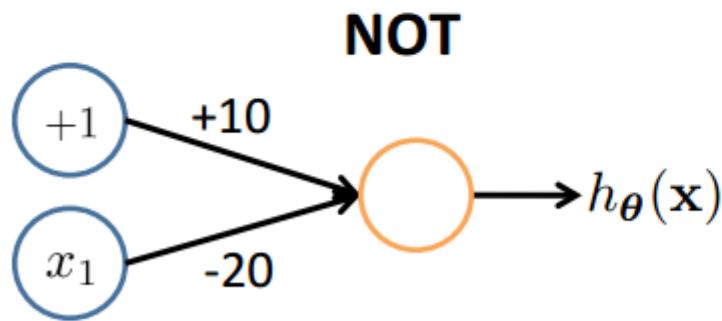
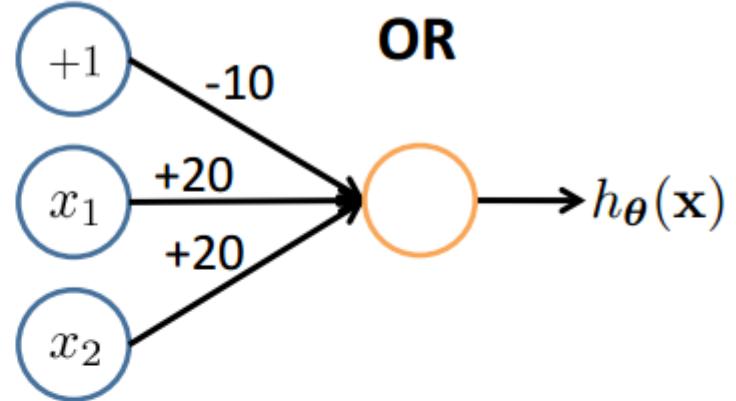
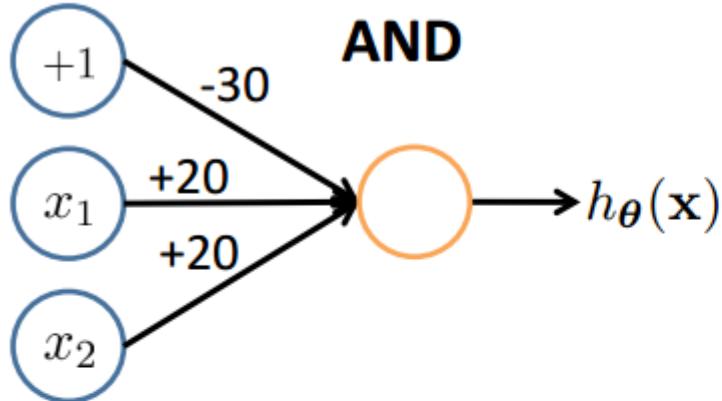
$$y = x_1 \text{ AND } x_2$$



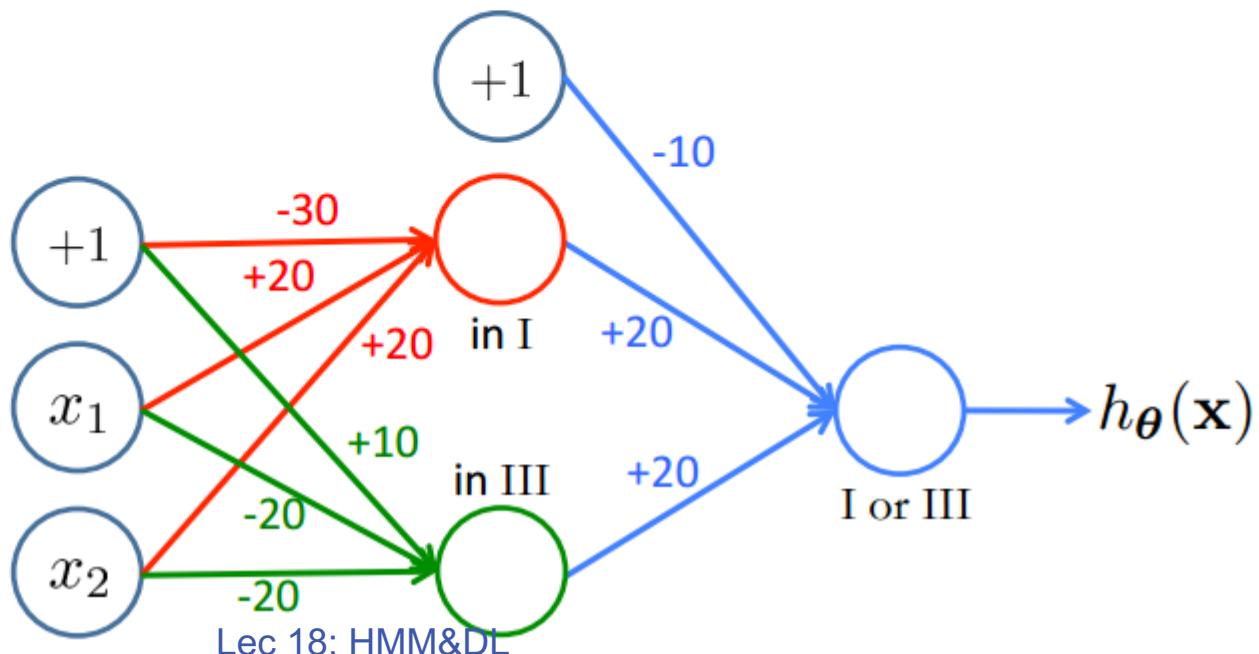
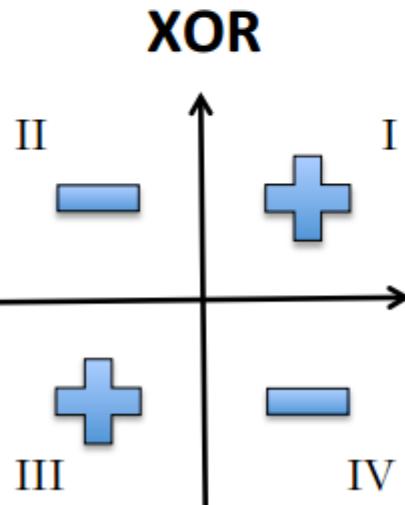
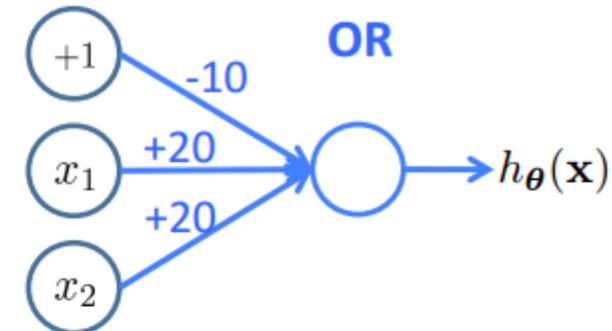
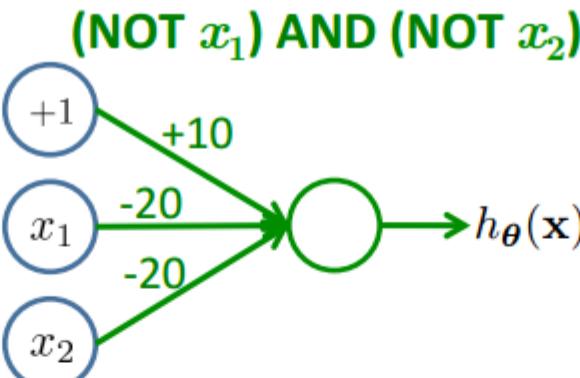
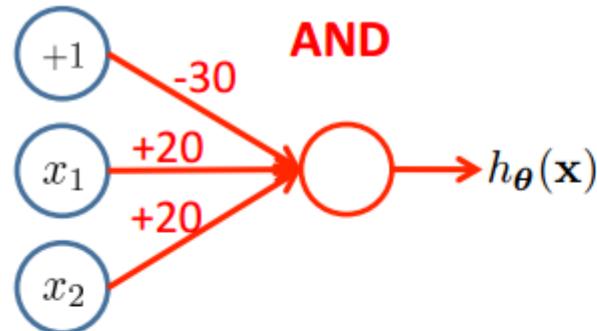
$$h_{\theta}(\mathbf{x}) = g(-30 + 20x_1 + 20x_2)$$



x_1	x_2	$h_{\theta}(\mathbf{x})$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

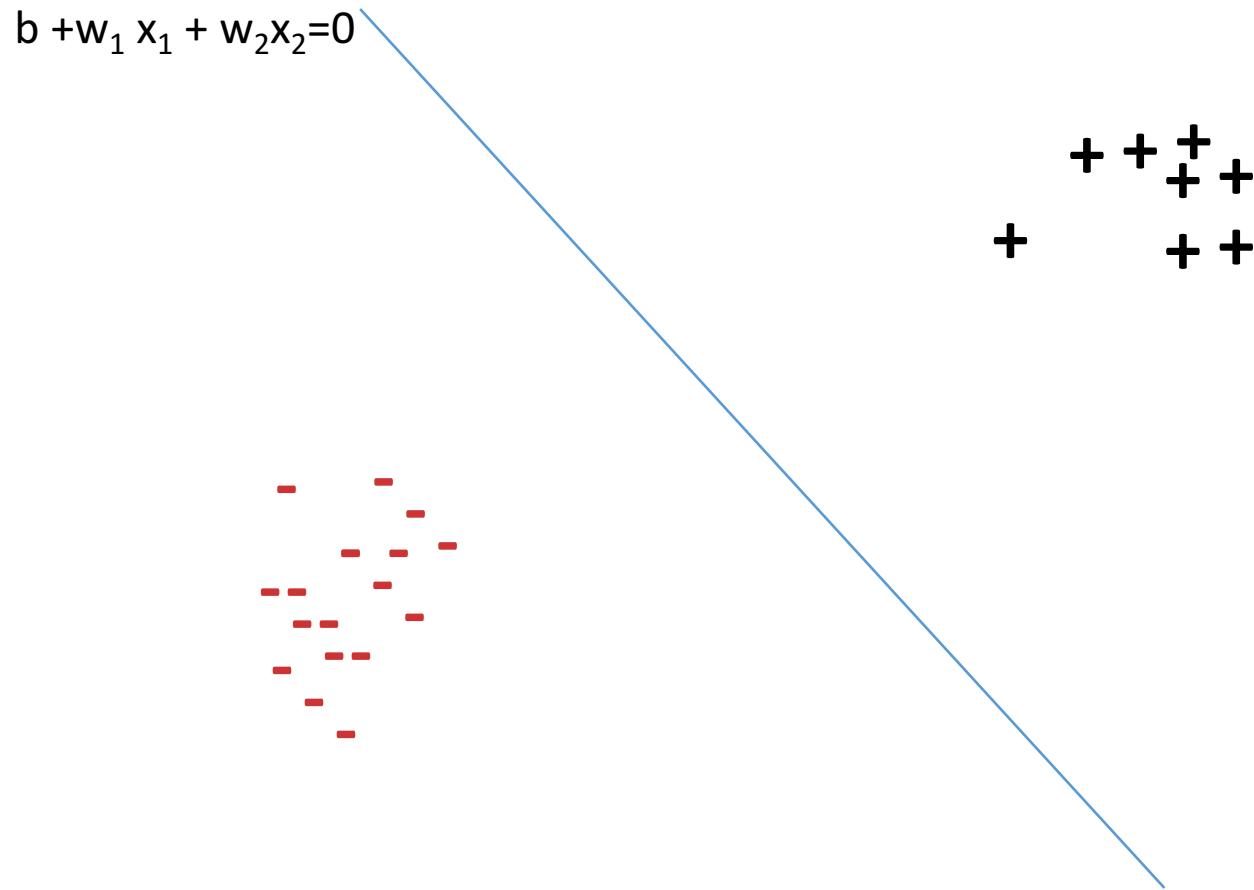


Combining Representations to Create Non-Linear Functions

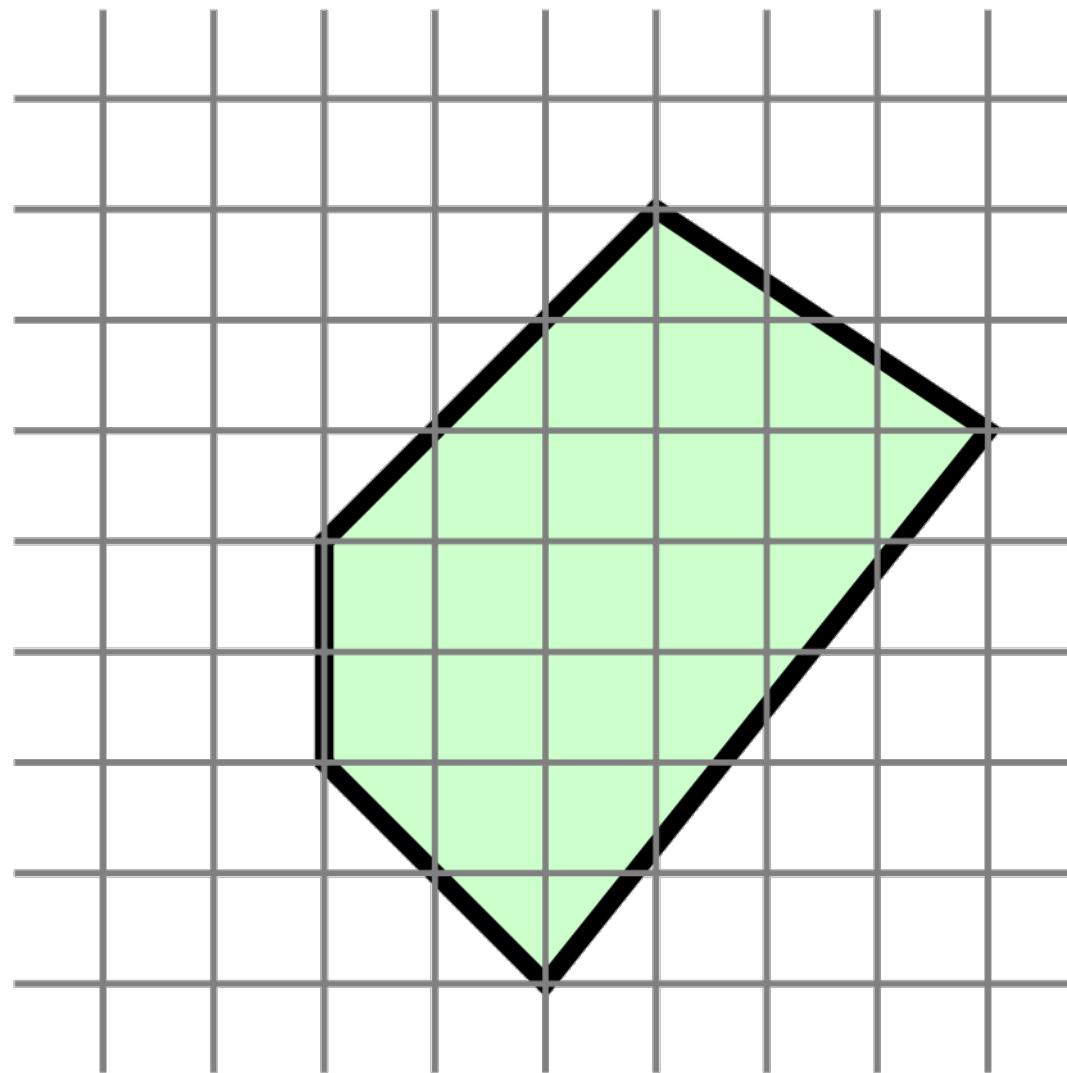


A single neuron with threshold activation

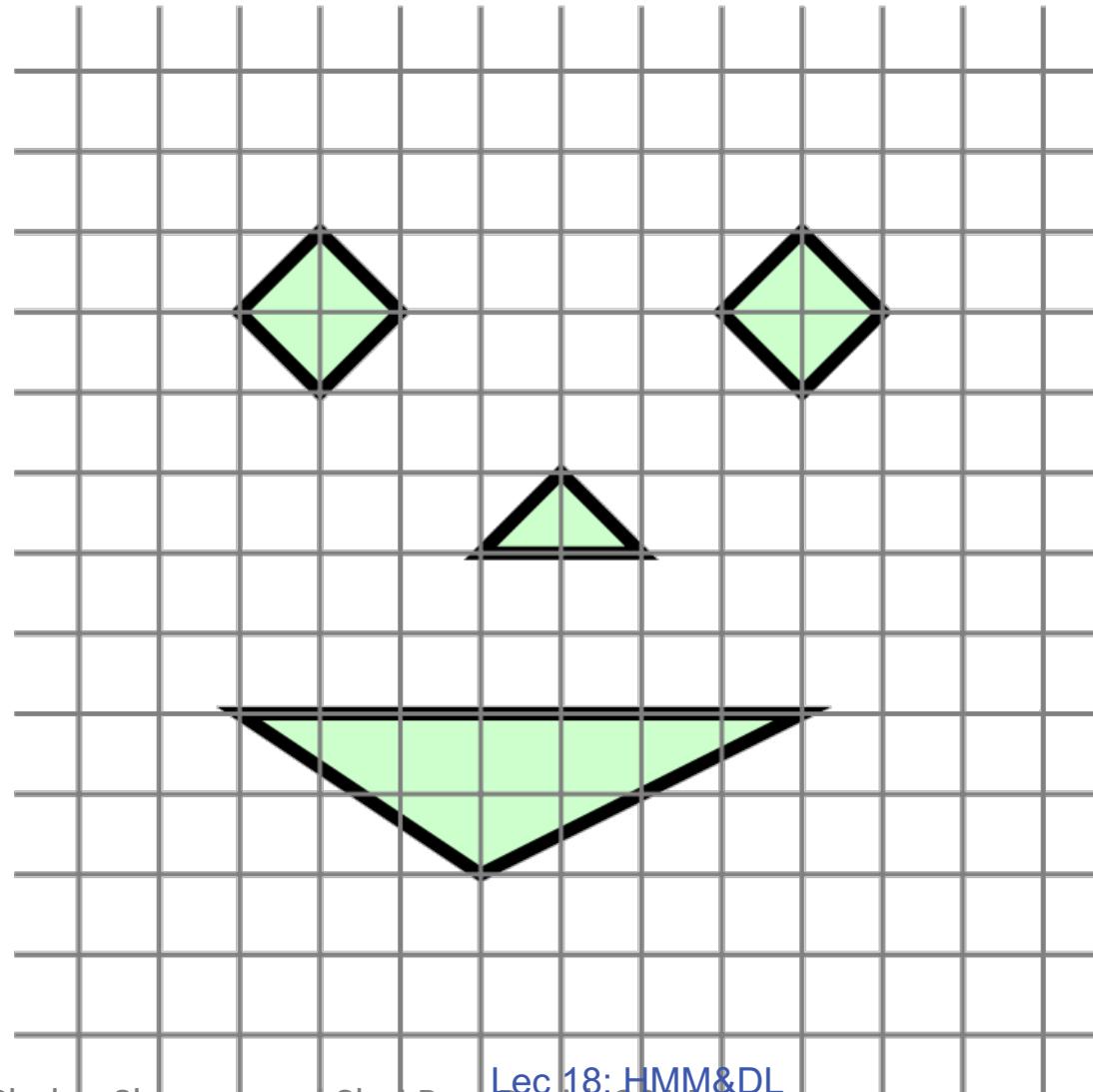
$$\text{Prediction} = \text{sgn}(b + w_1 x_1 + w_2 x_2)$$



Two layers, with threshold activations



Three layers with threshold activations

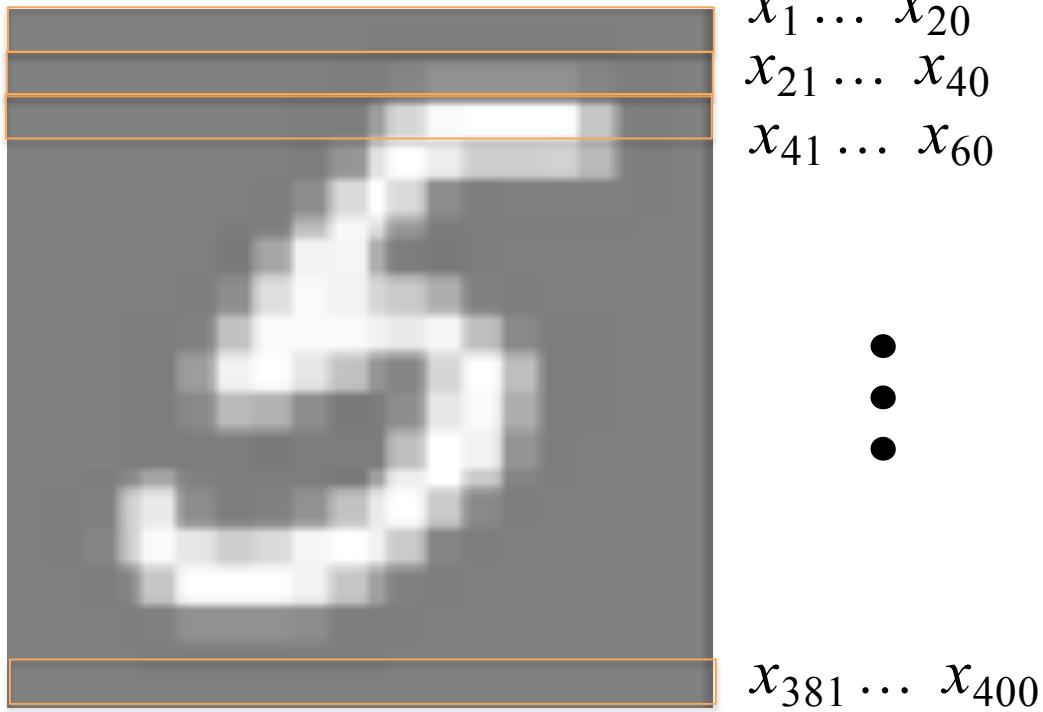


In general, unions
of convex polygons

Layering Representations



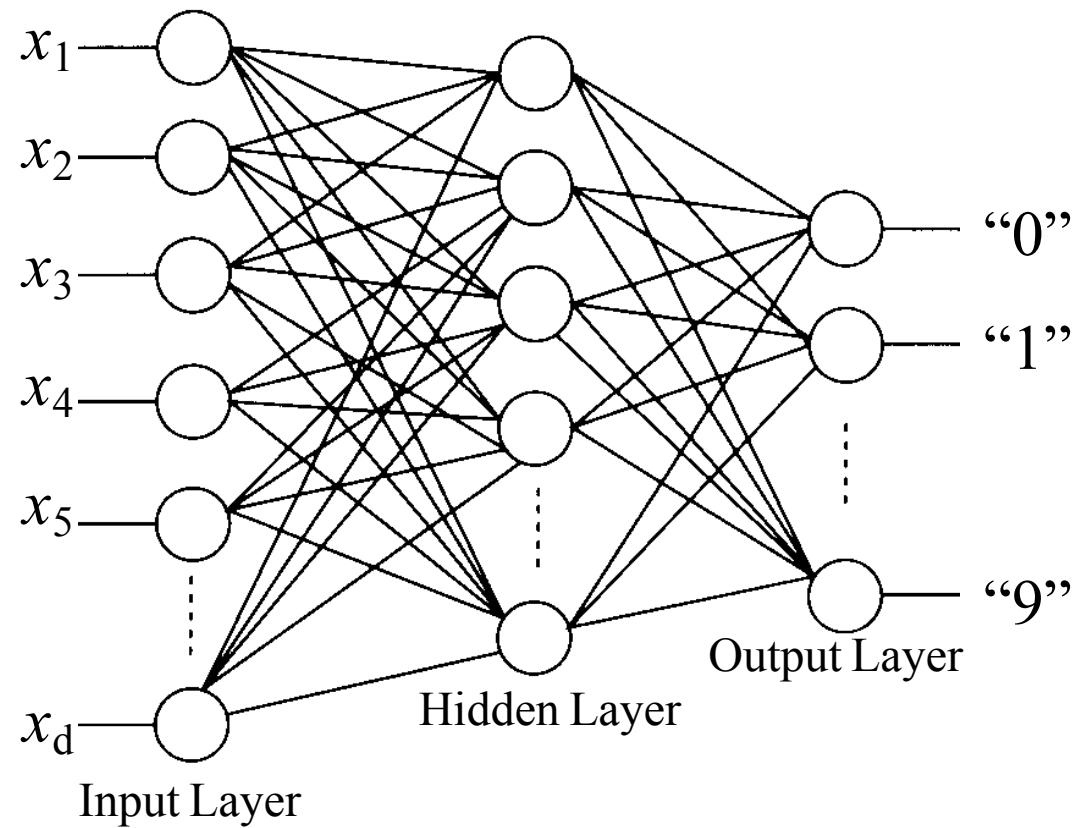
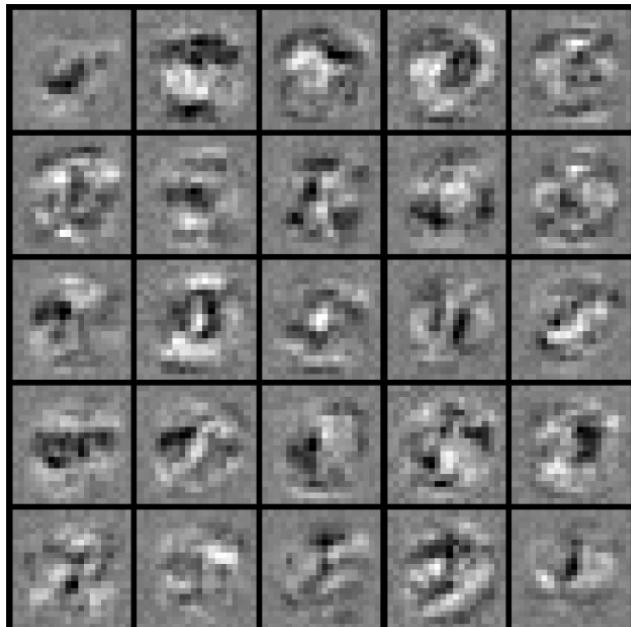
20×20 pixel images
 $d = 400$ 10 classes



Each image is “unrolled” into a vector x of pixel intensities

Layering Representations

7	9	6	5	8	7	4	4	1	0
0	7	3	3	2	4	8	4	5	1
6	6	3	2	9	1	3	3	2	6
1	3	7	1	5	6	5	2	4	4
7	0	9	8	7	5	8	9	5	4
4	6	6	5	0	2	1	3	6	9
8	5	1	8	9	3	8	7	3	6
1	0	2	8	2	3	0	5	1	5
6	7	8	2	5	3	9	7	0	0
7	9	3	9	8	5	7	2	9	8



Visualization of Hidden Layer
Lec 18: HMM&DL

This lecture

- ❖ Review: Neural Network
 - ❖ Learning NN
- ❖ Recursive and Recurrent NN

Why we need to learn this:

<https://medium.com/@karpathy/yes-you-should-understand-backprop-e2f06eab496b>

Recall: Learning as loss minimization

We have a classifier NN that is completely defined by its weights

Learn the weights by minimizing a loss L

$$\min_{\mathbf{w}} \sum_i L(NN(\mathbf{x}_i, \mathbf{w}), y_i)$$

Perhaps with a *regularizer*

So far, we saw that this strategy worked for:

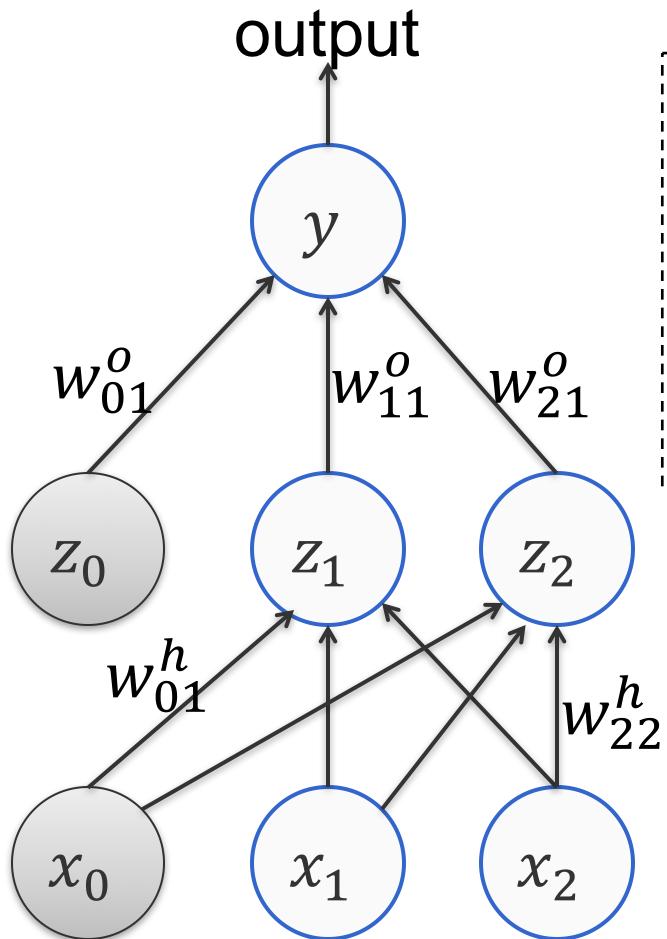
1. Logistic Regression
 2. Support Vector Machines
 3. Perceptron
 4. LMS regression
- Each
minimizes a
different loss
function

All of these are linear models

Same idea for non-linear models too!

Back to our running example

Given an input \mathbf{x} , how is the output predicted



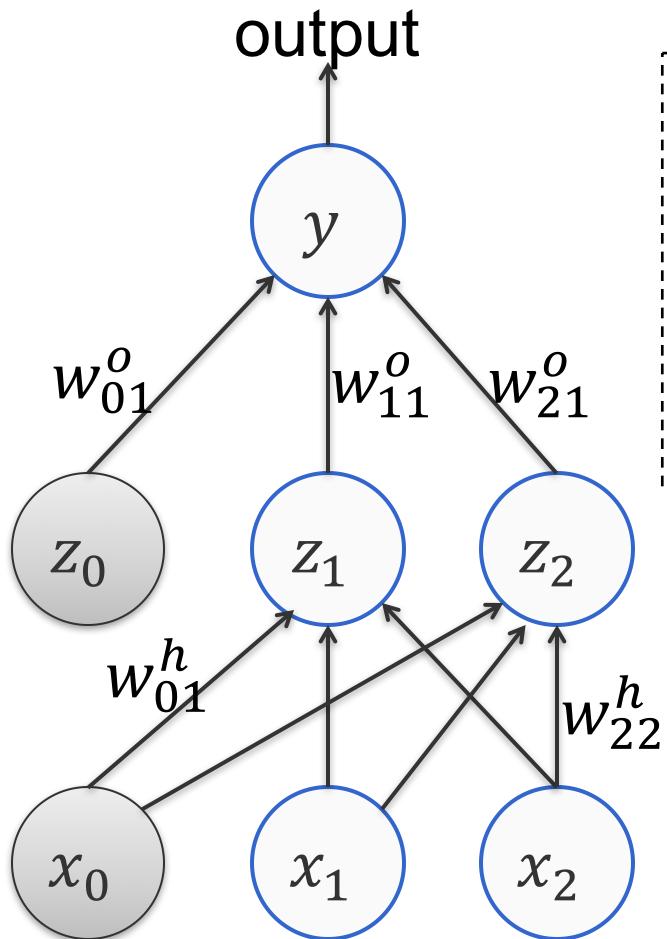
$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

Back to our running example

Given an input \mathbf{x} , how is the output predicted



$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

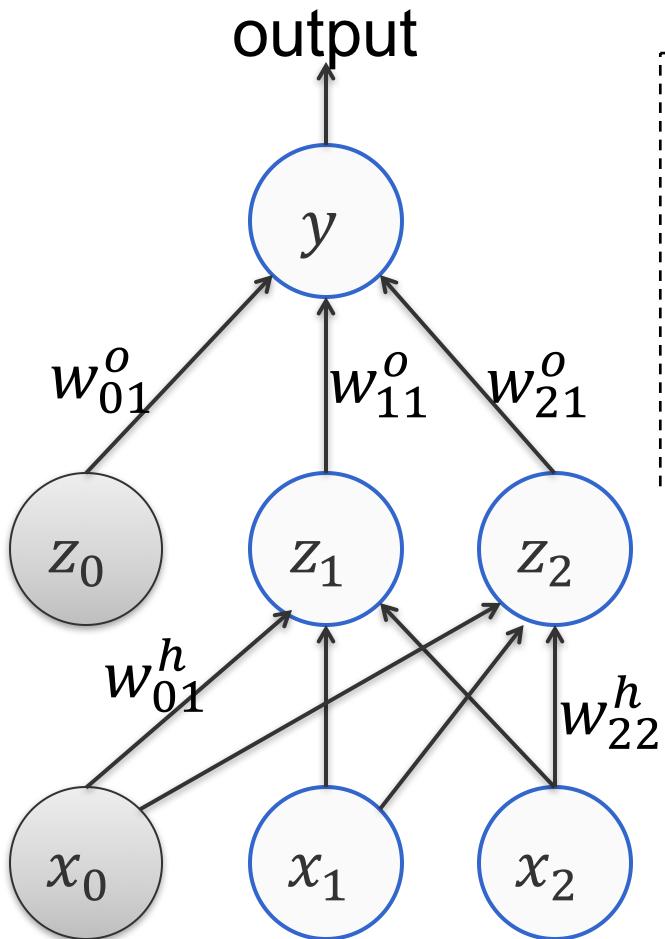
$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

Suppose the true label for this example is a number y^*

Back to our running example

Given an input \mathbf{x} , how is the output predicted



$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

Suppose the true label for this example is a number y^*

We can write the *square loss* for this example as:

$$L = \frac{1}{2} (y - y^*)^2$$

Learning as loss minimization

We have a classifier NN that is completely defined by its weights

Learn the weights by minimizing a loss L

$$\min_w \sum_i L(NN(x_i, w), y_i)$$

Perhaps with a *regularizer*

How do we solve the optimization problem?

$$\min_w \sum_i L(NN(x_i, w), y_i)$$

Stochastic gradient descent

Given a training set $S = \{(x_i, y_i)\}$

1. Initialize parameters w

2. For epoch = 1 ... T:

1. Shuffle the training set

2. For each training example (x_i, y_i) :

❖ Treat this example as the entire dataset

Compute the gradient of the loss $\nabla L(NN(x_i, w), y_i)$

❖ Update: $w \leftarrow w - \gamma_t \nabla L(NN(x_i, w), y_i)$

3. Return w

γ_t : learning rate,
many tweaks possible

$$\min_w \sum_i L(NN(x_i, w), y_i)$$

Stochastic gradient descent

Given a training set $S = \{(x_i, y_i)\}$, x

1. Initialize parameters w

2. For epoch = 1 ... T:

1. Shuffle the training set

2. For each training example (x_i, y_i) :

❖ Treat this example as the entire dataset

Compute the gradient of the loss $\nabla L(NN(x_i, w), y_i)$

❖ Update: $w \leftarrow w - \gamma_t \nabla L(NN(x_i, w), y_i)$

3. Return w

Have we solved everything?

The objective is not convex.
Initialization can be important

\circ_t : learning rate, many
tweaks possible

The derivative of the loss function?

$$\nabla L(NN(\mathbf{x}_i, \mathbf{w}), y_i)$$

If the neural network is a differentiable function, we can find the gradient

- ❖ Or maybe its sub-gradient
- ❖ This is decided by the activation functions and the loss function

It was easy for SVMs and logistic regression

- ❖ Only one layer

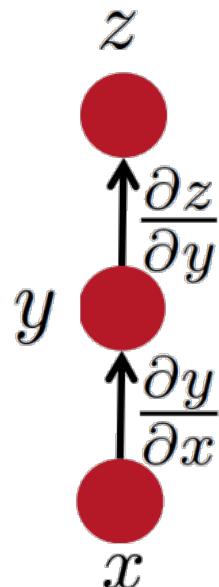
But how do we find the sub-gradient of a more complex function?

- ❖ Eg: A recent paper used a ~150 layer neural network for image classification!

We need an efficient algorithm: **Backpropagation**
Lec 18: HMM&DL

Reminder: Chain rule for derivatives

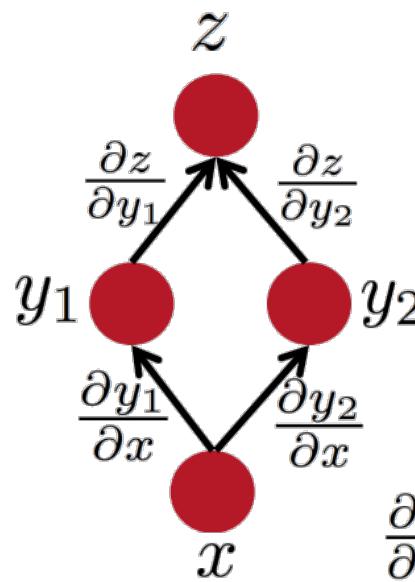
- ❖ If z is a function of y and y is a function of x
 - ❖ Then z is a function of x , as well
- ❖ Question: how to find $\frac{\partial z}{\partial x}$



$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

Reminder: Chain rule for derivatives

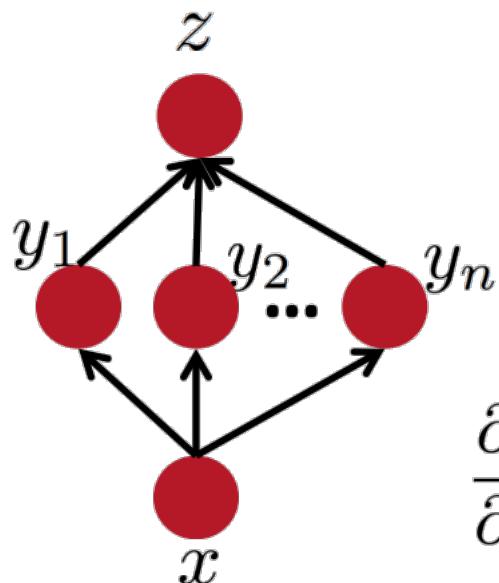
- ❖ If z = a function of y_1 + a function of y_2 , and the y_i 's are functions of x
 - ❖ Then z is a function of x , as well
- ❖ Question: how to find $\frac{\partial z}{\partial x}$



$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x} + \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial x}$$

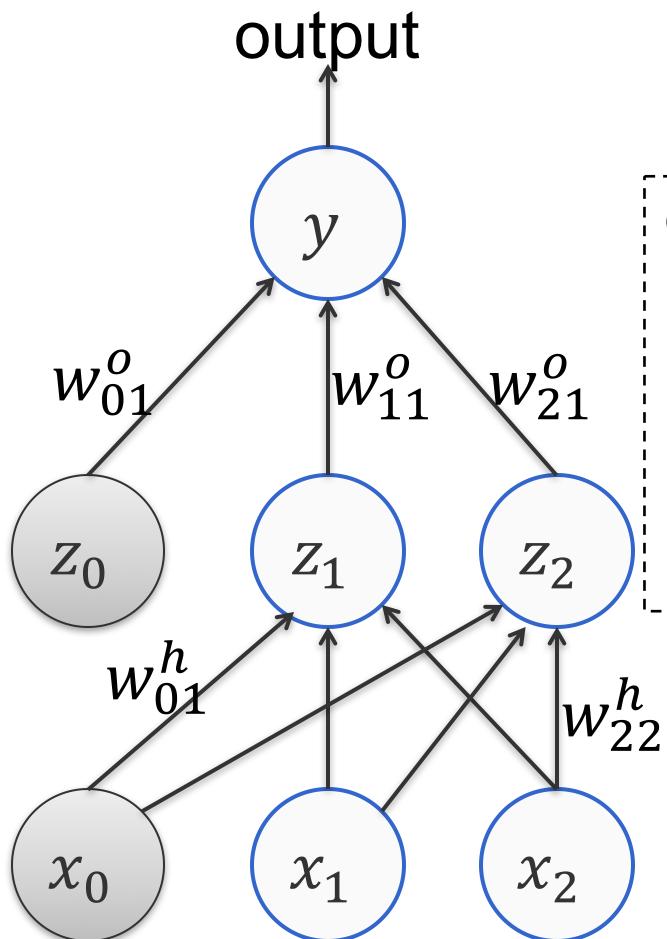
Reminder: Chain rule for derivatives

- ❖ If z is a sum of functions of y_i 's, and the y_i 's are functions of x
 - ❖ Then z is a function of x , as well
- ❖ Question: how to find $\frac{\partial z}{\partial x}$



$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$$

Backpropagation



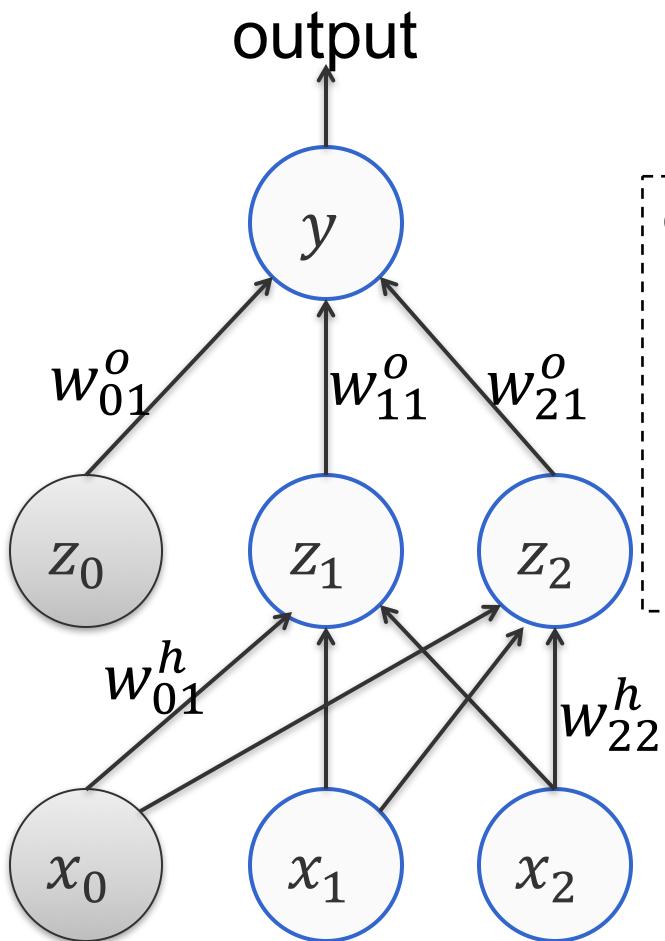
$$L = \frac{1}{2}(y - y^*)^2$$

$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

Backpropagation



$$L = \frac{1}{2}(y - y^*)^2$$

$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

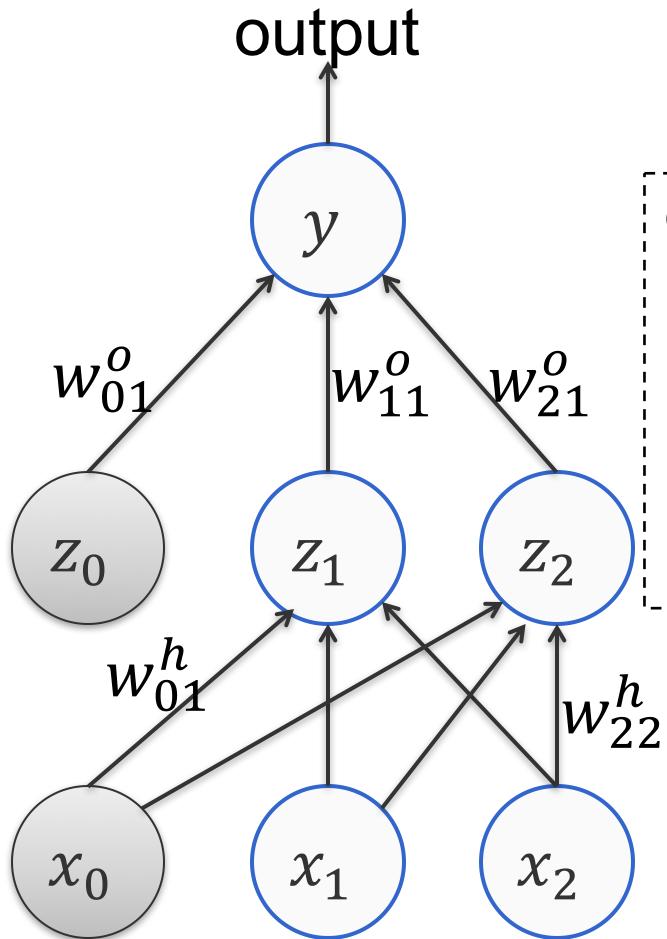
$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

We want to compute

$$\frac{\partial L}{\partial w_{ij}^o} \text{ and } \frac{\partial L}{\partial w_{ij}^h}$$

Backpropagation

Applying the chain rule to compute the gradient
(And remembering partial computations along
the way to speed up things)



$$L = \frac{1}{2}(y - y^*)^2$$

$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

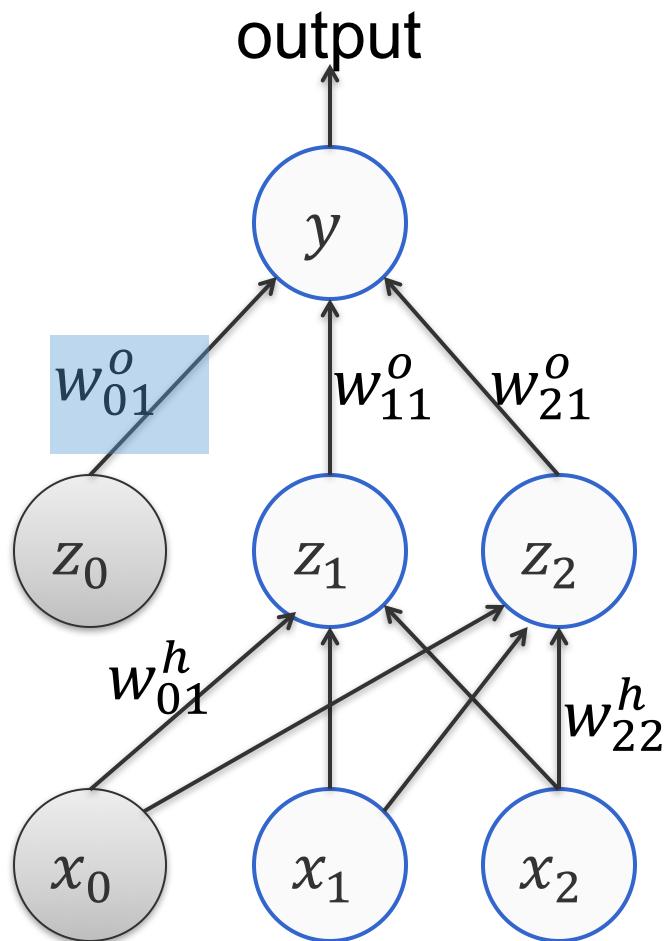
We want to compute

$$\frac{\partial L}{\partial w_{ij}^o} \text{ and } \frac{\partial L}{\partial w_{ij}^h}$$

$$L = \frac{1}{2}(y - y^*)^2$$

output $y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$

Output layer

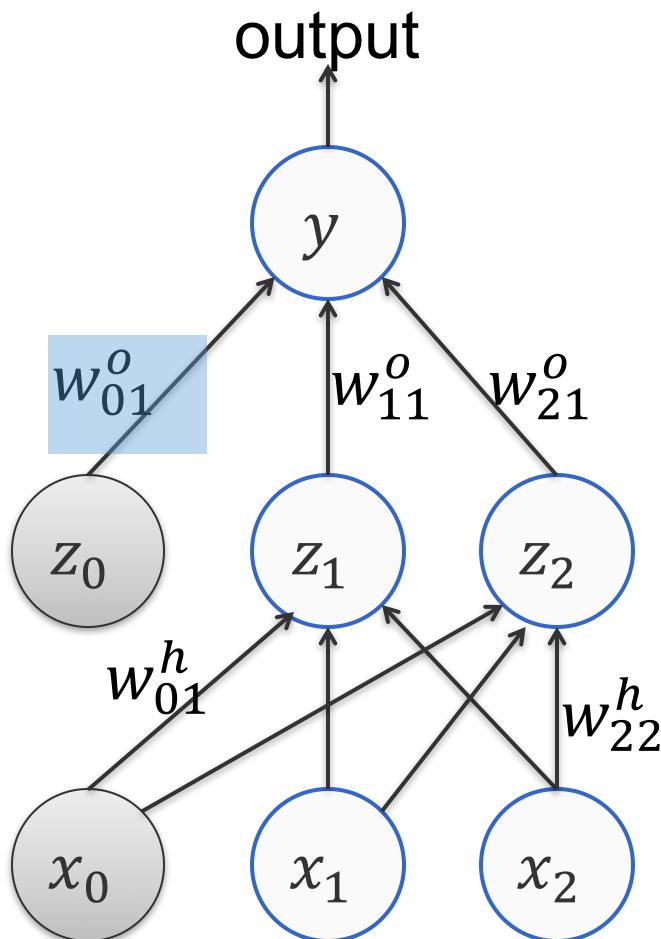


$$\frac{\partial L}{\partial w_{01}^o}$$

$$L = \frac{1}{2}(y - y^*)^2$$

output $y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$

Output layer

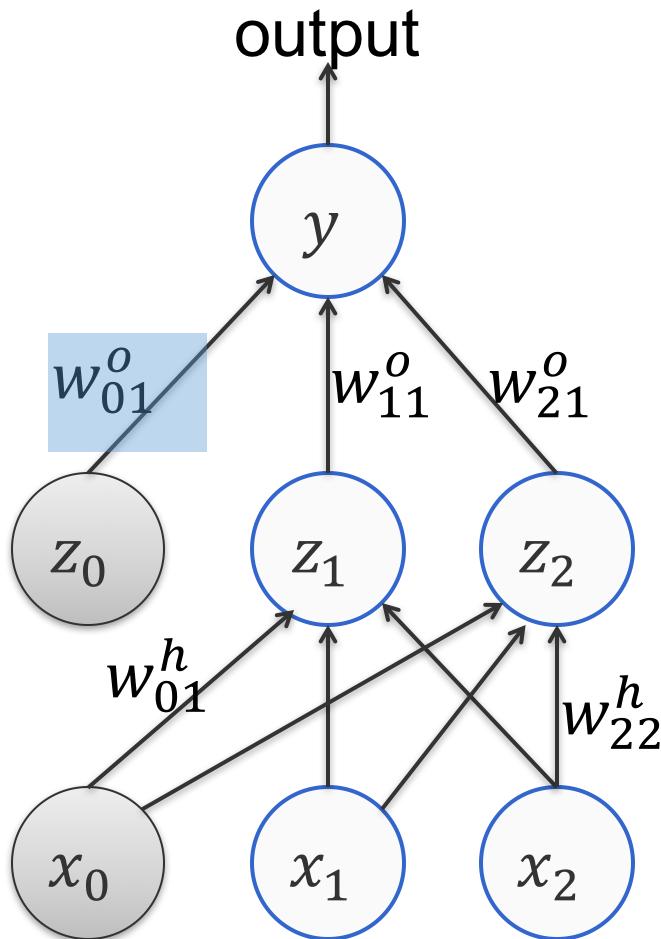


$$\frac{\partial L}{\partial w_{01}^o} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{01}^o}$$

$$L = \frac{1}{2}(y - y^*)^2$$

output $y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$

Output layer



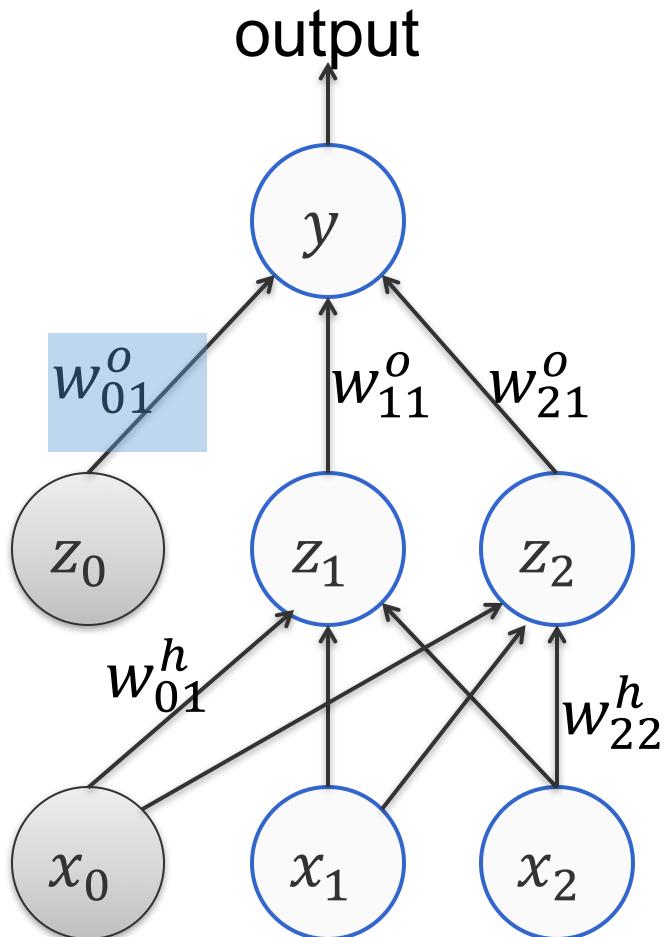
$$\frac{\partial L}{\partial w_{01}^o} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{01}^o}$$

$$\frac{\partial L}{\partial y} = y - y^*$$

$$L = \frac{1}{2}(y - y^*)^2$$

output $y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$

Output layer



$$\frac{\partial L}{\partial w_{01}^o} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{01}^o}$$

$$\frac{\partial L}{\partial y} = y - y^*$$

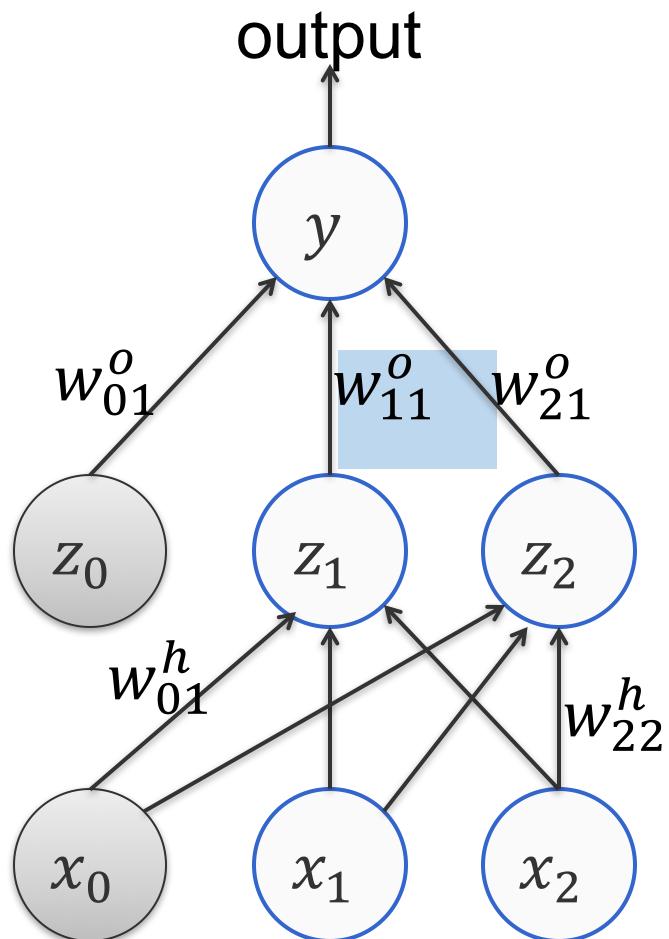
$$\frac{\partial y}{\partial w_{01}^o} = 1$$

Backpropagation example

$$L = \frac{1}{2}(y - y^*)^2$$

output $y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$

Output layer

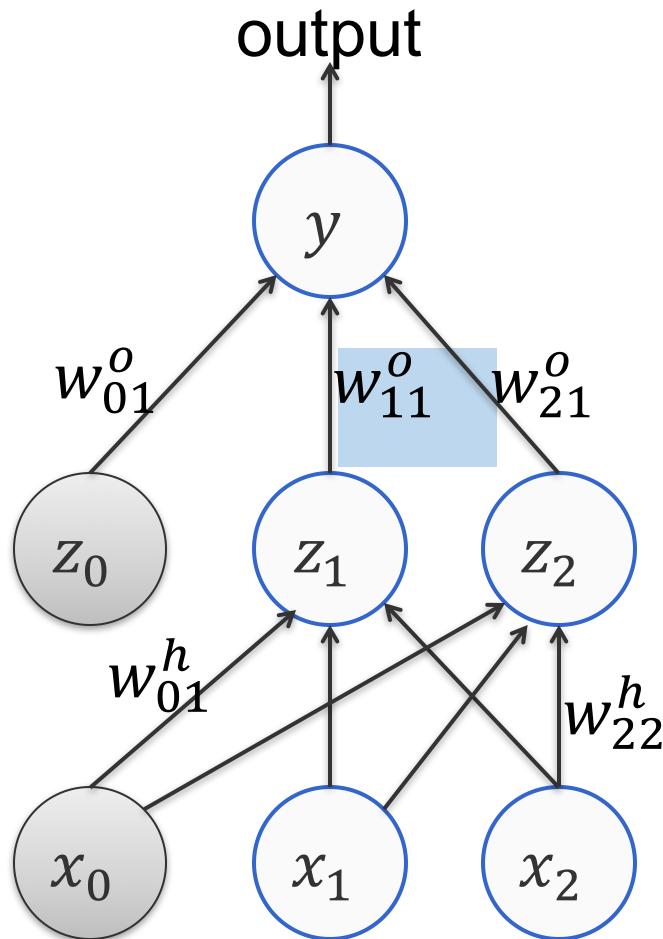


$$\frac{\partial L}{\partial w_{11}^o}$$

$$L = \frac{1}{2}(y - y^*)^2$$

output $y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$

Output layer

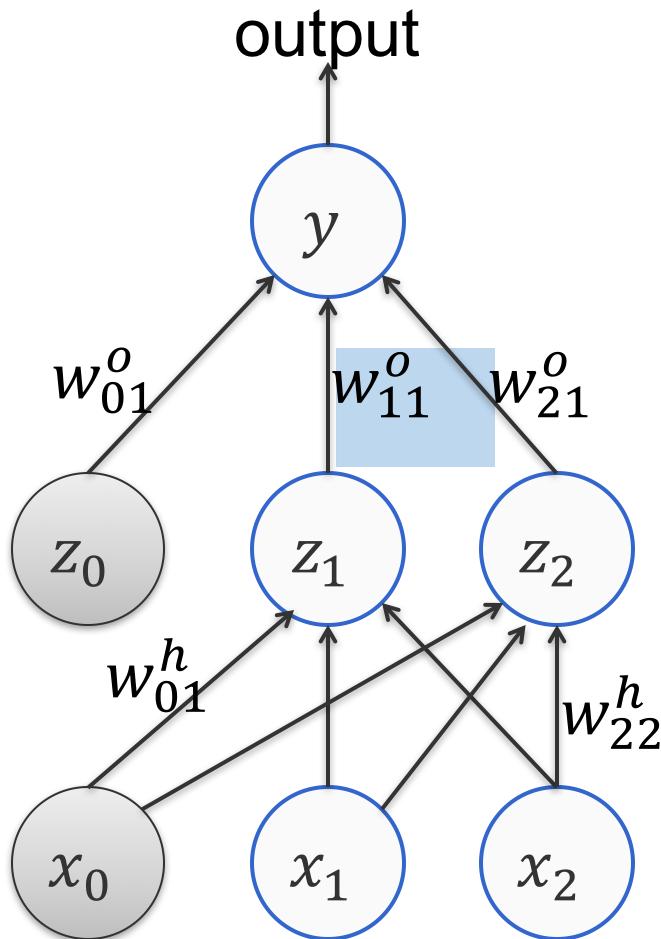


$$\frac{\partial L}{\partial w_{11}^o} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{11}^o}$$

$$L = \frac{1}{2}(y - y^*)^2$$

output $y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$

Output layer



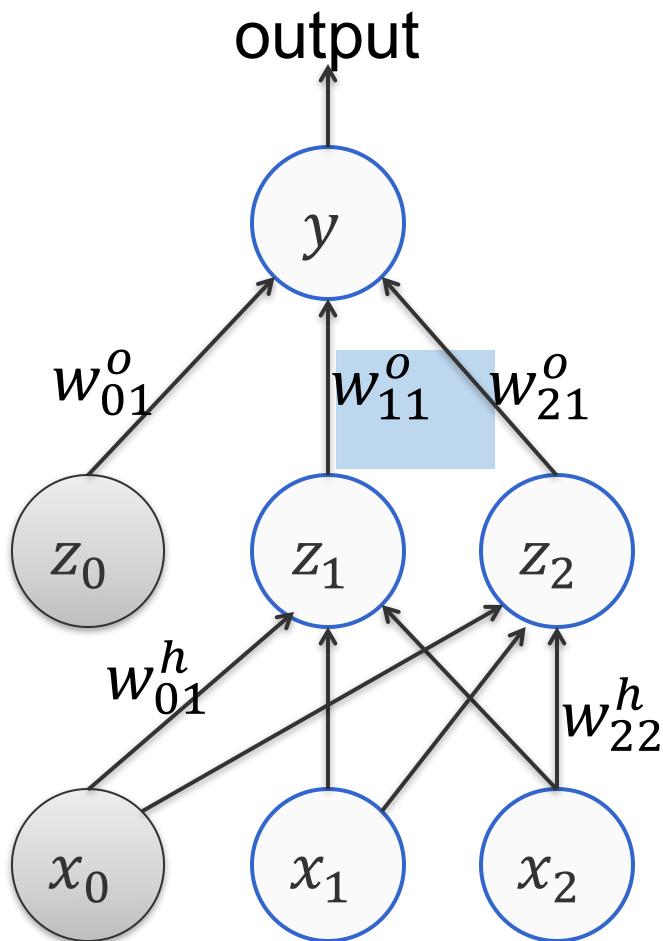
$$\frac{\partial L}{\partial w_{11}^o} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{11}^o}$$

$$\frac{\partial L}{\partial y} = y - y^*$$

$$L = \frac{1}{2}(y - y^*)^2$$

output $y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$

Output layer



$$\frac{\partial L}{\partial w_{11}^o} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{11}^o}$$

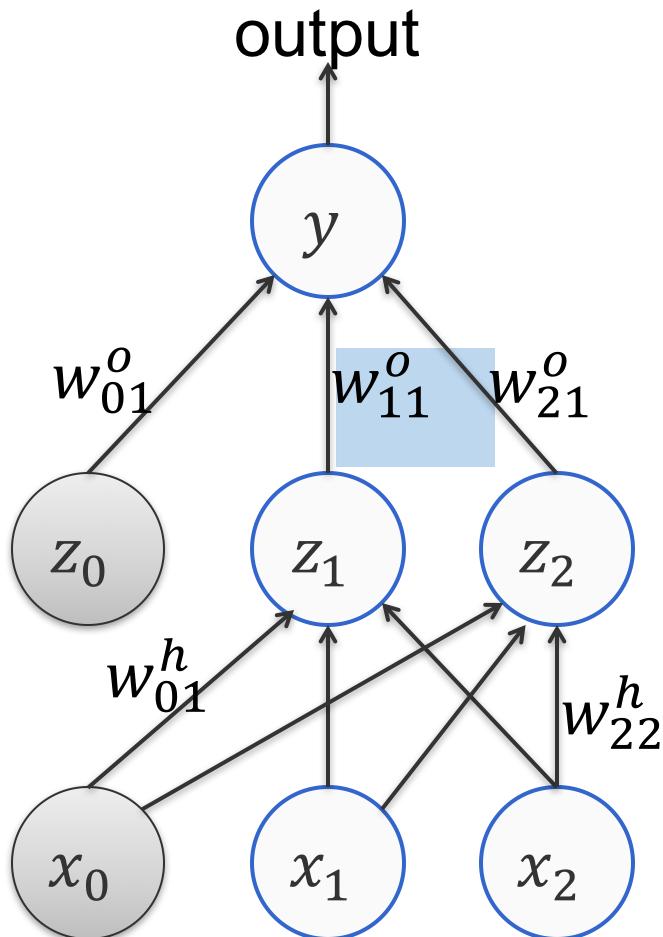
$$\frac{\partial L}{\partial y} = y - y^*$$

$$\frac{\partial y}{\partial w_{01}^o} = z_1$$

$$L = \frac{1}{2}(y - y^*)^2$$

output $y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$

Output layer



$$\frac{\partial L}{\partial w_{11}^o} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{11}^o}$$

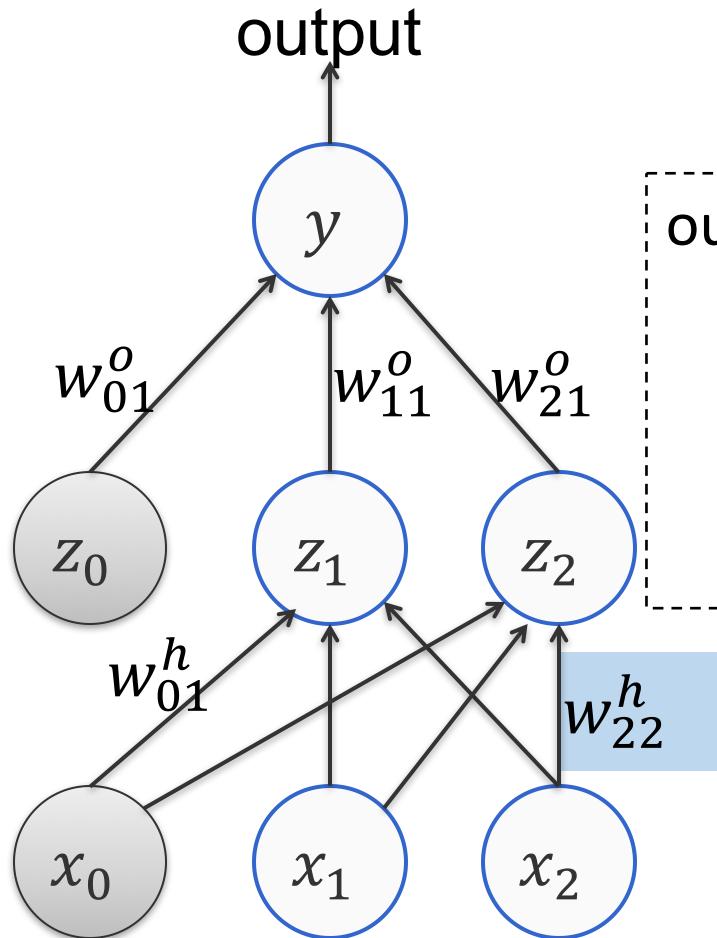
$$\frac{\partial L}{\partial y} = y - y^*$$

$$\frac{\partial y}{\partial w_{01}^o} = z_1$$

We have already computed this partial derivative for the previous case

Cache to speed up!

Hidden layer derivatives



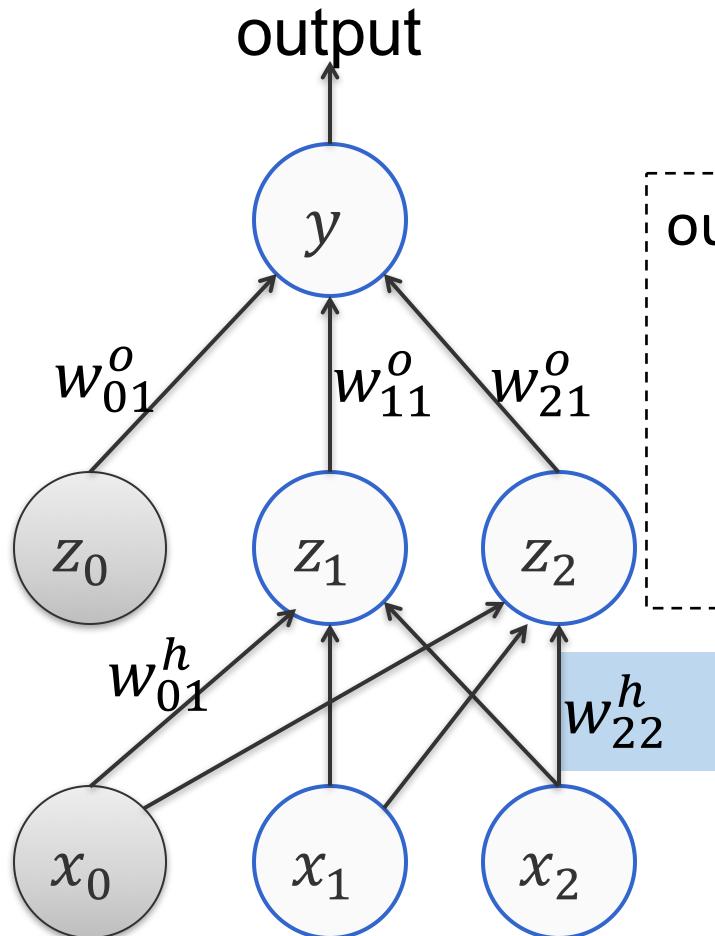
$$L = \frac{1}{2}(y - y^*)^2$$

$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

Hidden layer derivatives



$$L = \frac{1}{2}(y - y^*)^2$$

$$\text{output } y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

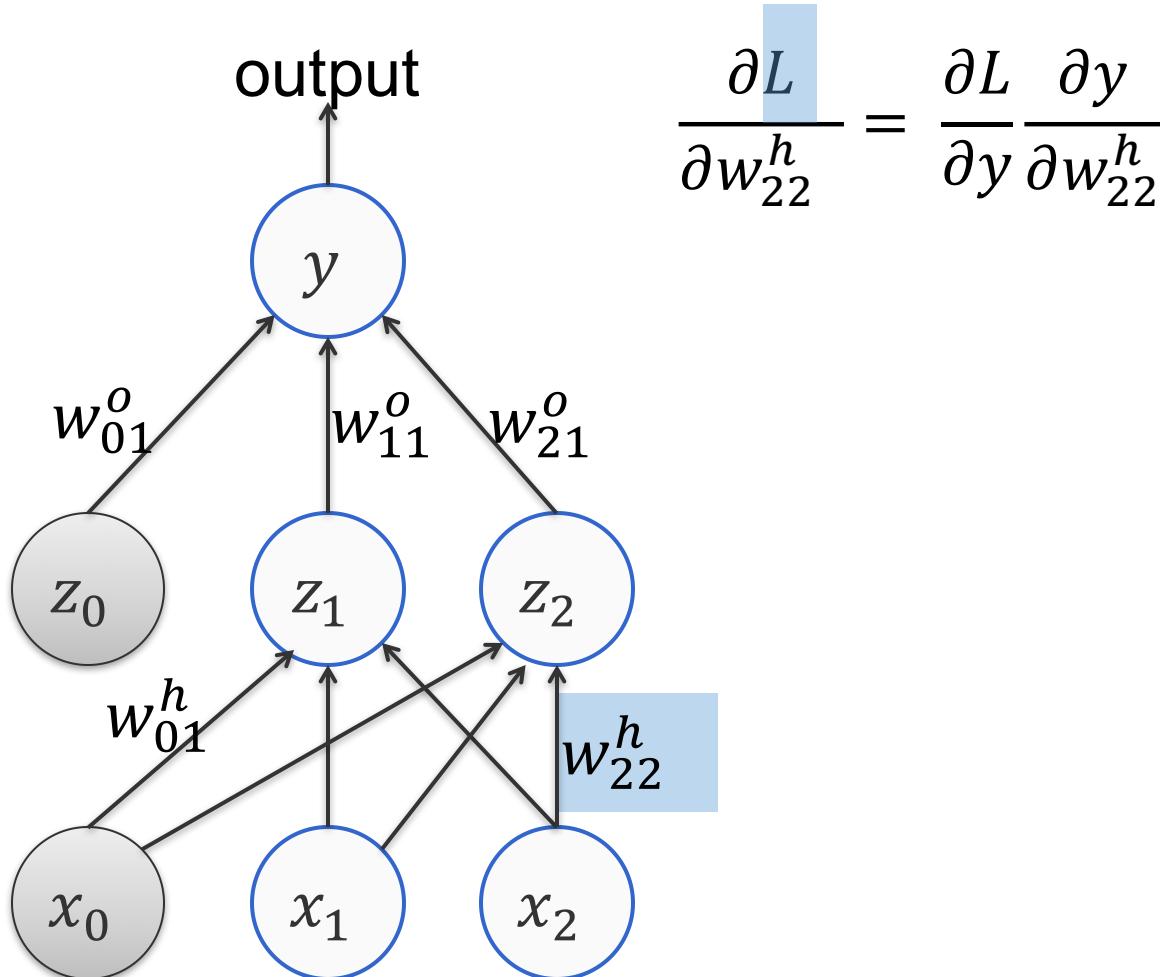
$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

$$z_1 = \sigma(w_{01}^h + w_{11}^h x_1 + w_{21}^h x_2)$$

We want $\frac{\partial L}{\partial w_{22}^h}$

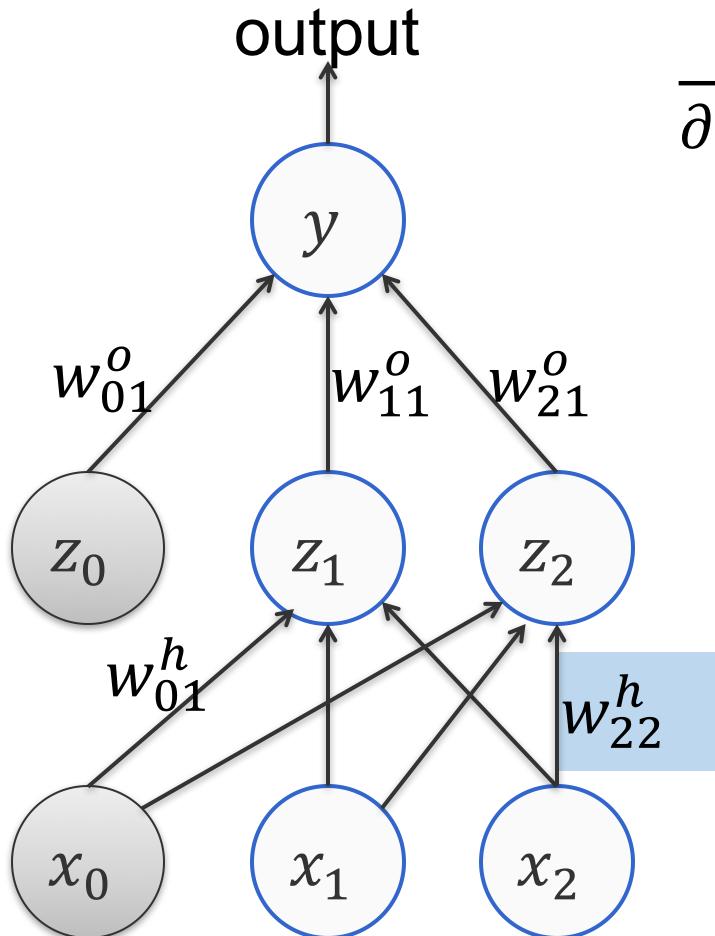
$$L = \frac{1}{2}(y - y^*)^2$$

Hidden layer derivatives



$$y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

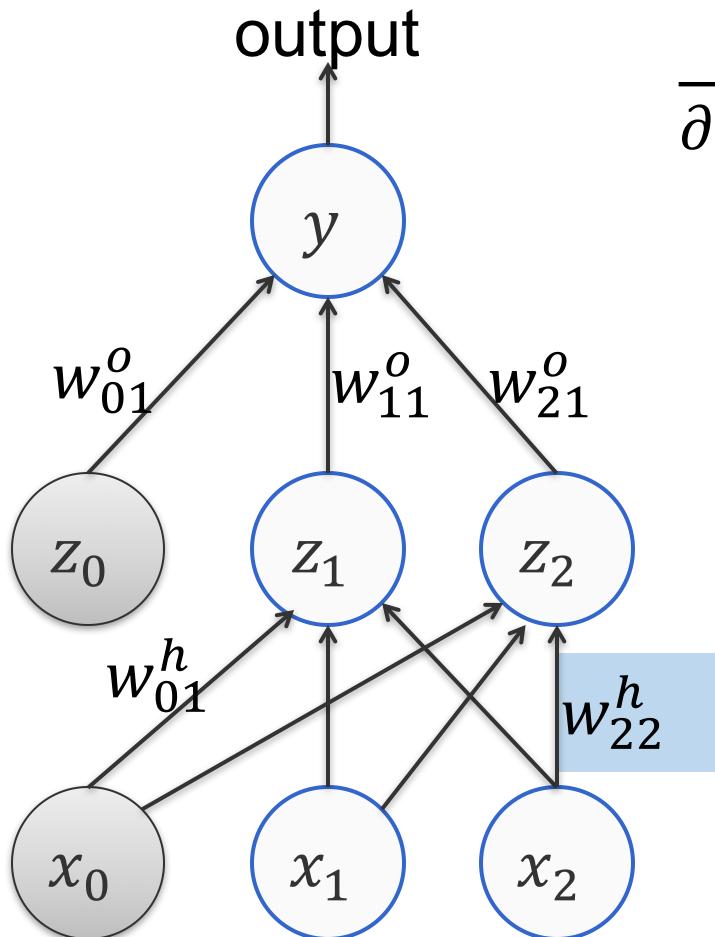
Hidden layer



$$\begin{aligned}\frac{\partial L}{\partial w_{22}^h} &= \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{22}^h} \\ &= \frac{\partial L}{\partial y} \frac{\partial}{\partial w_{22}^h} (w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2)\end{aligned}$$

$$y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

Hidden layer



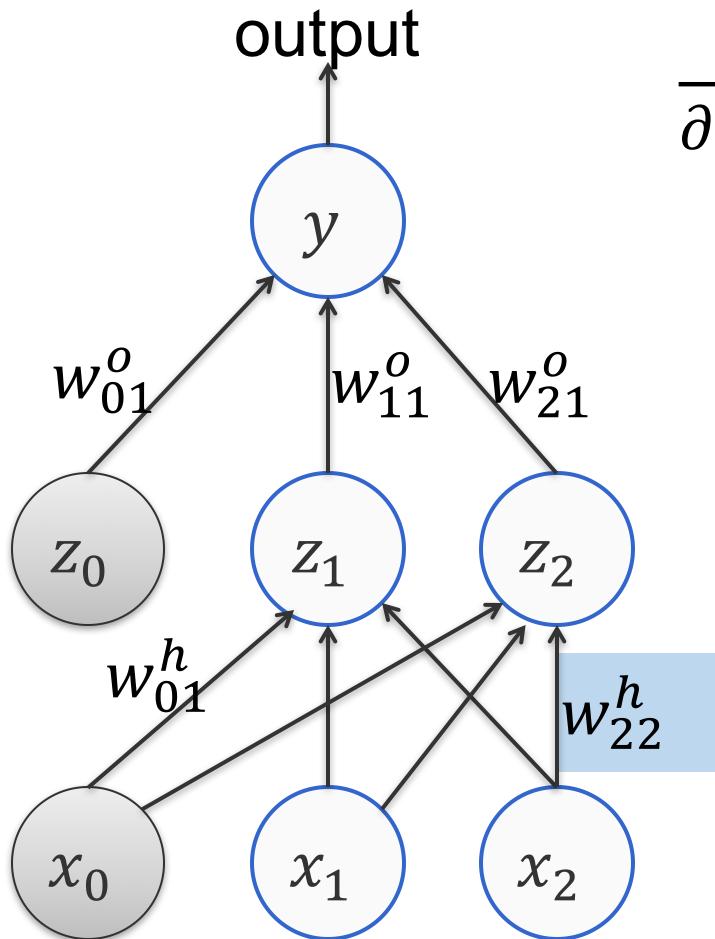
$$\begin{aligned}\frac{\partial L}{\partial w_{22}^h} &= \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{22}^h} \\ &= \frac{\partial L}{\partial y} \frac{\partial}{\partial w_{22}^h} (w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2) \\ &= \frac{\partial L}{\partial y} \left(w_{11}^o \frac{\partial}{\partial w_{22}^h} z_1 + w_{21}^o \frac{\partial}{\partial w_{22}^h} z_2 \right)\end{aligned}$$

0

z_1 is not a function of w_{22}^h

$$y = w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2$$

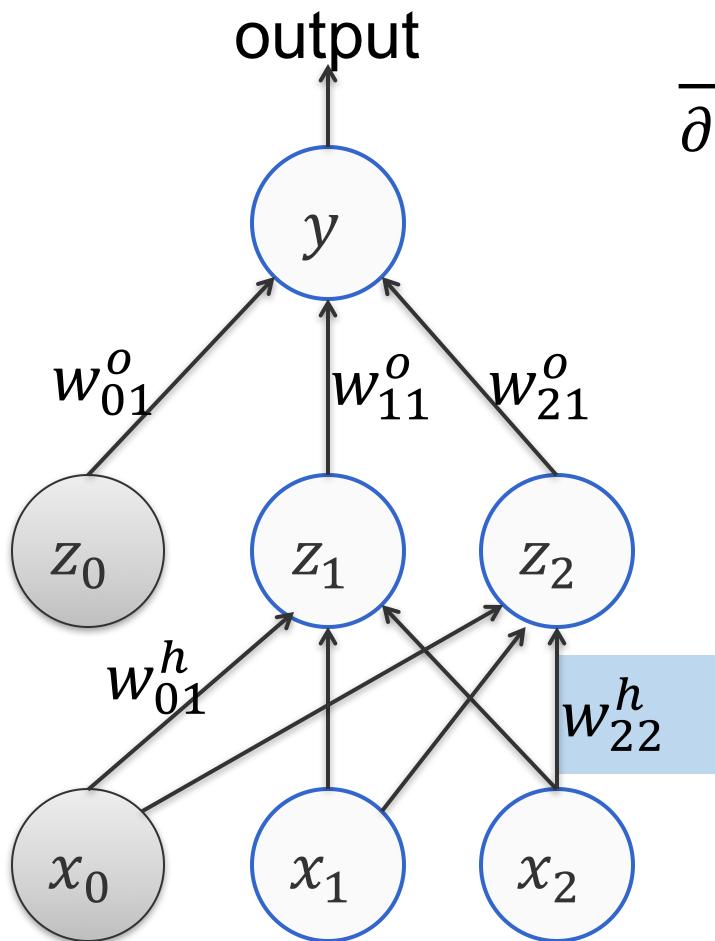
Hidden layer



$$\begin{aligned}\frac{\partial L}{\partial w_{22}^h} &= \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{22}^h} \\ &= \frac{\partial L}{\partial y} \frac{\partial}{\partial w_{22}^h} (w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2) \\ &= \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial w_{22}^h}\end{aligned}$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

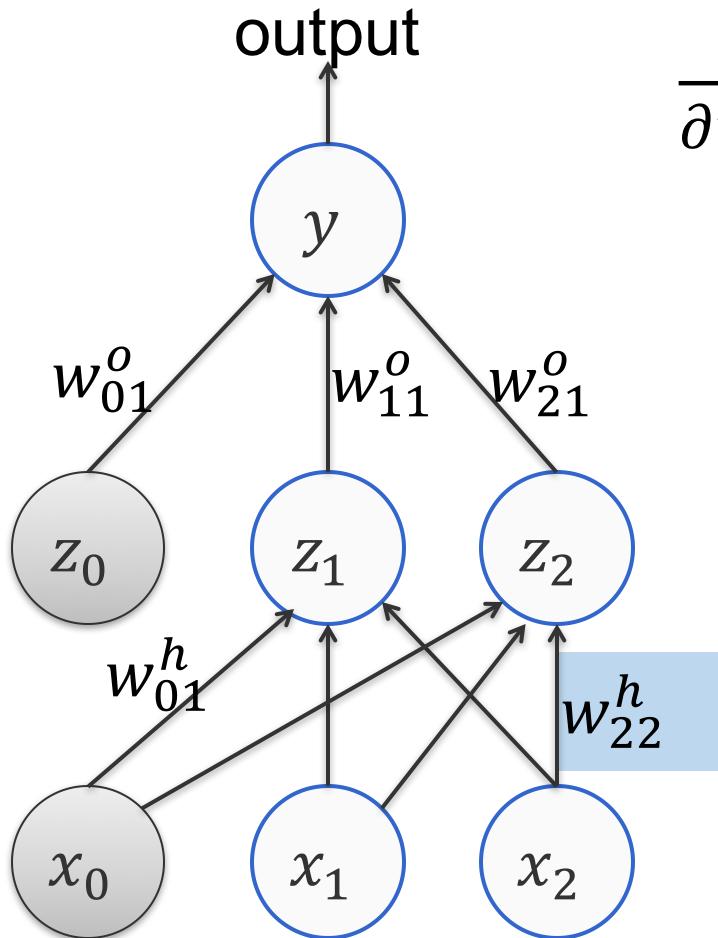
Hidden layer



$$\begin{aligned}\frac{\partial L}{\partial w_{22}^h} &= \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{22}^h} \\ &= \frac{\partial L}{\partial y} \frac{\partial}{\partial w_{22}^h} (w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2) \\ &= \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial w_{22}^h}\end{aligned}$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

Hidden layer

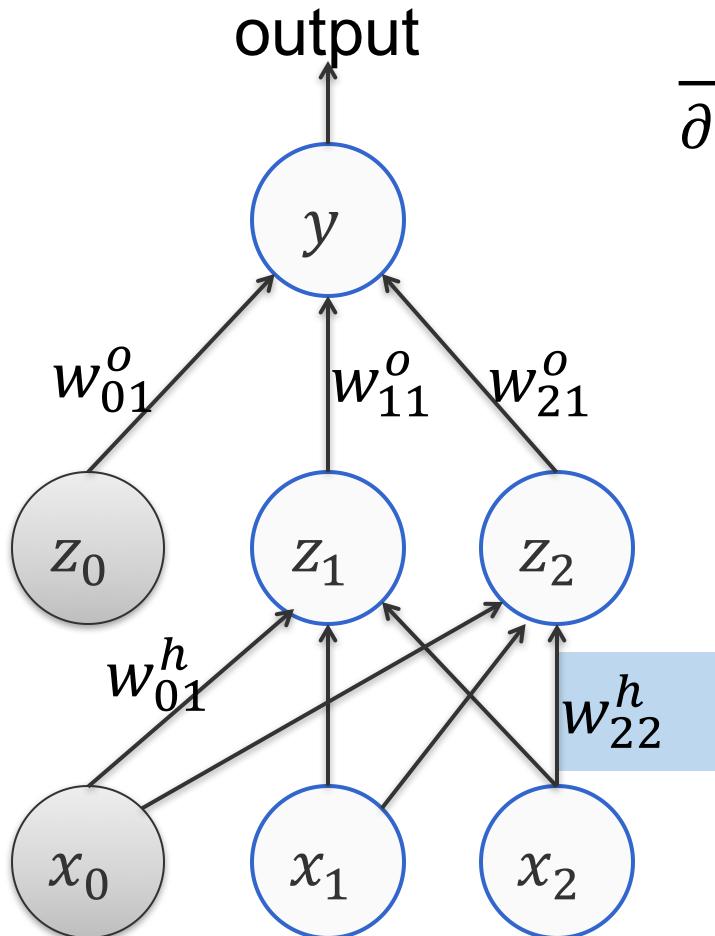


$$\begin{aligned}\frac{\partial L}{\partial w_{22}^h} &= \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{22}^h} \\ &= \frac{\partial L}{\partial y} \frac{\partial}{\partial w_{22}^h} (w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2) \\ &= \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial w_{22}^h}\end{aligned}$$

Call this s

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

Hidden layer



$$\begin{aligned}\frac{\partial L}{\partial w_{22}^h} &= \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{22}^h} \\ &= \frac{\partial L}{\partial y} \frac{\partial}{\partial w_{22}^h} (w_{01}^o + w_{11}^o z_1 + w_{21}^o z_2) \\ &= \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial w_{22}^h} \\ &= \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial s} \frac{\partial s}{\partial w_{22}^h}\end{aligned}$$

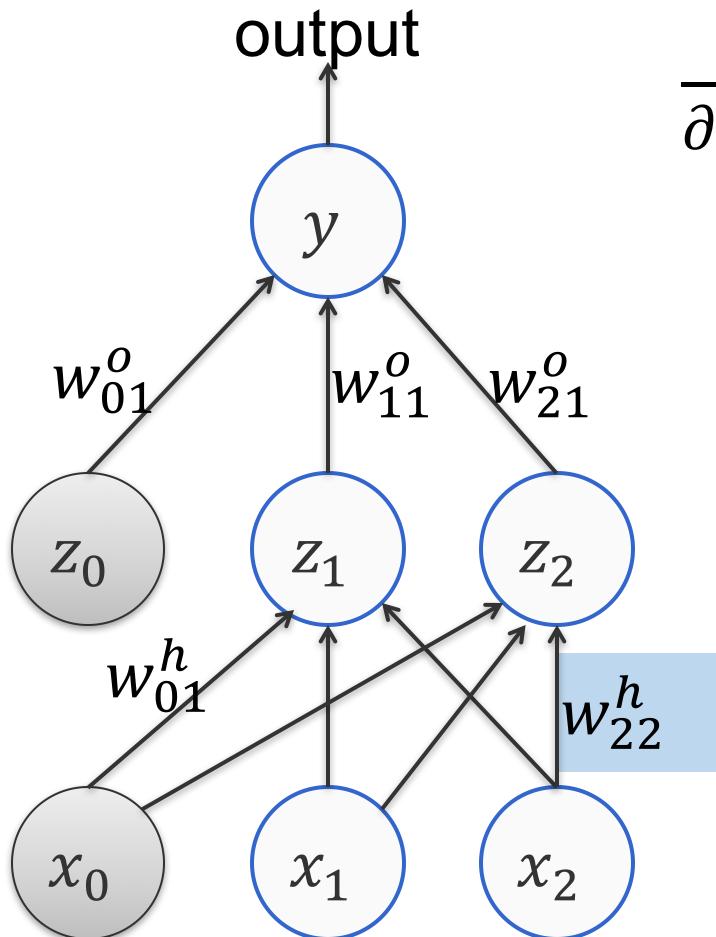
Call this s

Backpropagation example

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

Hidden layer

Call this s

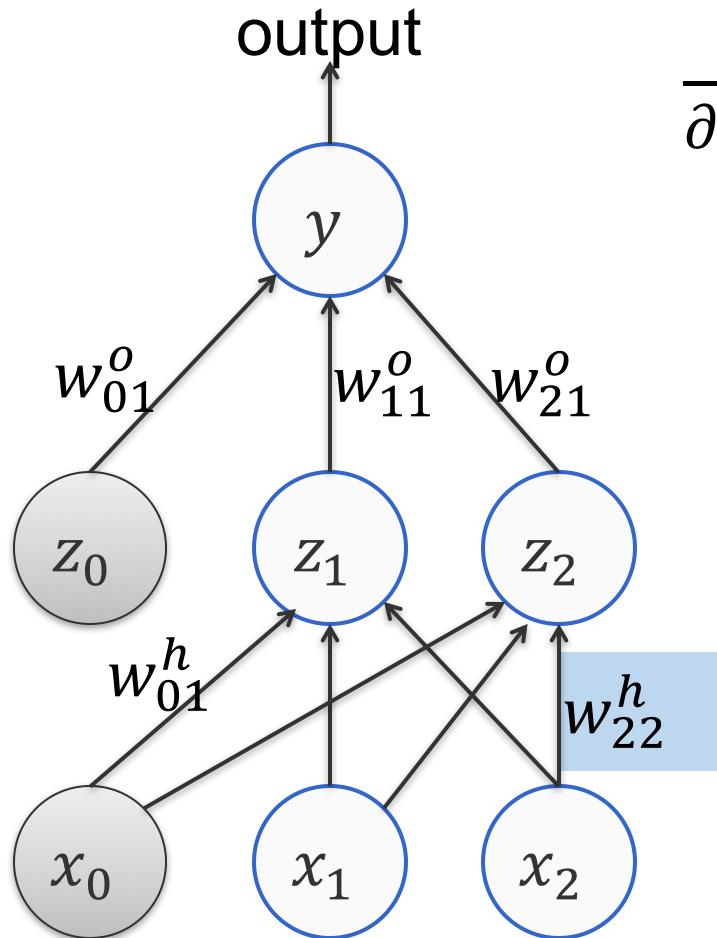


$$\frac{\partial L}{\partial w_{22}^h} = \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial s} \frac{\partial s}{\partial w_{22}^h} \text{ (From previous slide)}$$

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

Hidden layer

Call this s

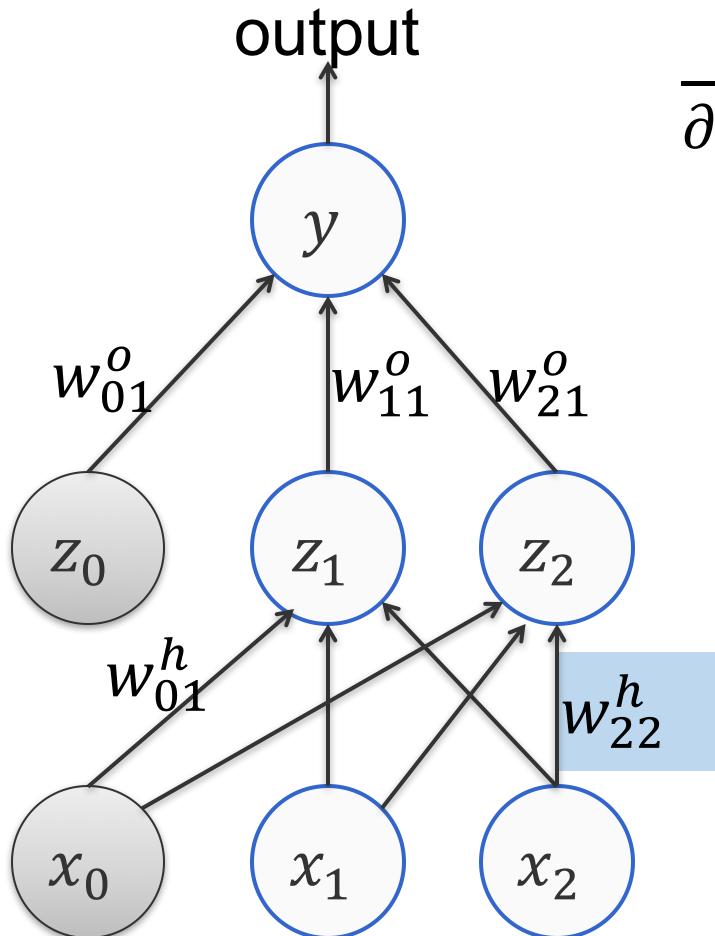


$$\frac{\partial L}{\partial w_{22}^h} = \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial s} \frac{\partial s}{\partial w_{22}^h}$$

Each of these partial derivatives is easy

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

Hidden layer



$$\frac{\partial L}{\partial w_{22}^h} = \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial s} \frac{\partial s}{\partial w_{22}^h}$$

Each of these partial derivatives is easy

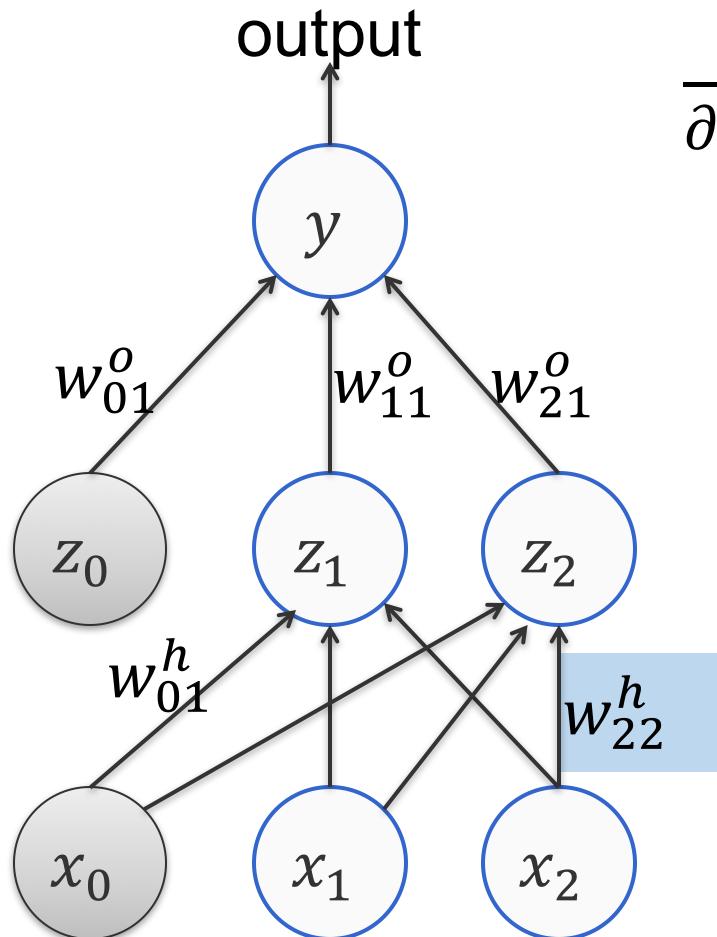
$$\frac{\partial L}{\partial y} = y - y^*$$

Call this s

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

Hidden layer

Call this s



$$\frac{\partial L}{\partial w_{22}^h} = \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial s} \frac{\partial s}{\partial w_{22}^h}$$

Each of these partial derivatives is easy

$$\frac{\partial L}{\partial y} = y - y^*$$

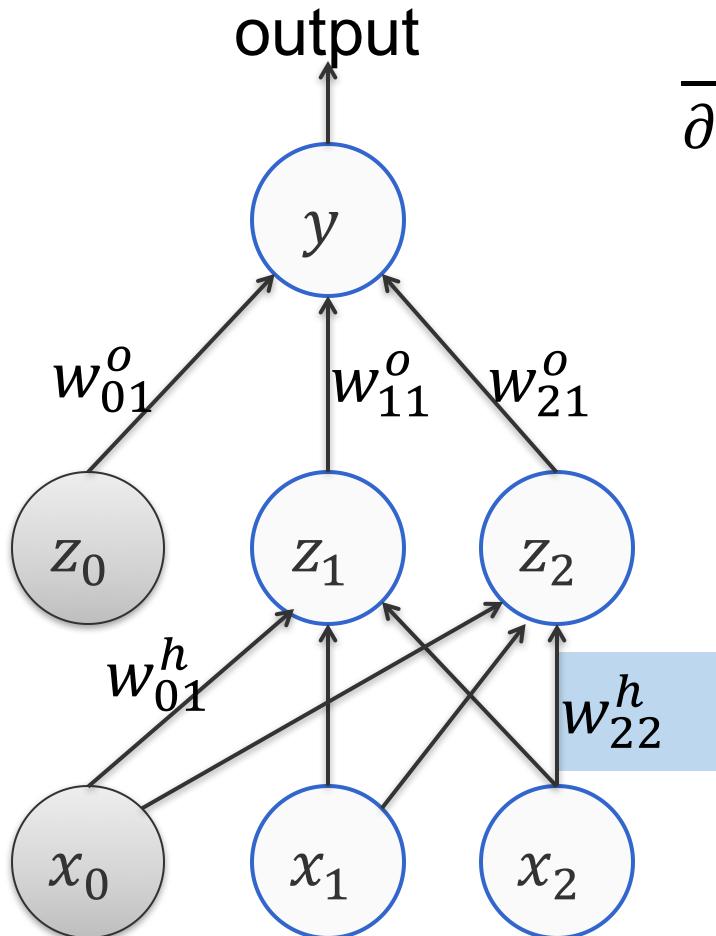
$$\frac{\partial z_2}{\partial s} = z_2(1 - z_2)$$

Why? Because $z_2(s)$ is the logistic function we have already seen

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

Hidden layer

Call this s



$$\frac{\partial L}{\partial w_{22}^h} = \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial s} \frac{\partial s}{\partial w_{22}^h}$$

Each of these partial derivatives is easy

$$\frac{\partial L}{\partial y} = y - y^*$$

$$\frac{\partial z_2}{\partial s} = z_2(1 - z_2)$$

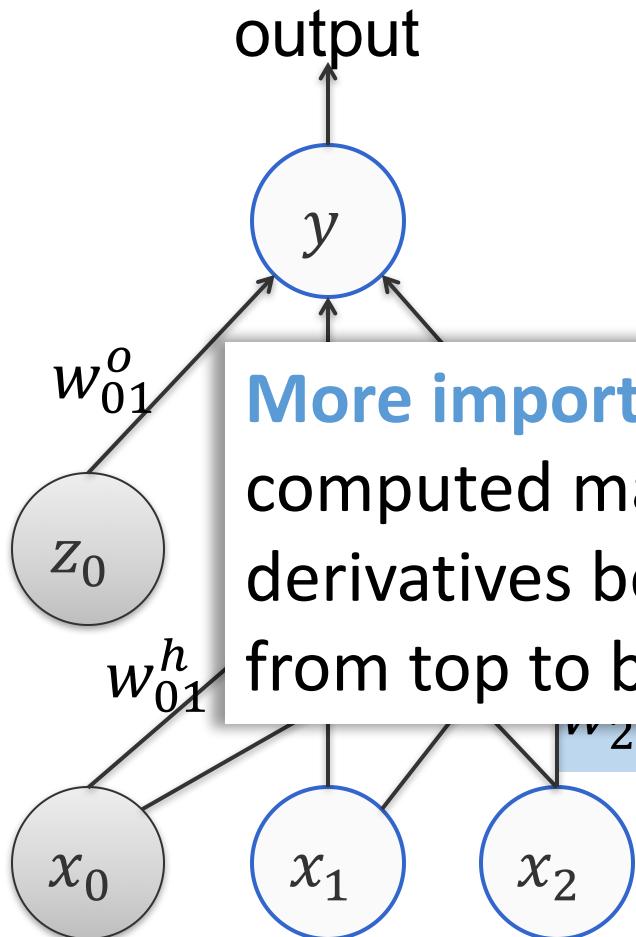
$$\frac{\partial s}{\partial w_{22}^h} = x_2$$

Why? Because $z_2(s)$ is the logistic function we have already seen

$$z_2 = \sigma(w_{02}^h + w_{12}^h x_1 + w_{22}^h x_2)$$

Hidden layer

Call this s



$$\frac{\partial L}{\partial w_{22}^h} = \frac{\partial L}{\partial y} w_{21}^o \frac{\partial z_2}{\partial s} \frac{\partial s}{\partial w_{22}^h}$$

Each of these partial derivatives is easy

More important: We have already computed many of these partial derivatives because we are proceeding from top to bottom (i.e. backwards)

$$\frac{\partial L}{\partial w_{22}^h} = -x_2$$

cause $z_2(s)$
istic
we have
een

The Backpropagation Algorithm

The same algorithm works for multiple layers

Repeated application of the chain rule for partial derivatives

- ❖ First perform forward pass from inputs to the output
- ❖ Compute loss

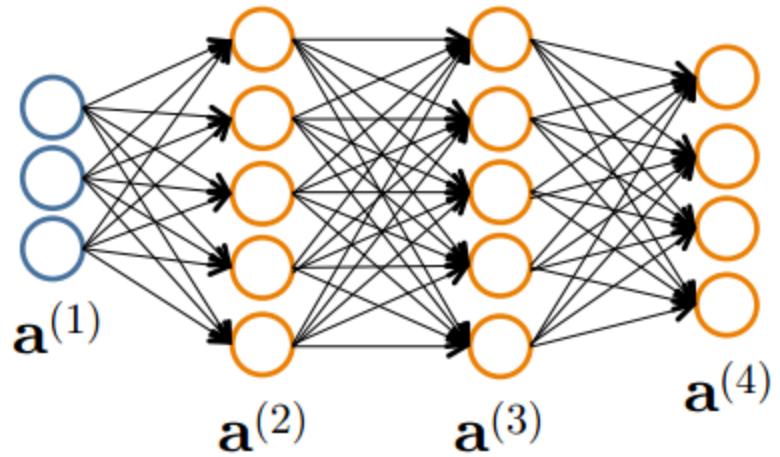
- ❖ From the loss, proceed backwards to compute partial derivatives using the chain rule
- ❖ Cache partial derivatives as you compute them
 - ❖ Will be used for lower layers

Forward Propagation

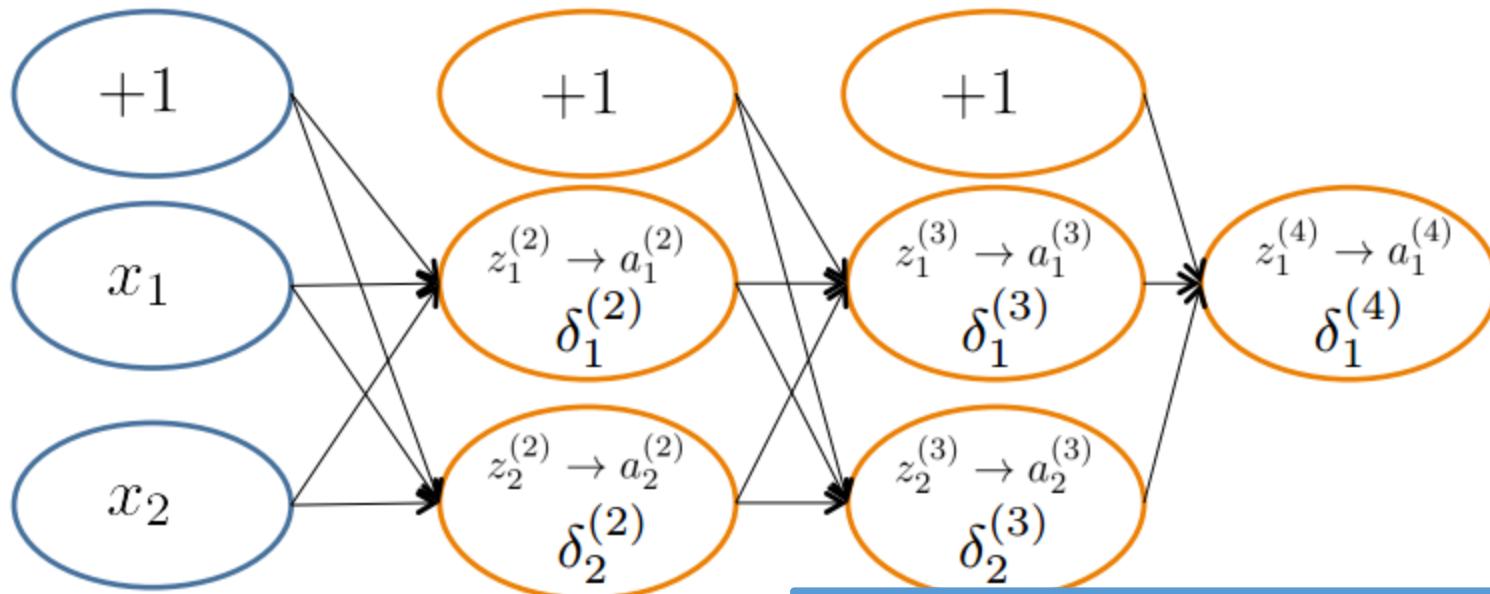
- Given one labeled training instance (\mathbf{x}, y) :

Forward Propagation

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \Theta^{(1)}\mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$ [add $a_0^{(2)}$]
- $\mathbf{z}^{(3)} = \Theta^{(2)}\mathbf{a}^{(2)}$
- $\mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$ [add $a_0^{(3)}$]
- $\mathbf{z}^{(4)} = \Theta^{(3)}\mathbf{a}^{(3)}$
- $\mathbf{a}^{(4)} = h_{\Theta}(\mathbf{x}) = g(\mathbf{z}^{(4)})$



Backpropagation: Compute Gradient



$$\frac{d}{dt} f(g(t)) = f'(g(t))g'(t) = \frac{df}{dg} \cdot \frac{dg}{dt}$$

$\delta_j^{(l)}$ = “error” of node j in layer l

Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where $\text{cost}(\mathbf{x}_i) = y_i \log h_{\Theta}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\Theta}(\mathbf{x}_i))$

Mechanizing learning

- ❖ Backpropagation gives you the gradient that will be used for gradient descent
 - ❖ SGD gives us a generic learning algorithm
 - ❖ Backpropagation is a generic method for computing partial derivatives
- ❖ A recursive algorithm that proceeds from the top of the network to the bottom
- ❖ Modern neural network libraries implement automatic differentiation using backpropagation
 - ❖ Allows easy exploration of network architectures
 - ❖ Don't have to keep deriving the gradients by hand each time

$$\min_w \sum_i L(NN(x_i, w), y_i)$$

Stochastic gradient descent

Given a training set $S = \{(x_i, y_i)\}$

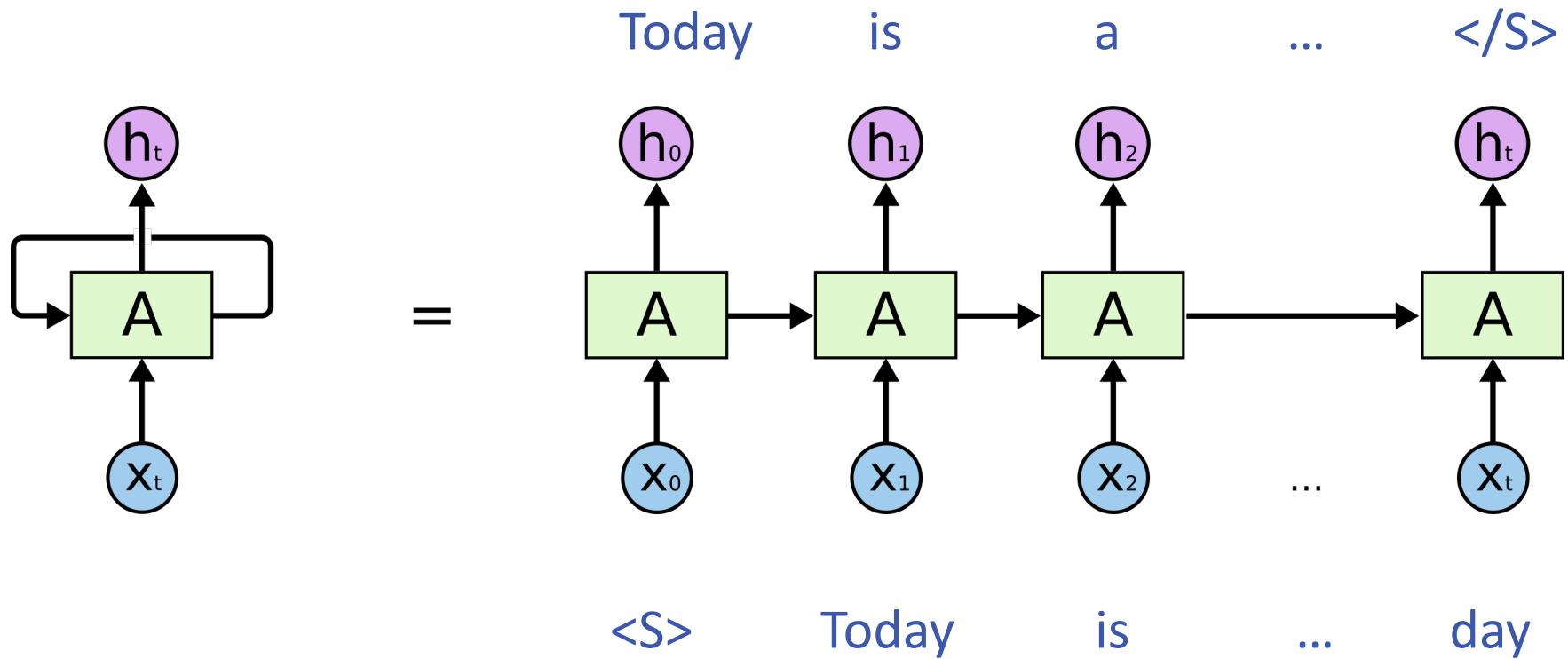
1. Initialize parameters w
2. For epoch = 1 ... T:
 1. Shuffle the training set
 2. For each training example $(x_i, y_i) \in S$:
 - ❖ Treat this example as the entire dataset
 - ❖ Compute the gradient of the loss $\nabla L(NN(x_i, w), y_i)$ using backpropagation
 - ❖ Update: $w \leftarrow w - \gamma_t \nabla L(NN(x_i, w), y_i)$
3. Return w

This lecture

- ❖ Neural Network
- ❖ Recurrent NN

How to deal with input with variant size?

- ❖ Use same parameters

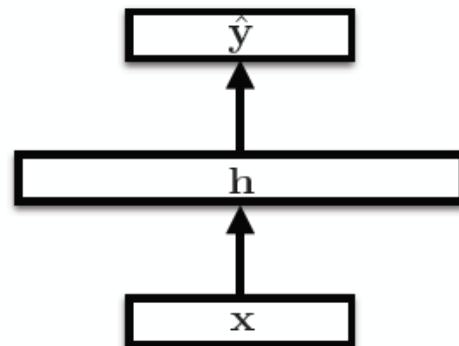


Recurrent Neural Networks

Feed-forward NN

$$\mathbf{h} = g(\mathbf{V}\mathbf{x} + \mathbf{c})$$

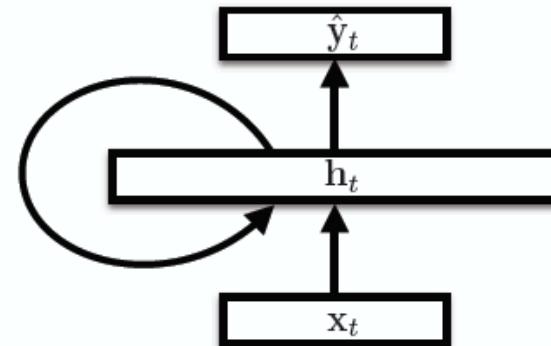
$$\hat{\mathbf{y}} = \mathbf{W}\mathbf{h} + \mathbf{b}$$



Recurrent NN

$$\mathbf{h}_t = g(\mathbf{V}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c})$$

$$\hat{\mathbf{y}}_t = \mathbf{W}\mathbf{h}_t + \mathbf{b}$$

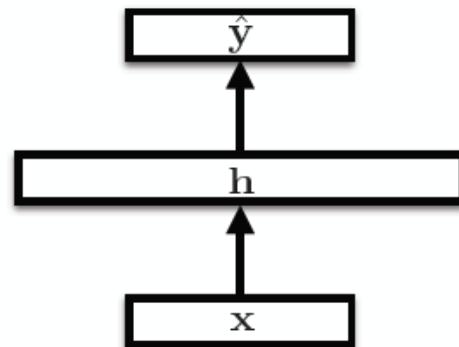


Recurrent Neural Networks

Feed-forward NN

$$\mathbf{h} = g(\mathbf{V}\mathbf{x} + \mathbf{c})$$

$$\hat{\mathbf{y}} = \mathbf{W}\mathbf{h} + \mathbf{b}$$

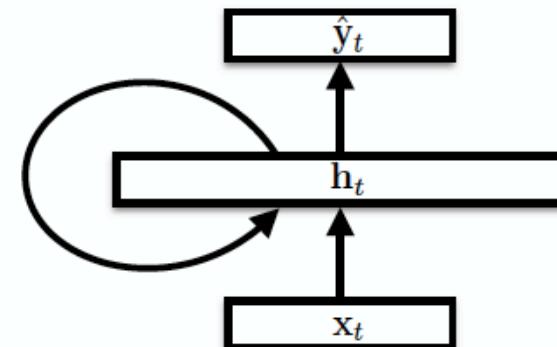


Recurrent NN

~~$$\mathbf{h}_t = g(\mathbf{V}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c})$$~~

$$\mathbf{h}_t = g(\mathbf{V}[\mathbf{x}_t; \mathbf{h}_{t-1}] + \mathbf{c})$$

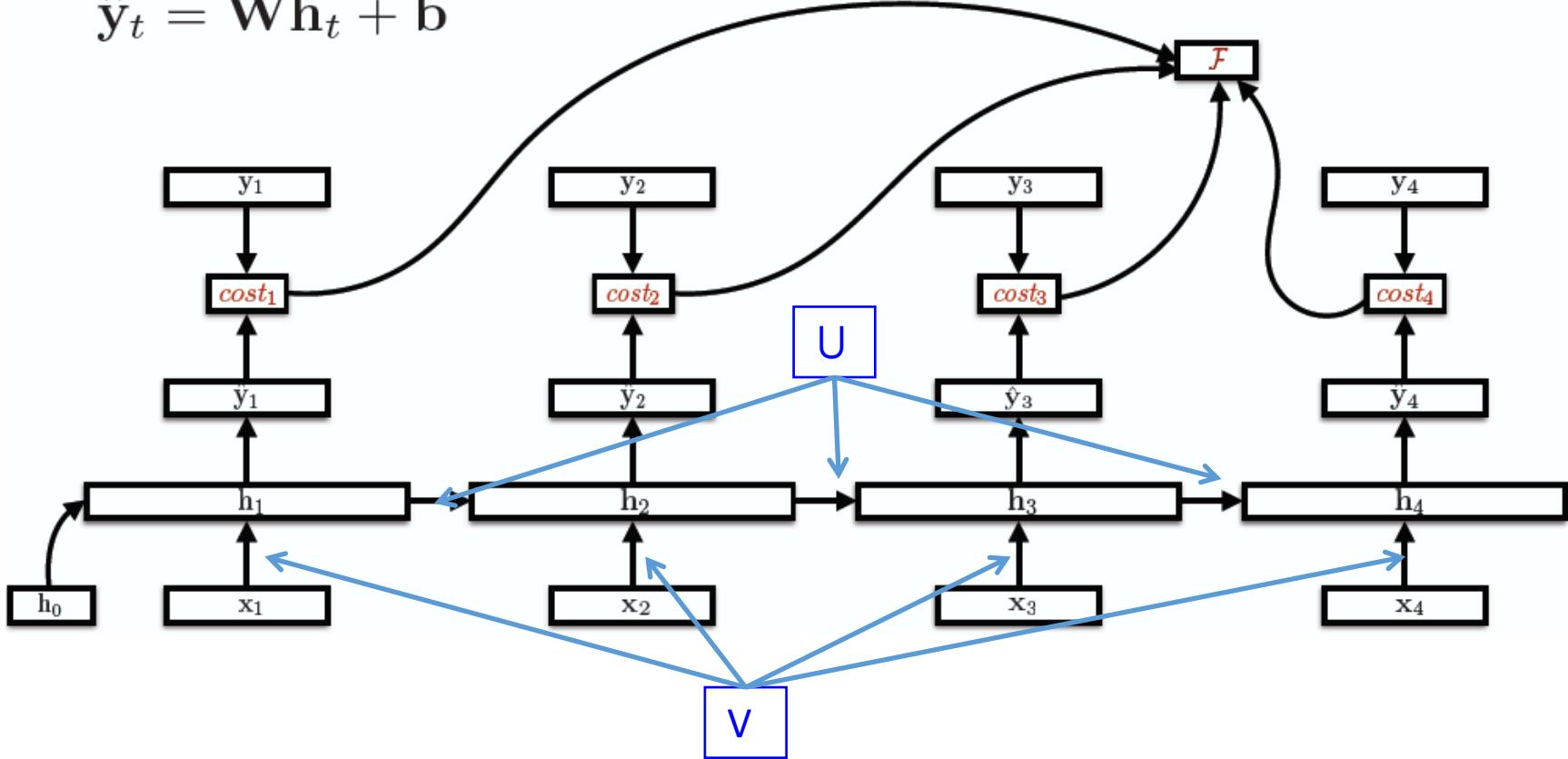
$$\hat{\mathbf{y}}_t = \mathbf{W}\mathbf{h}_t + \mathbf{b}$$



Unroll RNNs

$$\mathbf{h}_t = g(\mathbf{V}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c})$$

$$\hat{\mathbf{y}}_t = \mathbf{W}\mathbf{h}_t + \mathbf{b}$$

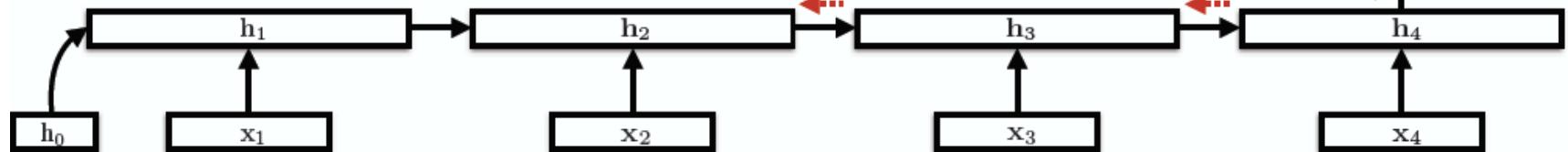


RNN training

❖ Back-propagation over time

$$\mathbf{h}_t = g(\mathbf{V}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c})$$

$$\hat{\mathbf{y}} = \mathbf{W}\mathbf{h}_{|\mathbf{x}|} + \mathbf{b}$$



What happens to gradients as you go back
in time?

$$\frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_4}{\partial \mathbf{h}_3} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}_4} \frac{\partial \mathcal{F}}{\partial \hat{\mathbf{y}}} \frac{\partial \mathcal{F}}{\partial \mathcal{F}}$$

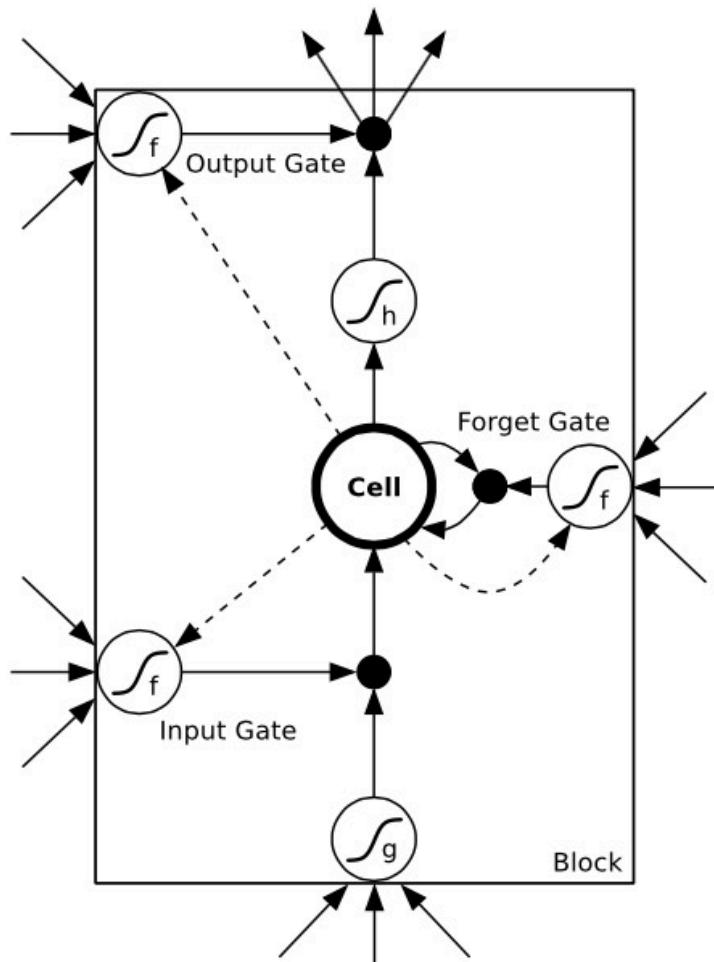
Vanishing Gradients

- ❖ For the traditional activation functions, each gradient term has the value in range $(-1, 1)$.
- ❖ Multiplying n of these small numbers to compute gradients
- ❖ The longer the sequence is, the more severe the problems are.

RNNs characteristics

- ❖ Model hidden states (input) dependencies
- ❖ Errors “back propagation over time”
- ❖ Feature learning methods
- ❖ Vanishing gradient problem:
cannot model long-distant dependencies of
the hidden states.

Long-Short Term Memory Networks (LSTMs)



$$\begin{pmatrix} i_t \\ f_t \\ o_t \\ g_t \end{pmatrix} = \begin{pmatrix} \sigma(W_i[x_t, h_t] + b_i) \\ \sigma(W_f[x_t, h_t] + b_f) \\ \sigma(W_o[x_t, h_t] + b_o) \\ f(W_g[x_t, h_t] + b_g) \end{pmatrix}$$

$$c_t = f_t * c_{t-1} + i_t * g_t$$

$$h_t = o_t * f(c_t)$$

Use gates to control the information to be added from the **input**, **forgot** from the previous memories, and **outputted**. σ and f are *sigmoid* and *tanh* function respectively, to map the value to $[-1, 1]$