

CS60050 Machine Learning

Neural Networks

Introduction, Backpropagation, Activation Functions

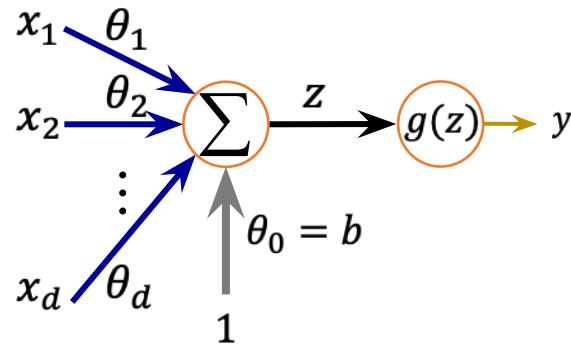
Somak Aditya, Sudeshna Sarkar
Department of CSE, IIT Kharagpur



Neural Networks

- Origins: Algorithms (loosely) inspired by the brain.
- Very widely used in 80s and early 90s; popularity diminished in late 90s. (*the famed AI Winter*)
- Recent resurgence: State-of-the-art technique for many applications
- Disclaimer: Artificial neural networks are not nearly as complex or intricate as the actual brain structure

Logistic Regression



$$\mathbf{w} = [\theta_1 \ \theta_2 \ \dots \ \theta_d]^T$$

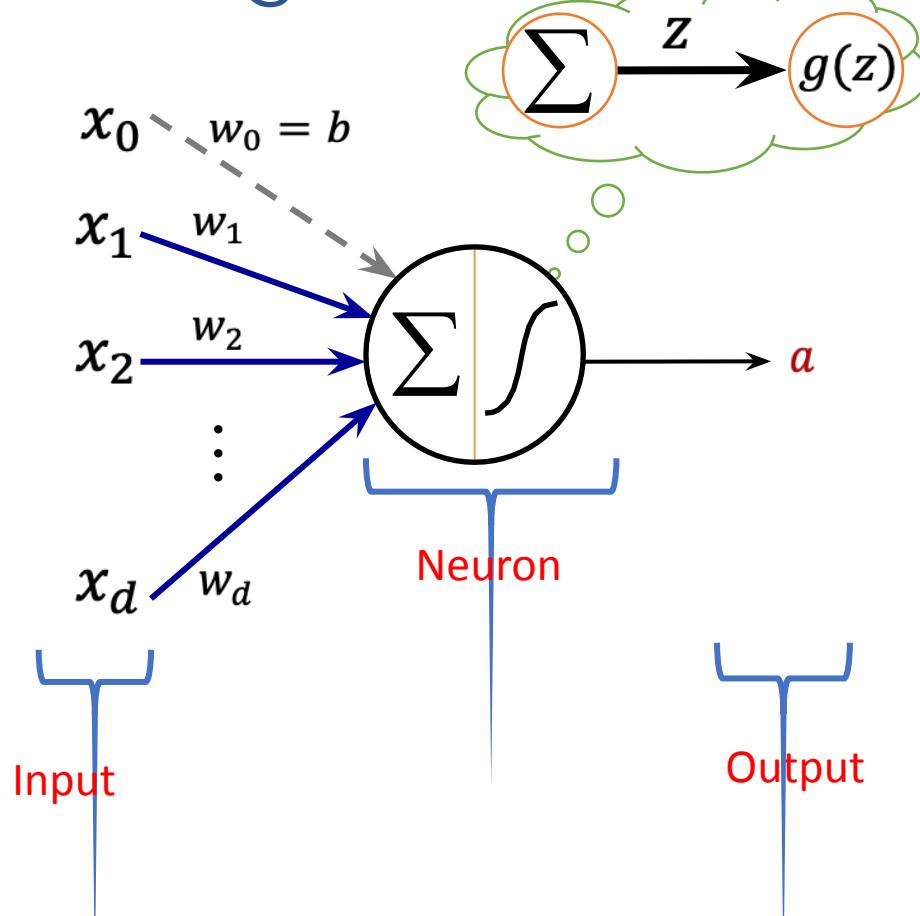
$$\mathbf{x} = [x_1 \ x_2 \ \dots \ x_d]^T$$

$$\mathbf{z} = b + \sum_{i=1}^d \theta_i x_i = [\boldsymbol{\theta}^T b] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$\mathbf{y} = g(\mathbf{z})$$

- Regression: $g(z) = z$
- Classification:
 - Binary: $g(z) = \sigma(z) = \frac{1}{1+e^{-z}}$
 - Multi-class (more on this later)

A Single Artificial Neuron



$$\mathbf{w} = [w_1 \ w_2 \ \dots \ w_d]^T \text{ and } \mathbf{x} = [x_1 \ x_2 \ \dots \ x_d]^T$$

$$\mathbf{z} = b + \sum_{i=1}^d w_i x_i = [\mathbf{w}^T \mathbf{b}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

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

Terminologies:-

\mathbf{x} : input, \mathbf{w} : weights, \mathbf{b} : bias

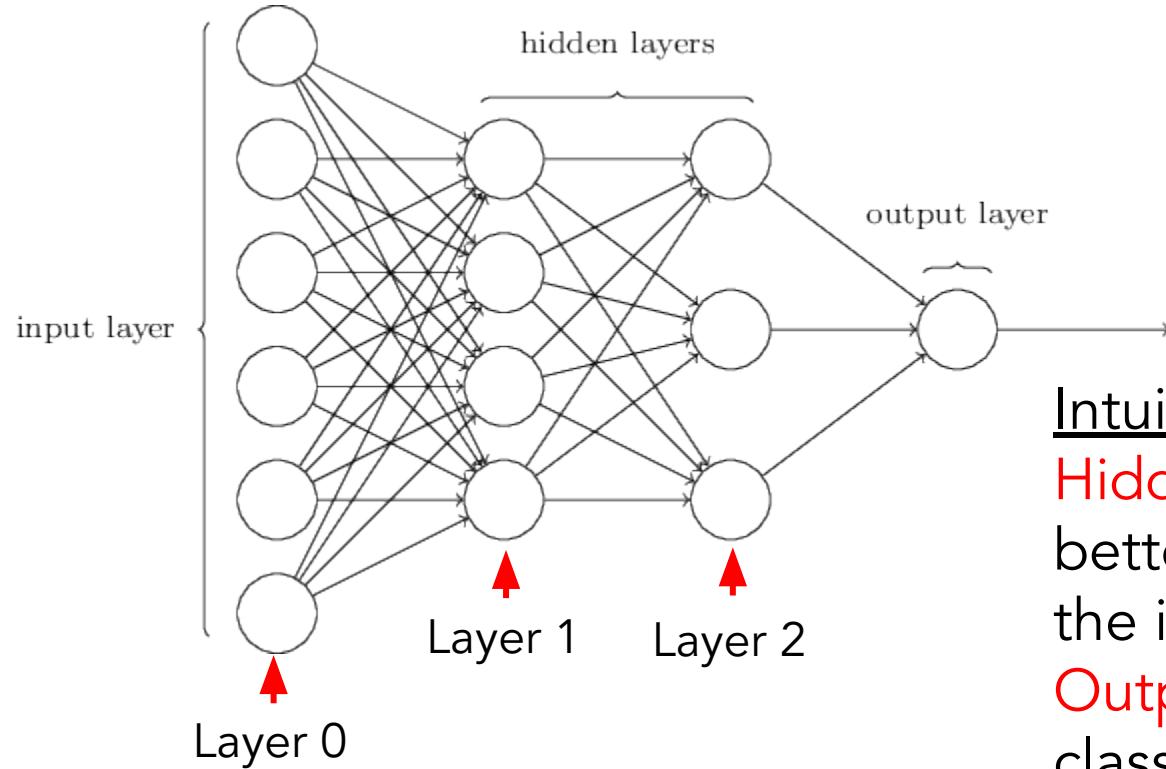
\mathbf{z} : pre-activation (input activation)

g : activation function

y : activation (output activation)

\mathbf{a} : activation at hidden units

Multilayer Perceptrons (MLP)



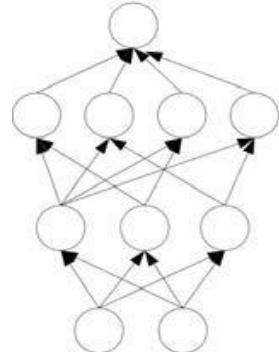
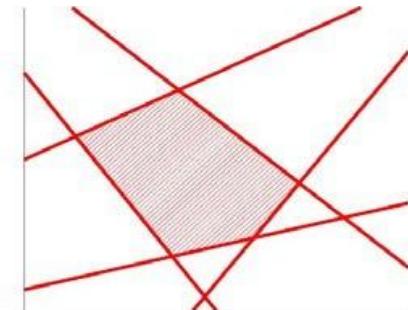
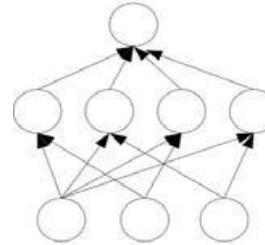
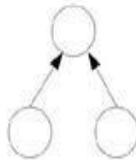
Intuition:-

Hidden Layer: Extracts better representation of the input data

Output layer: Does the classification

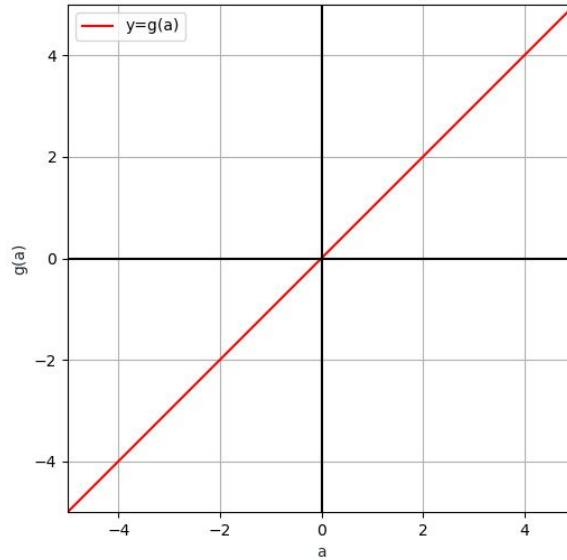
Beyond Single Layer

1 layer of trainable weights



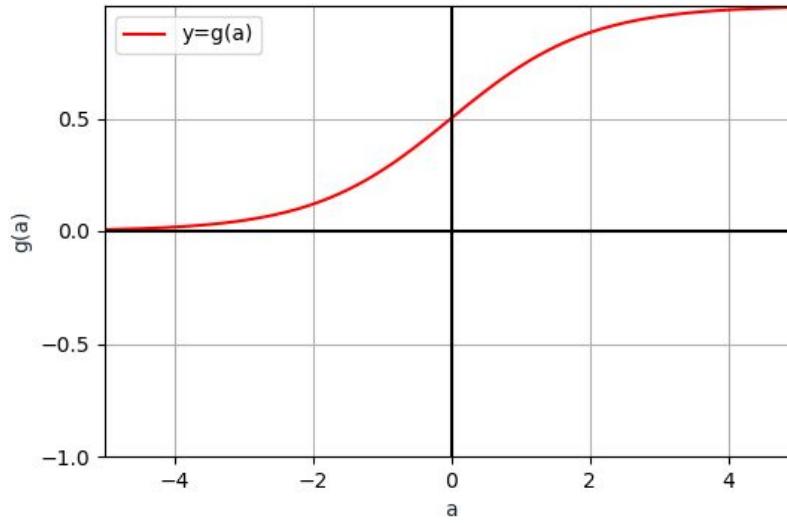
*What happens to bias?
What happens to variance?*

Common Activation Functions ($g(\cdot)$)



Linear activation function

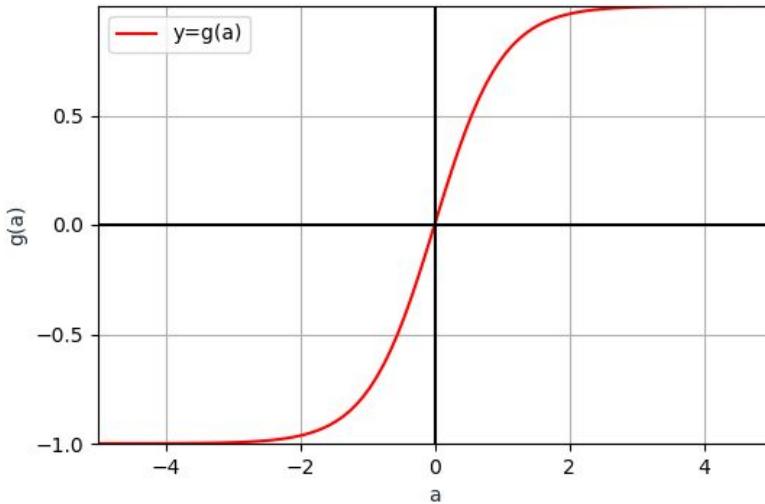
- $g(a) = a$
- Unbounded
- $g'(a) = 1$



Sigmoid activation function

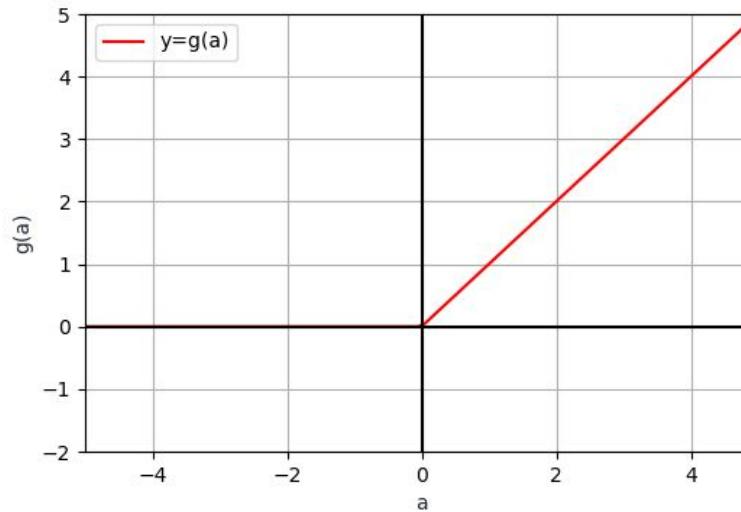
- $g(a) = \sigma(a) = \frac{1}{1+\exp(-a)}$
- Bounded (0, 1), Always positive
- $g'(a) = g(a)(1 - g(a))$

Common Activation Functions ($g(\cdot)$)



tanh activation function

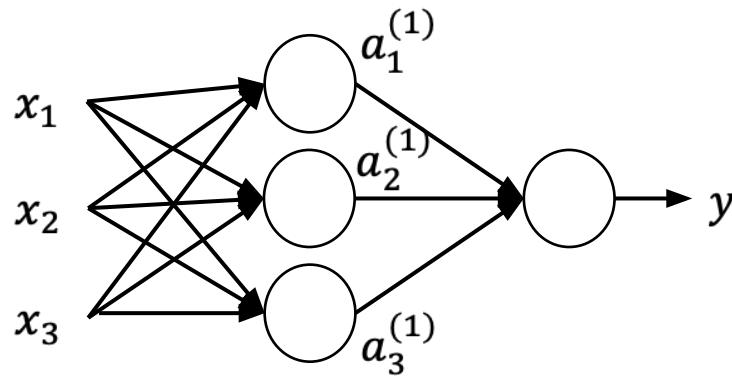
- $g(a) = \tanh(a) = \frac{\exp(a)-\exp(-a)}{\exp(a)+\exp(-a)}$
- Bounded (-1, 1)
- Can be positive or negative
- $g'(a) = 1 - g^2(a)$



ReLU activation function

- $g(a) = \max(0, a)$
- Bounded below by 0, But not upper-bounded
- $g'(a) = \begin{cases} 1, & a \geq 0 \\ 0, & a < 0 \end{cases}$

Why Non-Linear Activations**



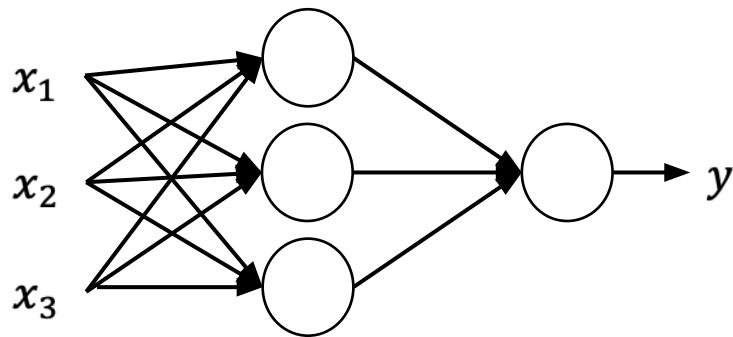
$$\begin{aligned} z^{(1)} &= W^{(1)}x + b^{(1)} \\ a^{(1)} &= g^{(1)}(z^{(1)}) \\ z^{(2)} &= W^{(2)}a^1 + b^{(2)} \\ y &= g^{(2)}(z^{(2)}) \end{aligned}$$

Assume if all are linear activation.
Task is to predict the yearly salary
 $y \in \$0, \dots, \$100,000$

$$\begin{aligned} a^{(1)} &= z^{(1)} = W^{(1)}x + b^{(1)} \\ y &= a^{(2)} = z^{(2)} = W^{(2)}a^1 + b^{(2)} \\ &= W^{(2)}(W^{(1)}x + b^{(1)}) + b^{(2)} \\ &= [W^{(2)}W^{(1)}]x + \{W^{(2)}b^{(1)} + b^{(2)}\} \\ &= W'x + b' \end{aligned}$$

$\Rightarrow y$ it is linear in x (the input)
 $\Rightarrow y$ is in \mathbb{R} (not what we need)

Vectorized Representation (*More on this later*)



$$z^{(1)} = W^{(1)}x + b^{(1)}$$

$$a^{(1)} = g^{(1)}(z^{(1)})$$

$$z^{(2)} = W^{(2)}a^{(1)} + b^{(2)}$$

$$y = g^{(2)}(z^{(2)})$$

$$z_1^{(1)} = W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)}$$

$$z_2^{(1)} = W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)}$$

$$z_3^{(1)} = W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)}$$

$3 \times 3, 1 \times 3$

$$a_1^{(1)} = g^{(1)}(z_1^{(1)})$$

$$a_2^{(1)} = g^{(1)}(z_2^{(1)})$$

$$a_3^{(1)} = g^{(1)}(z_3^{(1)})$$

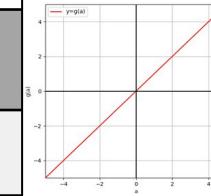
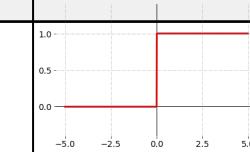
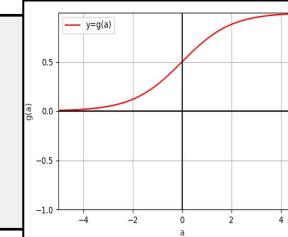
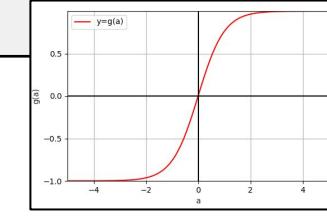
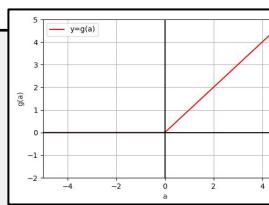
1×3

$3 \times 1, 1 \times 1$

$$z_1^{(2)} = W_{11}^{(2)}a_1^{(1)} + W_{12}^{(2)}a_2^{(1)} + W_{13}^{(2)}a_3^{(1)} + b^{(2)}$$

$$y = g^{(2)}(z^{(2)})$$

Common Activation Functions

Name	Function	Gradient	Graph
Linear	$y = g(a) = a$	$\frac{dy}{da} = 1$	
Binary step	$y = g(a) = \begin{cases} 1 & a \geq 0 \\ 0 & a < 0 \end{cases}$		
Sigmoid	$y = g(a) = \frac{1}{1 + e^{-a}}$		
Tanh	$y = g(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$		
ReLU	$y = g(a) = \begin{cases} a & a \geq 0 \\ 0 & a < 0 \end{cases}$		

Other Network Architectures



Pedestrian



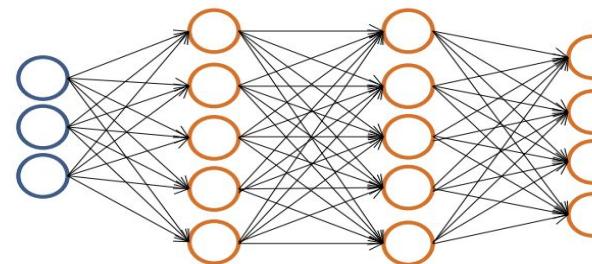
Car



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

Output Unit Activations

Regression: Linear

• Expected Output $\in \mathbb{R}$

$$\hat{y} = \mathbf{w}^T \mathbf{a} + b$$

Maximizing log-likelihood \Rightarrow

Minimizing squared error

$$J(\theta) = \|y - \hat{y}_\theta\|^2$$

Binary Classification: Sigmoid

Expected Output $\in [0,1]$

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{a} + b)$$

$$J(\theta) = -\log p(y|x)$$

$$= -\log \sigma((2y - 1)(\mathbf{w}^T \mathbf{a} + b))$$

Multi-Class Classification: Softmax

Expected Output $\in [0,1]^C$

Say, 3 classes.

- If output is Class 1, [1 0 0]
- If output is Class 2, [0 1 0]
- If output is Class 3, [0 0 1]

Need to produce a vector \hat{y} with $\hat{y}_i = p(y = i|x)$

$$\text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

$$\log \text{softmax}(z)_i = z_i - \log \sum_j \exp(z_j)$$

Loss Functions

Regression

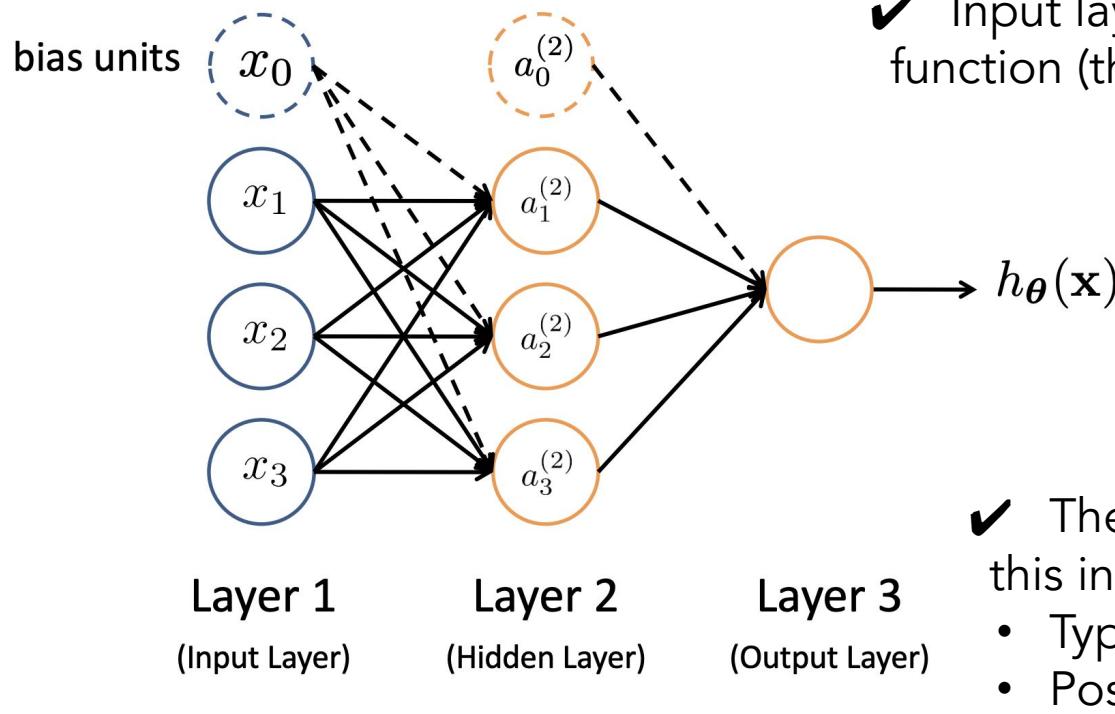
- Squared error: $L(y, \hat{y}) = (y - \hat{y})^2$

Classification

- Cross entropy: $L(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_k y_k \log \hat{y}_k$
- Easy to scale to k classes.

Forward Pass: Computing a Neural Network's Output

Neural Network – Forward Pass

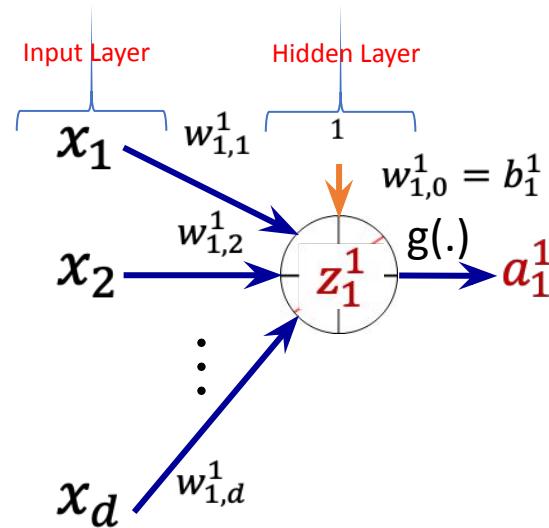


- ✓ Input layer units are set by some exterior function (think of these as sensors).
- ✓ Working forward through the network, the input function of each unit is applied to compute the input value.
- ✓ The activation function transforms this input function into a final value
 - Typically this is a nonlinear function
 - Possibly “Thresholding”

Feed-Forward Process

- Input layer units are set by some exterior function (think of these as sensors), which causes their output links to be activated at the specified level
- Working forward through the network, the input function of each unit is applied to compute the input value
 - Usually this is just the weighted sum of the activation on the links feeding into this node
- The activation function transforms this input function into a final value – Typically this is a nonlinear function, often a sigmoid function corresponding to the “threshold” of that node

Multilayer Neural Network



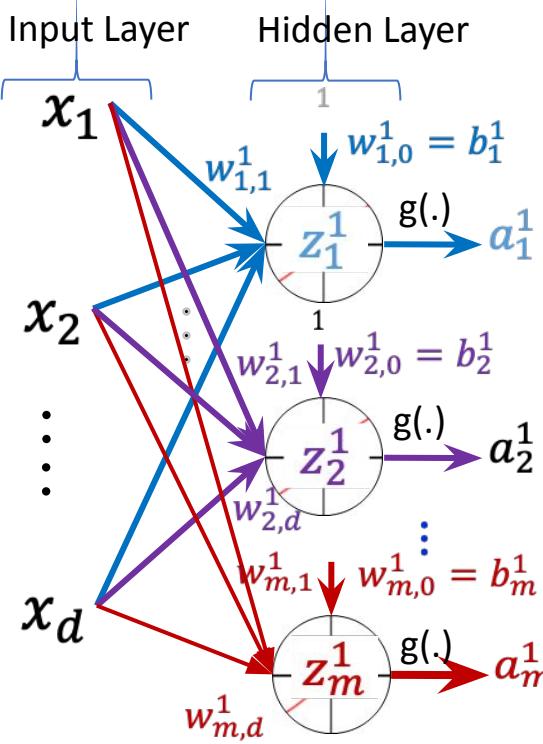
$$z_1^1 = b_1^1 + \sum_{i=1}^d w_{1,i}^1 x_i = [\mathbf{w}_1^1 \mathbf{b}_1^1] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$a_1^1 = g(z_1^1)$$

$$[x_1 \ x_2 \ \dots \ x_d]^T$$

$w_{i \leftarrow j}^l$ = Connection from $l - 1^{\text{th}}$ layer, j^{th} neuron
TO l^{th} layer, i^{th} neuron

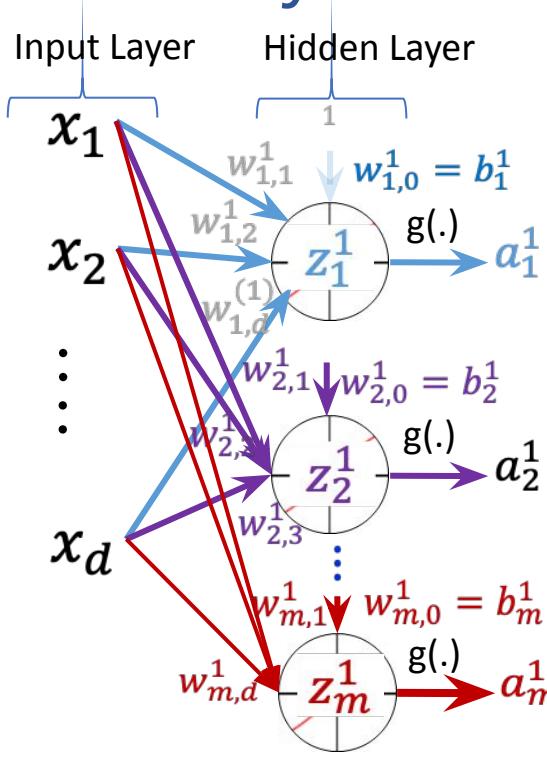
Multilayer Neural Network



$$\begin{aligned} a^{(0)} &= x \\ z^{(1)} &= \mathbf{w}^{(1)} \mathbf{a}^{(0)} \\ a^{(1)} &= g(z^{(1)}) \end{aligned}$$

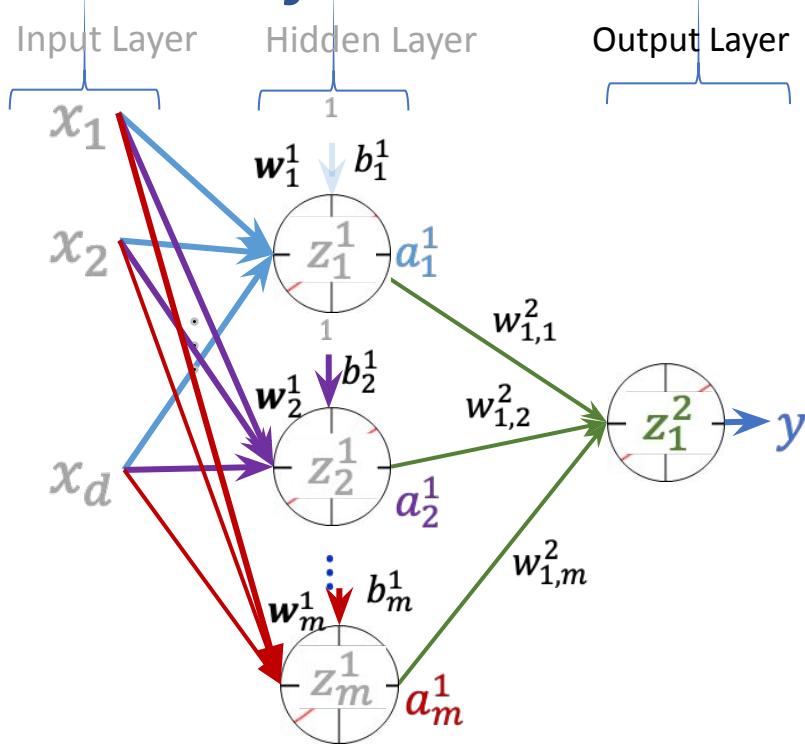
$W^1 : m \times n$ matrix
 $b^1 : m \times 1$ column vector
 $X : d \times 1$ column vector
 $Z^1 : m \times 1$ column vector
 $A^1 : m \times 1$ column vector

Multilayer Neural Network



$$\begin{aligned}
 a^{(0)} &= x \\
 z^{(1)} &= w^{(1)} a^{(0)} \\
 a^{(1)} &= g(z^{(1)}) \\
 \\
 \begin{bmatrix} z_1^1 \\ z_2^1 \\ \vdots \\ z_m^1 \end{bmatrix} &= \begin{bmatrix} w_1^1 & b_1^1 \\ w_2^1 & b_2^1 \\ \vdots & \vdots \\ w_m^1 & b_m^1 \end{bmatrix} \begin{bmatrix} x \\ 1 \end{bmatrix} \\
 z^1 &= [w^1 b^1] \begin{bmatrix} x \\ 1 \end{bmatrix} \\
 a^1 &= g(z^{(1)}) \\
 \\
 z_1^1 &= [w_1^1 b_1^1] \begin{bmatrix} x \\ 1 \end{bmatrix} \\
 z_2^1 &= [w_2^1 b_2^1] \begin{bmatrix} x \\ 1 \end{bmatrix} \\
 &\vdots \\
 z_M^1 &= [w_m^1 b_m^1] \begin{bmatrix} x \\ 1 \end{bmatrix} \\
 \\
 \begin{bmatrix} a_1^1 \\ a_2^1 \\ \vdots \\ a_m^1 \end{bmatrix} &= \begin{bmatrix} g(z_1^1) \\ g(z_2^1) \\ \vdots \\ g(z_m^1) \end{bmatrix}
 \end{aligned}$$

Multilayer Neural Network



Output Layer Pre-activation

$$z_1^{(2)} = [w_1^{(2)} \ b_1^{(2)}] \begin{bmatrix} a^{(1)} \\ 1 \end{bmatrix}$$

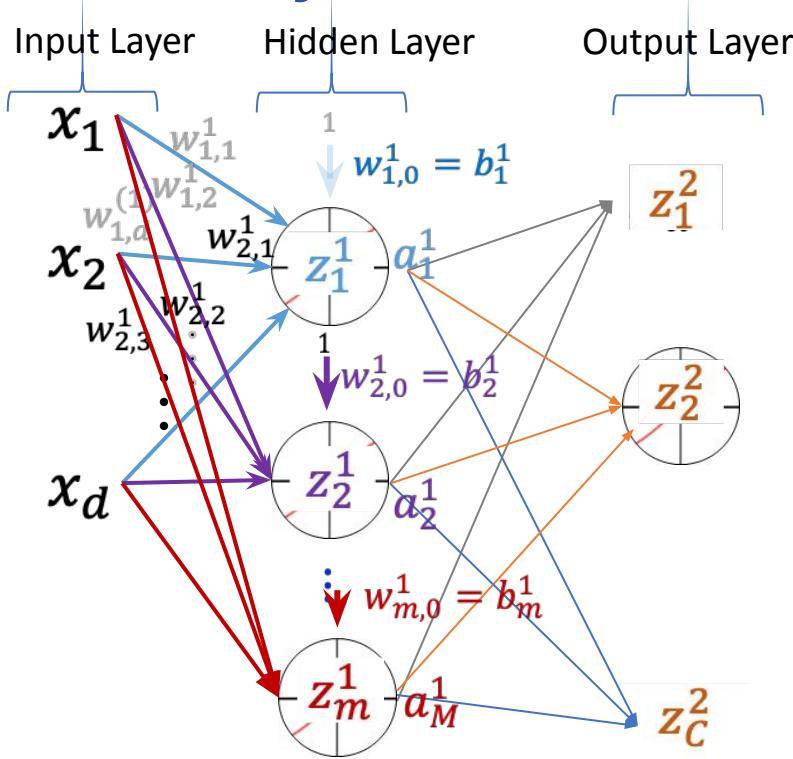
Output Layer Activation

$$y_1 = o(z_1^{(2)})$$

output

- Sigmoid for 2-class classification
- Softmax for multi-class classification
- Linear for regression

Multilayer Neural Network

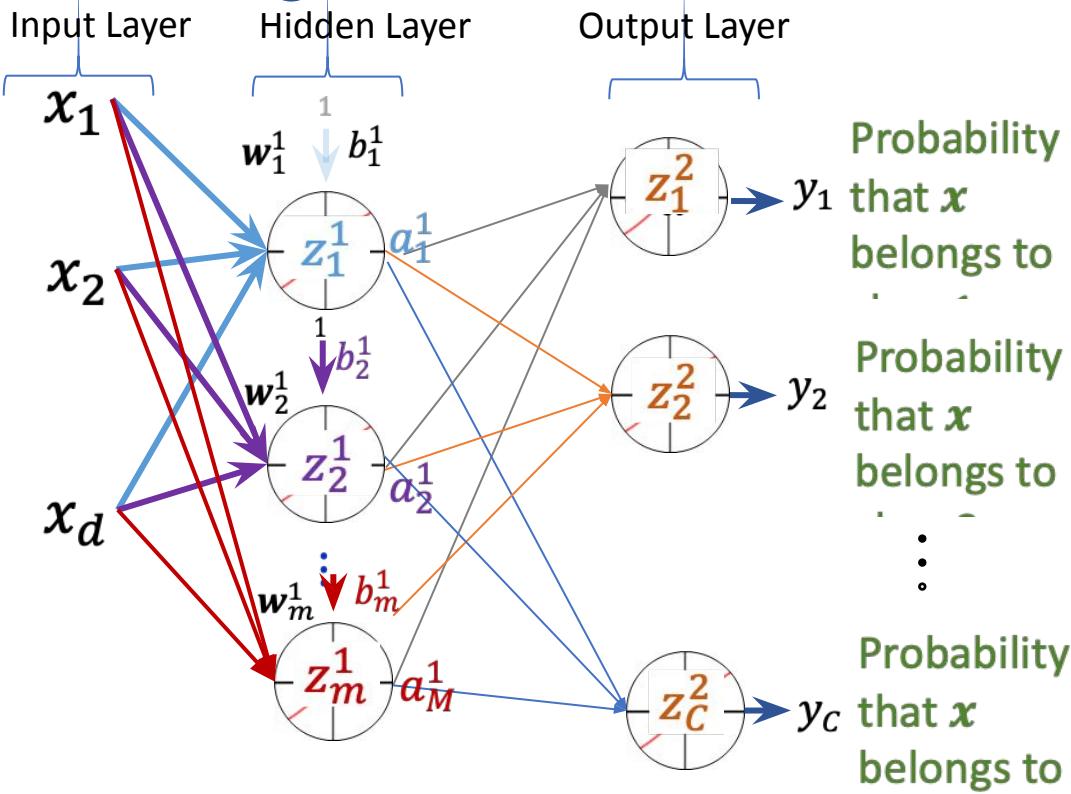


$$\rightarrow y_1 = o_1(z_1^{(2)}) = \frac{\exp(z_1^{(2)})}{\sum_c \exp(z_c^{(2)})}$$

$$\rightarrow y_2 = o_2(z_2^{(2)}) = \frac{\exp(z_2^{(2)})}{\sum_c \exp(z_c^{(2)})}$$

$$\rightarrow y_c = o_c(z_c^{(2)}) = \frac{\exp(z_c^{(2)})}{\sum_c \exp(z_c^{(2)})}$$

Training a Neural Network – Loss Function



Aim to maximize the probability corresponding to the correct class for any example x

$$\max y_c \equiv \max (\log y_c) \\ \equiv \min (-\log y_c)$$

$$\equiv - \sum_i \prod_{i=c} \log(y_i)$$

known as cross-entropy loss

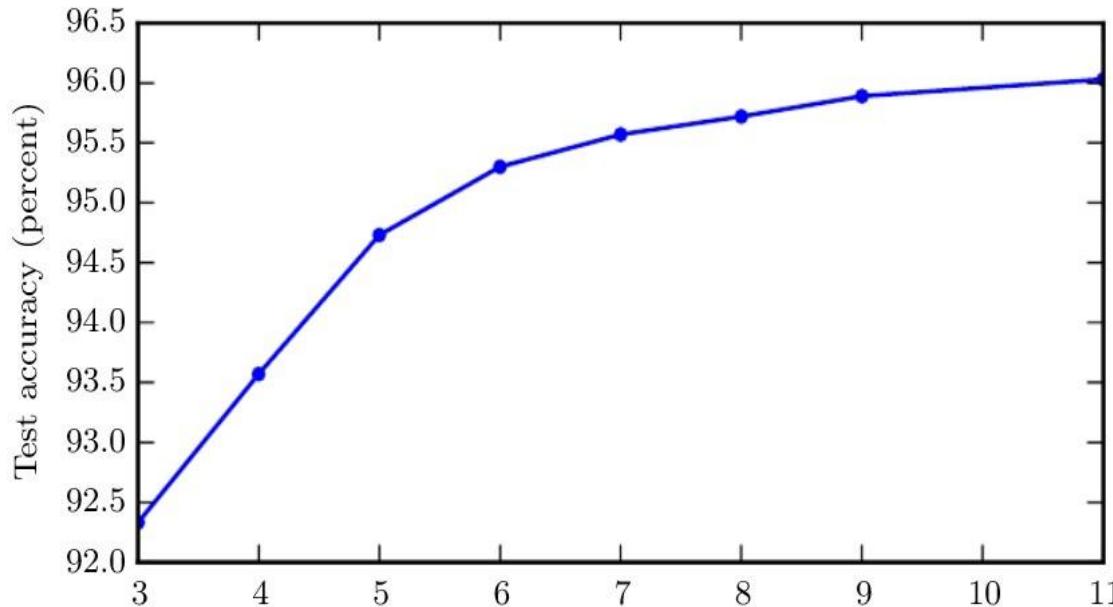
Universality

- Theoretical result [Cybenko, 1989]: 2-layer net with linear output with some squashing non-linearity in hidden units can approximate any continuous function over compact domain to arbitrary accuracy (given enough hidden units!)
- Implication: Regardless of function we are trying to learn, we know a large MLP can represent this function
- But not guaranteed that our training algorithm will be able to learn that function
- Gives no guidance on how large the network will be (exponential size in worst case)

A visual proof that neural nets can compute any function

- See Chapter 4 of [Neural Networks and Deep Learning](#) by Michael Nielsen
- <http://neuralnetworksanddeeplearning.com/chap4.html>

Advantages of Depth

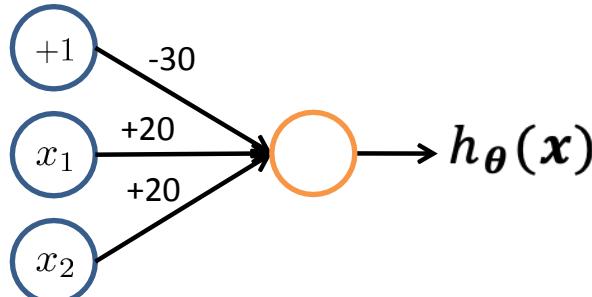


Representing Boolean Functions

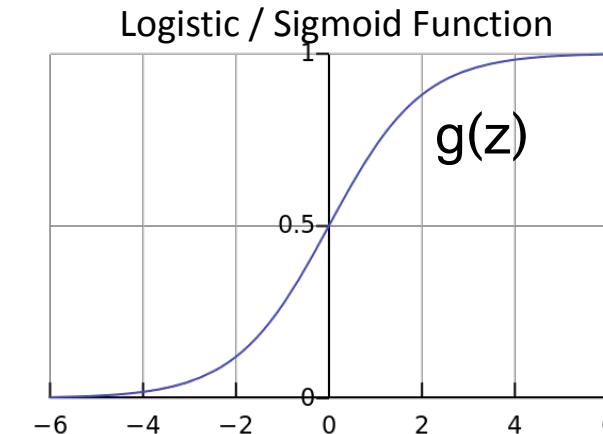
Simple example: AND

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

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

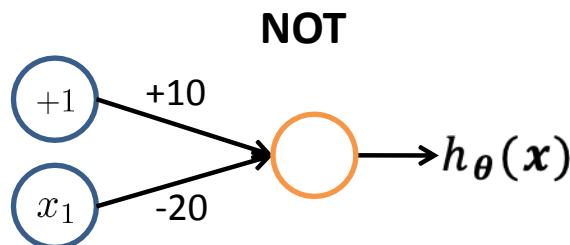
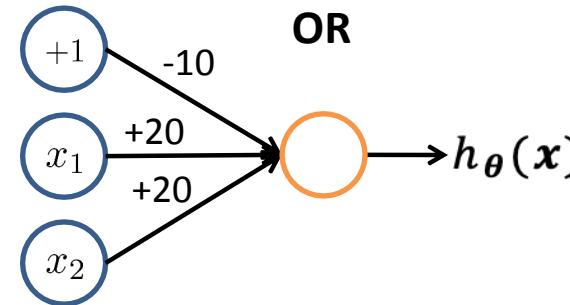
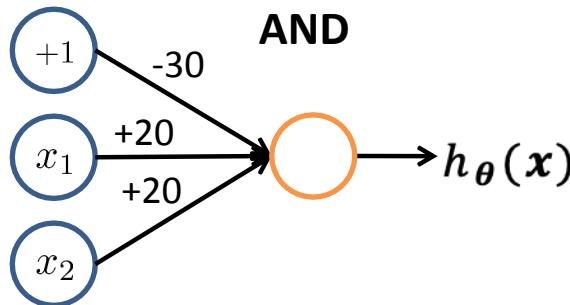


$$h_{\theta}(x) = g(-30 + 20x_1 + 20x_2)$$

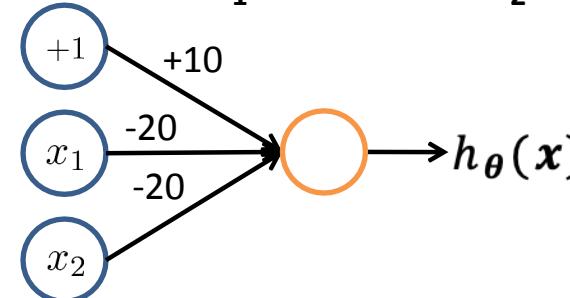


x_1	x_2	$h_{\theta}(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

Representing Boolean Functions

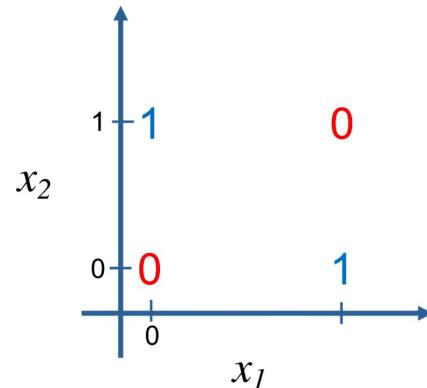
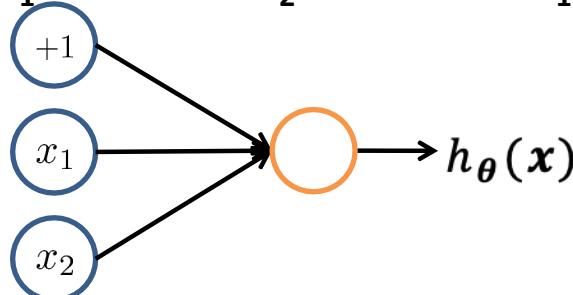


NOR: $(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$

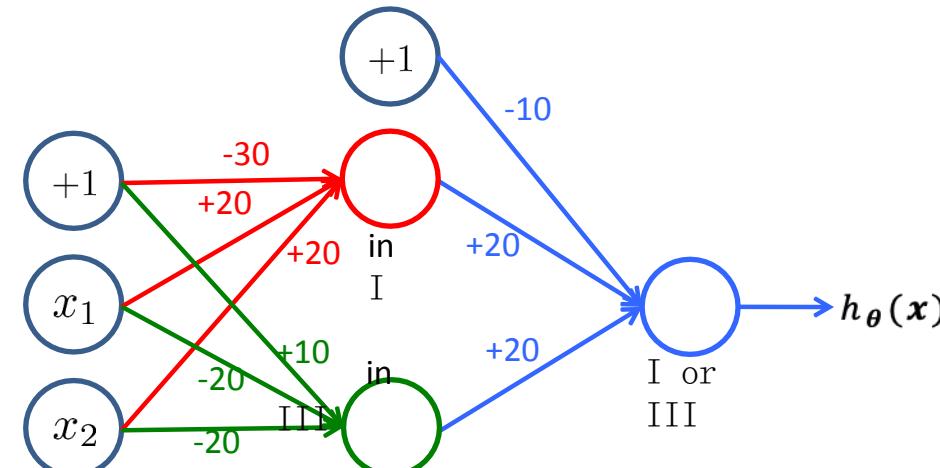
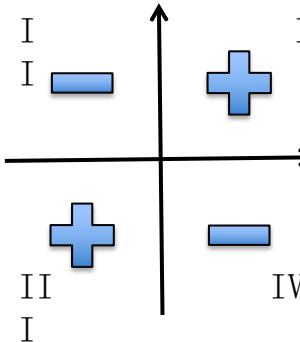
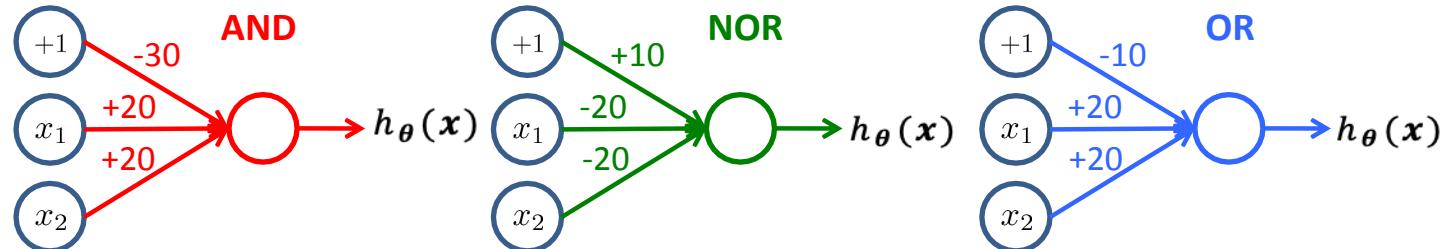


Representing Boolean Functions

XOR: $(x_1 \text{ AND } (\text{NOT } x_2)) \text{ OR } ((\text{NOT } x_1) \text{ AND } x_2)$

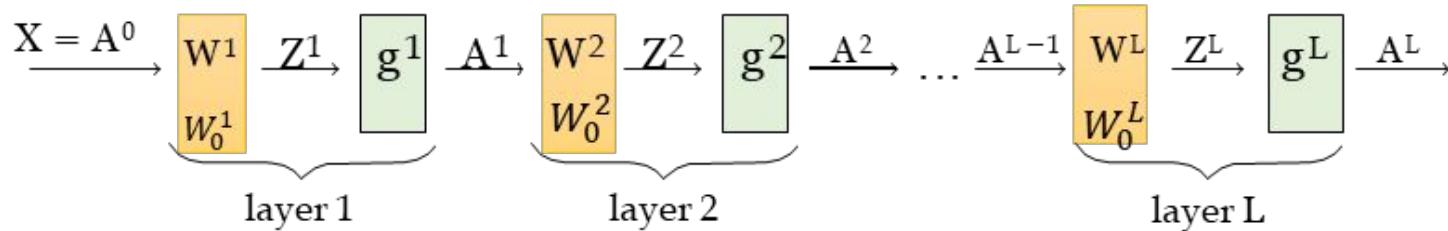


Combining Representations to Create Non Linear Functions



Feedforward Networks and Backpropagation

Forward Pass (Block Diagram)



Hidden layer pre-activation:

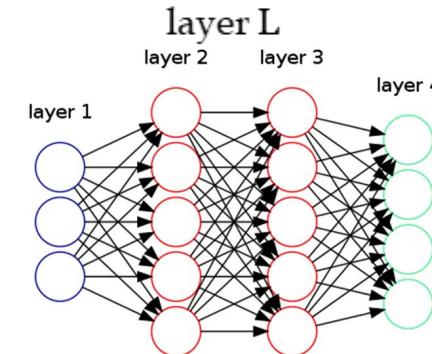
$$\text{For } l = 1, \dots, L; \mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{W}_0^{(l)}$$

Hidden layer activation:

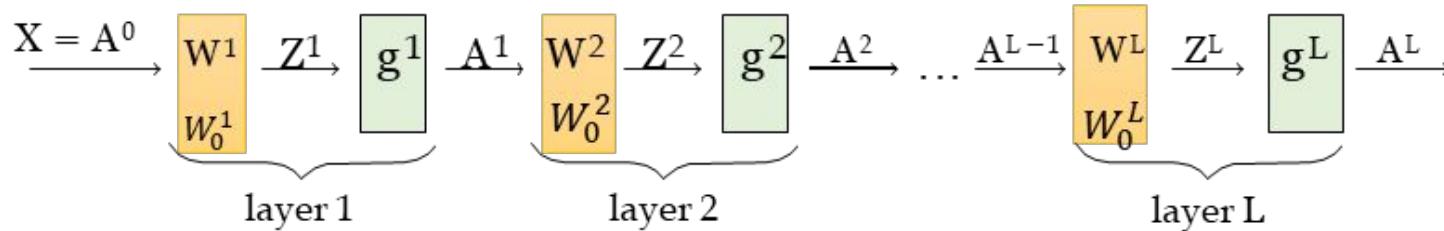
$$\text{For } l = 1, \dots, L - 1; \mathbf{a}^{(l)} = g^l(\mathbf{z}^{(l)})$$

Output layer activation:

$$\text{For } l = L; \mathbf{y} = \mathbf{a}^{(L)} = o(\mathbf{z}^{(L)}) = \mathbf{f}(\mathbf{x}, \boldsymbol{\theta})$$

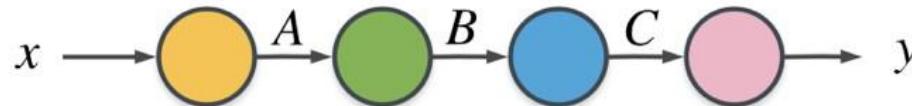


Error back-propagation



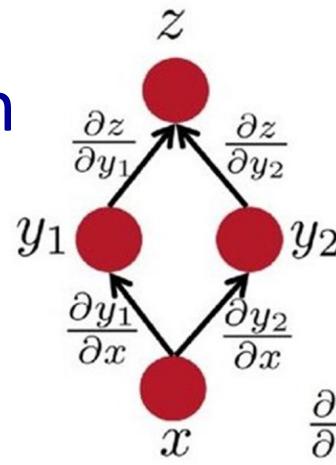
1. Compute Loss
2. Compute the derivative of the L w.r.t. the final output of the network A^L
3. Compute the derivative of L w.r.t. the pre-activation Z^L
4. Compute the derivative of L wrt W^L
5. Then compute the derivative of the L w.r.t. the final output of the network A^{L-1}
6. Then compute the derivative of L w.r.t. the pre-activation Z^{L-1}
7. Compute the derivative of L wrt W^{L-1}
-

Chain Rule



$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial C} \times \frac{\partial C}{\partial B} \times \frac{\partial B}{\partial A} \times \frac{\partial A}{\partial x}$$

Multiple Path



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

Exercise:

$$y_1 = w_1 x$$

$$y_2 = w_2 x$$

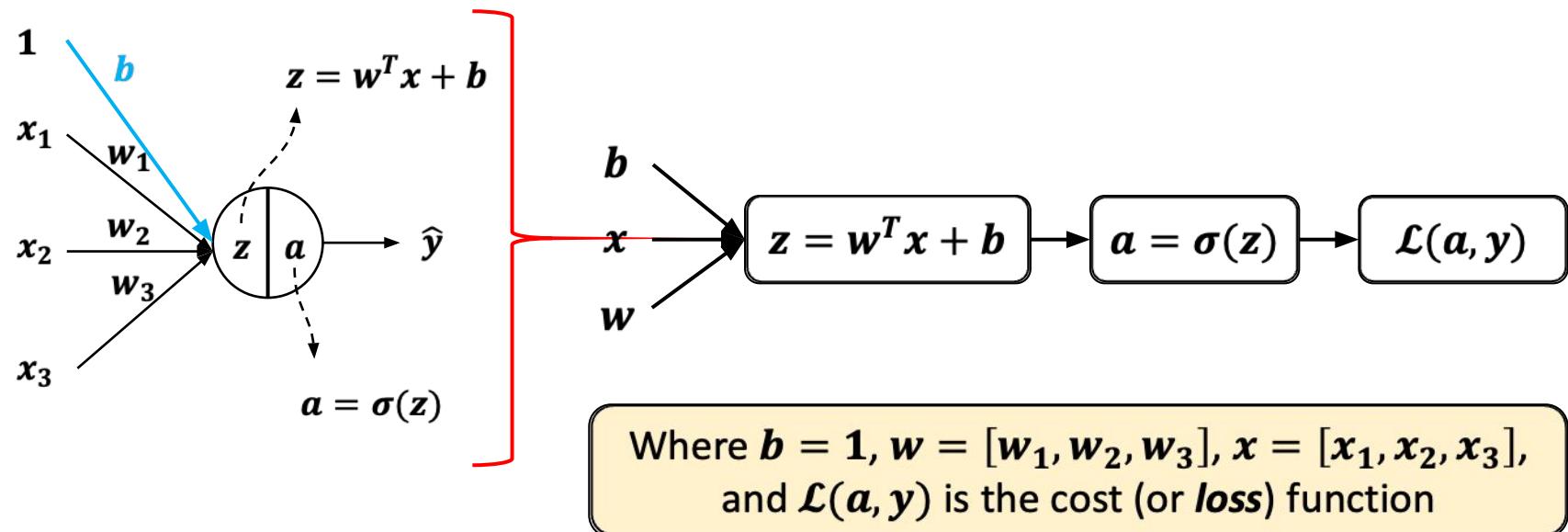
$$z = y_1 + y_2$$

What is $\frac{\partial z}{\partial x}$?

Backpropagation is just repeated application of the chain rule.

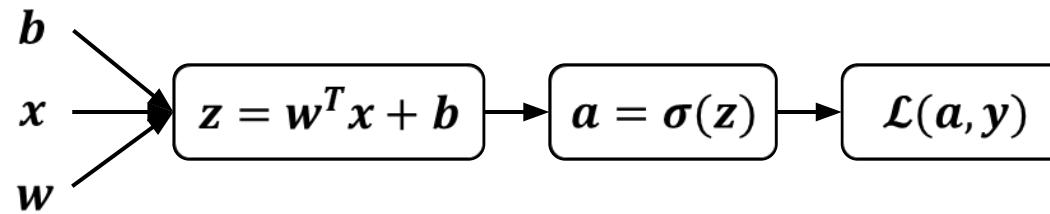
The Computation Graph of Logistic Regression

Let us translate logistic regression into a computation graph

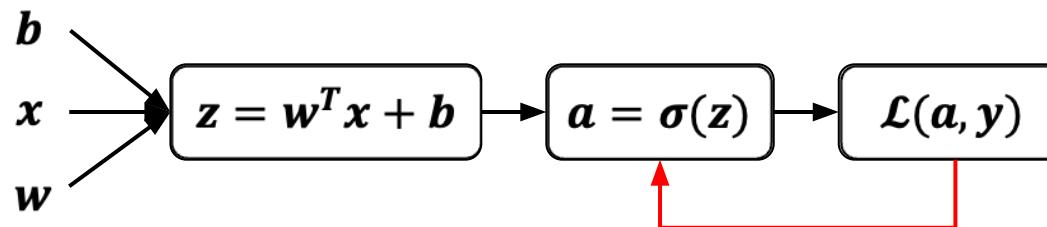


Forward Propagation

- The loss function can be computed by moving from left to right



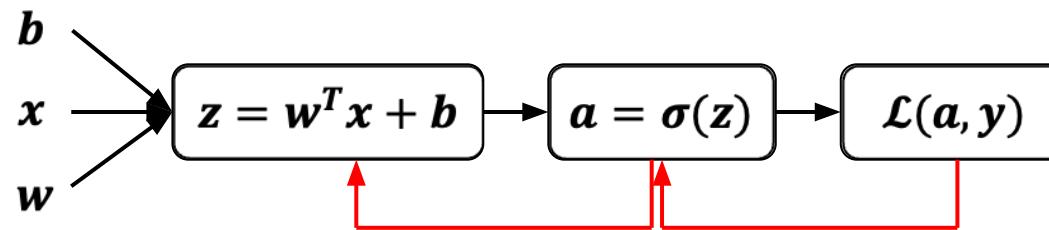
Backward Propagation



Partial derivative of \mathcal{L} with respect to a

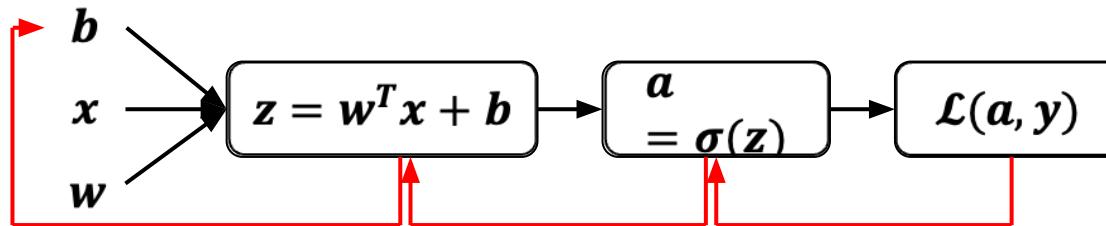
$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial a} &= \frac{\partial}{\partial a} (-y \log(a) - (1-y) \log(1-a)) \\ &= \frac{-y}{a} + \frac{(1-y)}{(1-a)}\end{aligned}$$

Backward Propagation



$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial z} &= \frac{\partial \mathcal{L}}{\partial a} \times \frac{\partial a}{\partial z} = \left(\frac{-y}{a} + \frac{(1-y)}{1-a} \right) \times \frac{\partial a}{\partial z} = \left(\frac{-y}{a} + \frac{(1-y)}{1-a} \right) \times a(1-a) \\ &= a - y\end{aligned}$$

Backward Propagation

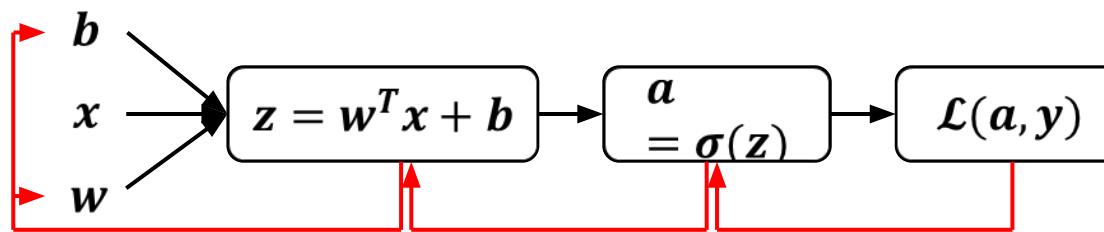


$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial b}$$

Partial derivative of \mathcal{L} with respect to b

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial b} = (a - y) \times \frac{\partial z}{\partial b} = (a - y) \times 1 = (a - y)$$

Backward Propagation



$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w} = (a - y) \times \frac{\partial z}{\partial w} = (a - y)x$$

Backward Propagation: Summary

- Here is the summary of the gradients in logistic regression:

$$\mathbf{dz} = \frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial a}{\partial z}$$

$$\mathbf{db} = \frac{\partial \mathcal{L}}{\partial b} = a - y$$

$$\mathbf{dw} = \frac{\partial \mathcal{L}}{\partial w} = (a - y)x$$

Introduction

Activation

Forward Pass

Backpropagation

Backward Pass MLP

Optimizing the Neural Network

$$z_1^{(2)} = w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3$$

$$a_1^{(2)} = g(z_1^{(2)})$$

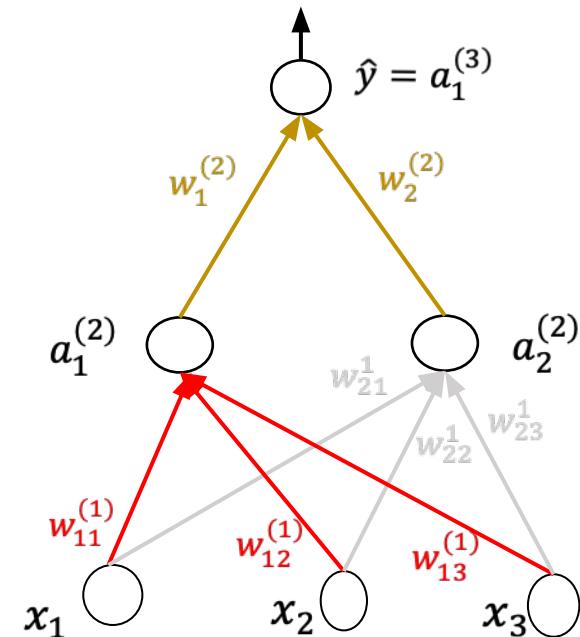
$$z_2^{(2)} = w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3$$

$$a_2^{(2)} = g(z_2^{(2)})$$

$$z_1^{(3)} = w_1^{(2)} a_1^{(2)} + w_2^{(2)} a_2^{(2)}$$

$$\hat{y} = g(z_1^{(3)}) = a_1^{(3)}$$

- Loss $L = (\hat{y} - y)^2$
- Compute $\frac{\partial L}{\partial w_1^{(2)}}, \frac{\partial L}{\partial w_2^{(2)}}$
- Compute $\frac{\partial L}{\partial w_{11}^{(2)}}, \frac{\partial L}{\partial w_{12}^{(2)}}, \dots, \frac{\partial L}{\partial w_{23}^{(2)}}$



Top Layer Updates

Output $\hat{y} = a_1^{(3)} = g(w_1^{(2)} a_1^{(2)} + w_2^{(2)} a_2^{(2)})$, Loss $L = (\hat{y} - y)^2$

Compute $\frac{\partial L}{\partial w_1^{(2)}}, \frac{\partial L}{\partial w_2^{(2)}}$

$$\frac{\partial L}{\partial w_1^{(2)}} = \frac{\partial (\hat{y} - y)^2}{\partial w_1^{(2)}}$$

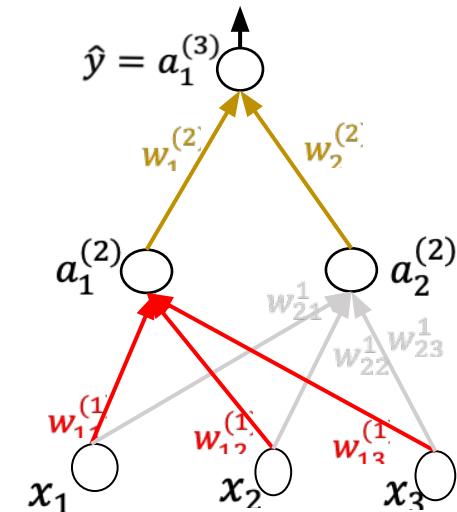
$$= 2(\hat{y} - y) \frac{\partial (g(w_1^{(2)} a_1^{(2)} + w_2^{(2)} a_2^{(2)}))}{\partial w_1^{(2)}}$$

$$= 2(\hat{y} - y) a_1^{(3')} a_1^{(2)}$$

Weight update (Gradient Descent): $w_1^{(2)} = w_1^{(2)} - \eta \frac{\partial L}{\partial w_1^{(2)}}$

Assume $\delta = 2(\hat{y} - y) a_1^{(3')}$, the updates become

- $w_1^{(2)} = w_1^{(2)} - \eta \delta a_1;$ $w_2^{(2)} = w_2^{(2)} - \eta \delta a_2;$



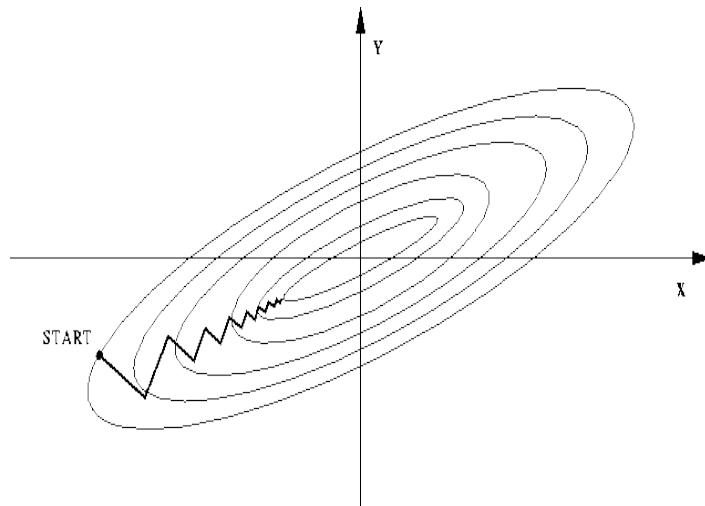
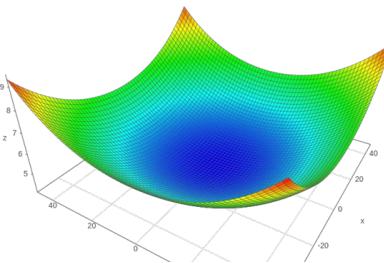
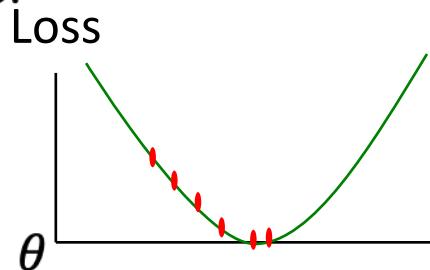
Background: Gradient descent

Partial derivatives give us the slope (i.e. direction to move) in that dimension

Approach:

- pick a starting point (θ)
- repeat:
 - pick a dimension
 - move a small amount in that dimension towards decreasing loss (using the derivative)

$$\theta_j = \theta_j - \eta \frac{d}{d\theta_j} \text{Loss}(\theta)$$



Next Layer Updates

Loss $L = (\hat{y} - y)^2$, Compute $\frac{\partial L}{\partial w_1^{(2)}}, \frac{\partial L}{\partial w_2^{(2)}}$

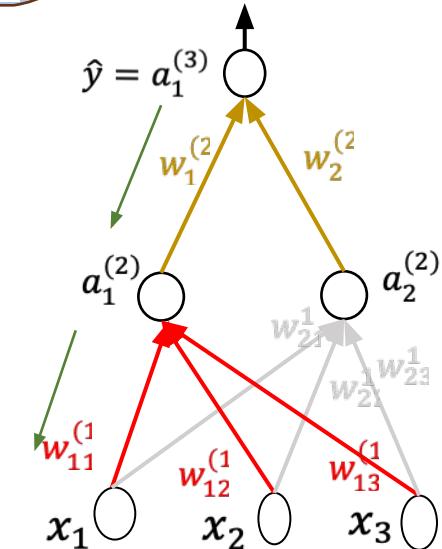
- $$\frac{\partial L}{\partial w_{11}^{(1)}} = \frac{\partial(\hat{y} - y)^2}{\partial w_{11}^{(1)}}$$

$$= 2(\hat{y} - y) \frac{\partial(g(w_1^{(2)} a_1^{(2)} + w_2^{(2)} a_2^{(2)}))}{\partial w_{11}^{(1)}}$$

$$= 2(\hat{y} - y) a_1^{3'} \left(w_1^2 \frac{\partial(g(w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3))}{\partial w_{11}^{(1)}} + 0 \right)$$

$$= 2(\hat{y} - y) a_1^{3'} w_1^{(2)} g'(z_1) x_1$$

$$= 2(\hat{y} - y) a_1^{3'} w_1^{(2)} a'_1 x_1$$



Next Layer Updates

$$\hat{y} = w_1^2 a_1 + w_2^2 a_2$$

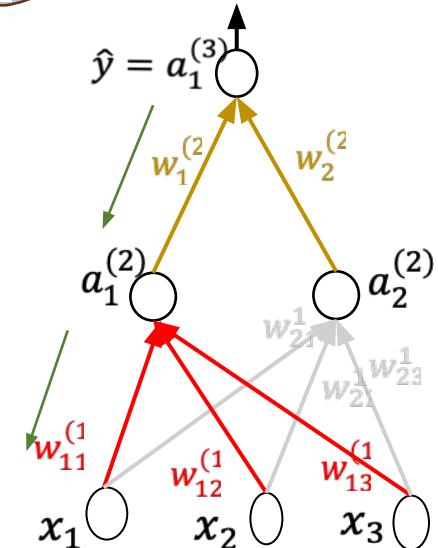
$$a_1 = g(z_1)$$

$$z_1 = w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3$$

- Loss $L = (\hat{y} - y)^2$, Compute $\frac{\partial L}{\partial w_1^{(2)}}, \frac{\partial L}{\partial w_2^{(2)}}$

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \frac{\partial(\hat{y}-y)^2}{\partial w_{11}^{(1)}} = 2(\hat{y} - y) a_1^{3'} w_1^{(2)} a_1^{(2)' } x_1$$

- Weight Update $w_{11}^{(1)} = w_{11}^{(1)} - \eta \frac{\partial L}{\partial w_{11}^{(1)}}$
 $w_{11}^{(1)} = w_{11}^{(1)} - \eta \delta a_1^{3'} w_1^{(2)} a_1^{(2)' } x_1$



All Updates Together Now

$$\text{Loss } L = (\hat{y} - y)^2$$

- For top layer $\frac{\partial L}{\partial w_1^{(2)}} = \delta a_1^{(1)}, \frac{\partial L}{\partial w_2^{(2)}} = \delta a_2^{(1)}$

$$\frac{\partial L}{\partial w_i^{(2)}} = \delta \quad a_i^{(1)}$$

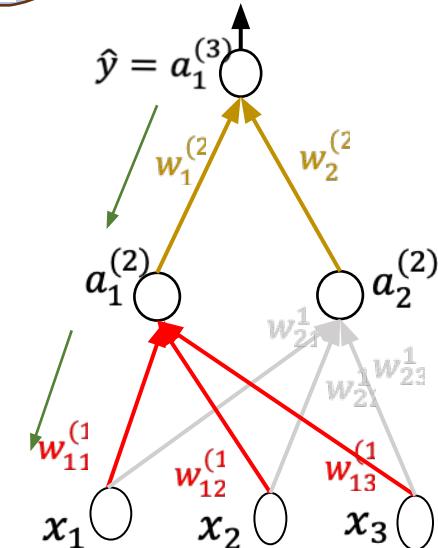
Local Error Local Input

- Next Layer $\frac{\partial L}{\partial w_{ij}^{(1)}} = \delta w_i^{(2)} a'_i x_j$

Assume, $\delta_i = \delta w_i^{(2)} a'_i$. (δ_j already contains the error term)

$$\frac{\partial L}{\partial w_{ij}^{(1)}} = \delta_i a'_j$$

Local Error Local Input



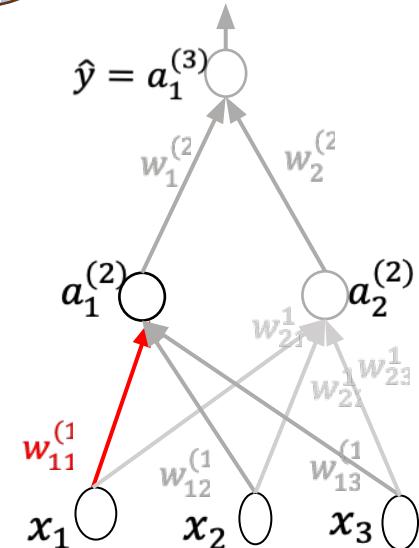
All Updates Together Now

$$\text{Loss } L = (\hat{y} - y)^2$$

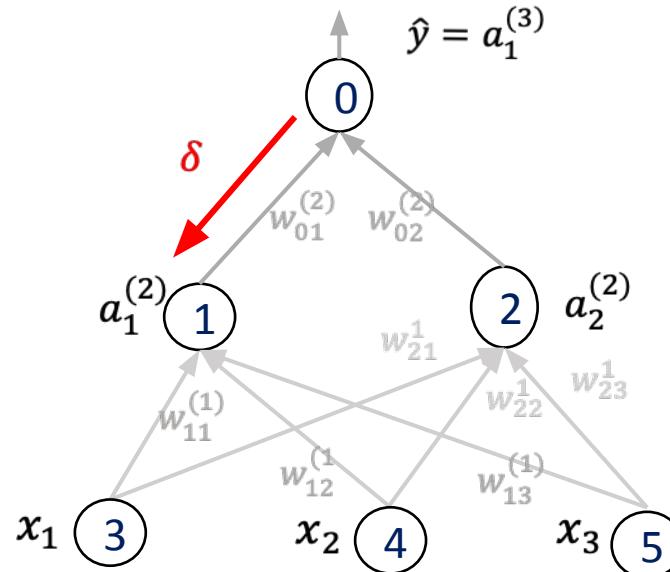
- Let w_{ij} be the connection from node j to node i

$$\frac{\partial L}{\partial w_{ij}^{(l)}} = \delta_i^{(l+1)} a_j^{(l)}$$

δ_i is the local error (going from i backwards) and a_j is the local input at node j .



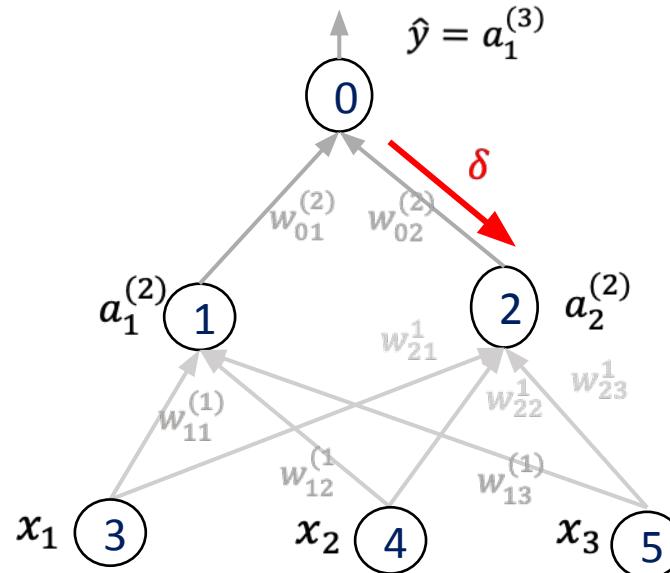
Cleaning Up: Top Layer



Local error from 0: $\delta_0 = 2(\hat{y} - y)a_1^{(3)'}$, Local input from 1: $a_1^{(1)}$

$$\frac{\partial L}{\partial w_{01}^{(l)}} = \delta_0 a_1^{(1)}$$

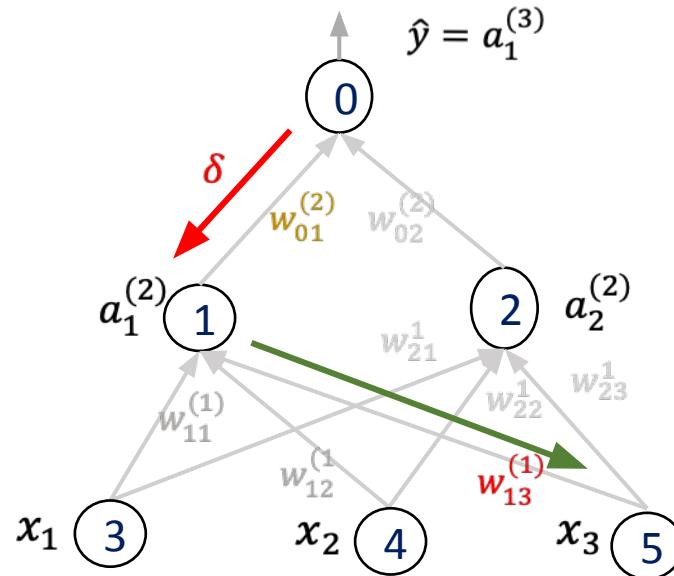
Cleaning Up: Top Layer



Local error from 0: $\delta_0 = 2(\hat{y} - y)a_1^{(3)'}$, Local input from 1: $a_2^{(1)}$

$$\frac{\partial L}{\partial w_{01}^{(2)}} = \delta_0 a_2^{(1)}$$

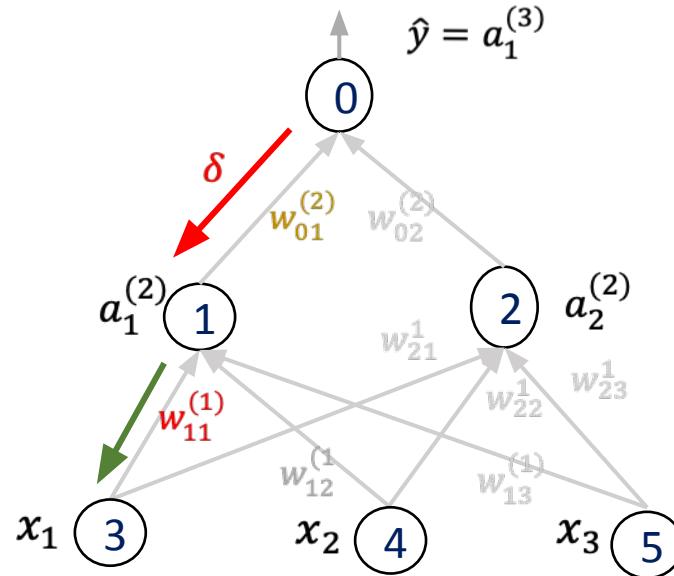
Cleaning Up: Next Layer



Local error from 1: $\delta_1 = (\delta)(w_{01}^{(2)})' a_1^{(2)'}$, Local input from 5: x_3

$$\frac{\partial L}{\partial w_{13}^{(1)}} = \delta_1 x_3$$

Cleaning Up: Next Layer



Local error from 1: $\delta_1 = (\delta)(w_{01}^{(2)}) \mathbf{a}_1^{(2)'}$, Local input from 3: x_1

$$\frac{\partial L}{\partial w_{11}^{(1)}} = \delta_1 x_1$$

Lets Vectorize

- Let $W^{(2)} = \begin{bmatrix} w_{01}^{(2)} \\ w_{02}^{(2)} \end{bmatrix}$,
- Let $W^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{21}^{(1)} \\ w_{12}^{(1)} & w_{22}^{(1)} \\ w_{13}^{(1)} & w_{23}^{(1)} \end{bmatrix}$
- Let $A^{(1)} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, A^{(2)} = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$

Forward computation (Vectorized)

1. Compute $Z^{(1)} = A^{(1)^\top} W^{(1)}$
2. Applying element-wise non-linearity $A^{(2)} = g(Z^{(2)})$
3. Compute Output $Z^{(2)} = A^{(2)^\top} W^{(2)}, \hat{y} = g(Z^{(2)})$
4. Compute Loss on example $L = (\hat{y} - y)^2$

Backward Pass (Vectorized)

1. Top: Compute δ
2. Gradient with respect to $W^{(2)} = \delta(A^{(2)})$
3. Compute $\delta_1 = (W^{(2)^\top} \delta) \odot A^{(2)'} \text{ (Hadamard product)}$
4. Gradient w.r.t $W^{(1)} = \delta_1(A^{(1)})$
5. Update $W^{(2)} := W^{(2)} - \eta \delta \ (A^{(2)})$
6. Update $W^{(1)} := W^{(1)} - \eta \delta_1(A^{(1)})$

Backpropagation

Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$

(Used to accumulate gradient)

For each training instance (\mathbf{x}_i, y_i) :

Set $\mathbf{a}^{(1)} = \mathbf{x}_i$

Compute $\{\mathbf{a}^{(2)}, \dots, \mathbf{a}^{(L)}\}$ via forward propagation

Compute $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y_i$

Compute errors $\{\boldsymbol{\delta}^{(L-1)}, \dots, \boldsymbol{\delta}^{(2)}\}$

Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n} \Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

$D^{(l)}$ is the matrix of partial derivatives of $J(\Theta)$

Note: Can vectorize $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ as $\Delta^{(l)} = \Delta^{(l)} + \boldsymbol{\delta}^{(l+1)} \mathbf{a}^{(l)\top}$

Training a NN via Gradient Descent with Backprop

Given: training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

Initialize all $\Theta^{(l)}$ randomly (NOT to 0!)

Loop // each iteration is called an epoch

Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$ (Used to accumulate gradient)

For each training instance (\mathbf{x}_i, y_i) :

Set $\mathbf{a}^{(1)} = \mathbf{x}_i$

Compute $\{\mathbf{a}^{(2)}, \dots, \mathbf{a}^{(L)}\}$ via forward propagation

Compute $\delta^{(L)} = \mathbf{a}^{(L)} - y_i$

Compute errors $\{\delta^{(L-1)}, \dots, \delta^{(2)}\}$

Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n} \Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

Update weights via gradient step $\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \alpha D_{ij}^{(l)}$

Until weights converge or max #epochs is reached

Training a Neural Network

- Pick a network architecture (connectivity pattern between nodes)
- # input units = # of features in dataset
- # output units = # classes
- Choose number of hidden layers and number of units in each layer.

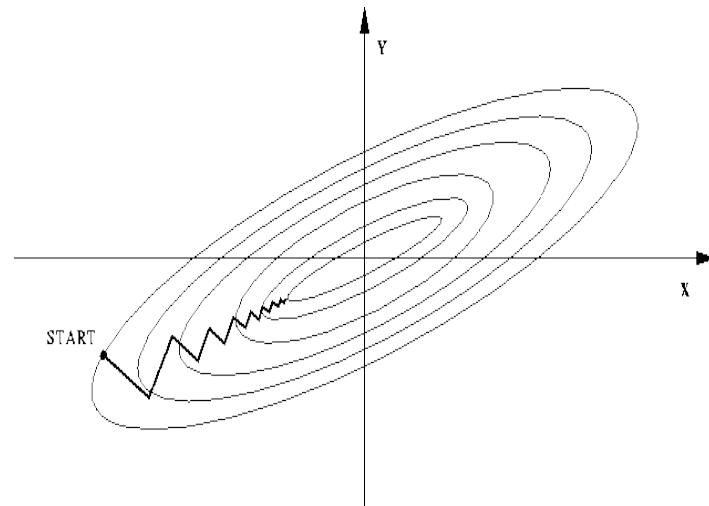
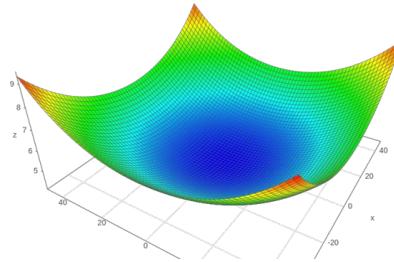
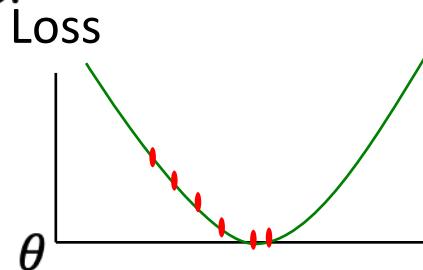
Gradient descent

Partial derivatives give us the slope (i.e. direction to move) in that dimension

Approach:

- pick a starting point (θ)
- repeat:
 - pick a dimension
 - move a small amount in that dimension towards decreasing loss (using the derivative)

$$\theta_j = \theta_j - \eta \frac{d}{d\theta_j} \text{Loss}(\theta)$$



Batch, Stochastic and Minibatch

- Optimization algorithms that use the entire training set to compute the gradient are called batch or deterministic gradient methods. Ones that use a single training example for that task are called stochastic or online gradient methods
- Most of the algorithms we use for deep learning fall somewhere in between!
- These are called minibatch or minibatch stochastic methods

Batch, Stochastic and Mini-batch Stochastic Gradient Descent

Algorithm 1 Batch Gradient Descent at Iteration k

Require: Learning rate ϵ_k

Require: Initial Parameter θ

- 1: **while** stopping criteria not met **do**
- 2: Compute gradient estimate over N examples:
- 3: $\hat{\mathbf{g}} \leftarrow +\frac{1}{N} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
- 4: Apply Update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$
- 5: **end while**

Mini-batch

Algorithm 8.1 Stochastic gradient descent (SGD) update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

end while

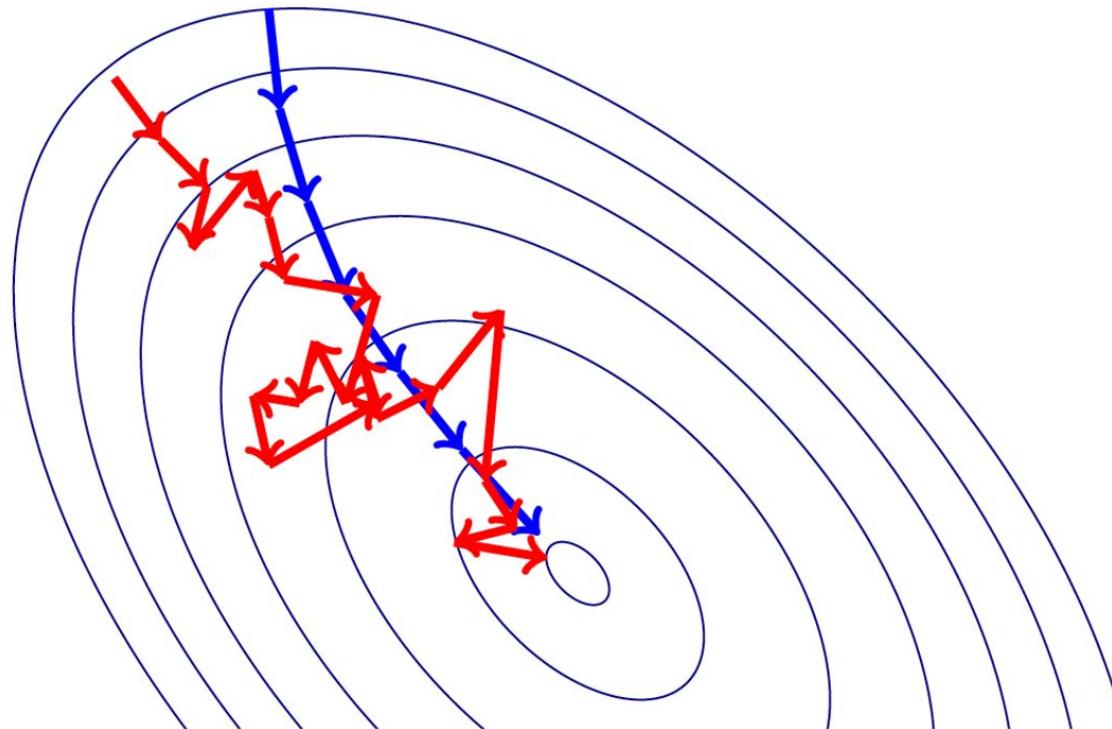
Algorithm 2 Stochastic Gradient Descent at Iteration k

Require: Learning rate ϵ_k

Require: Initial Parameter θ

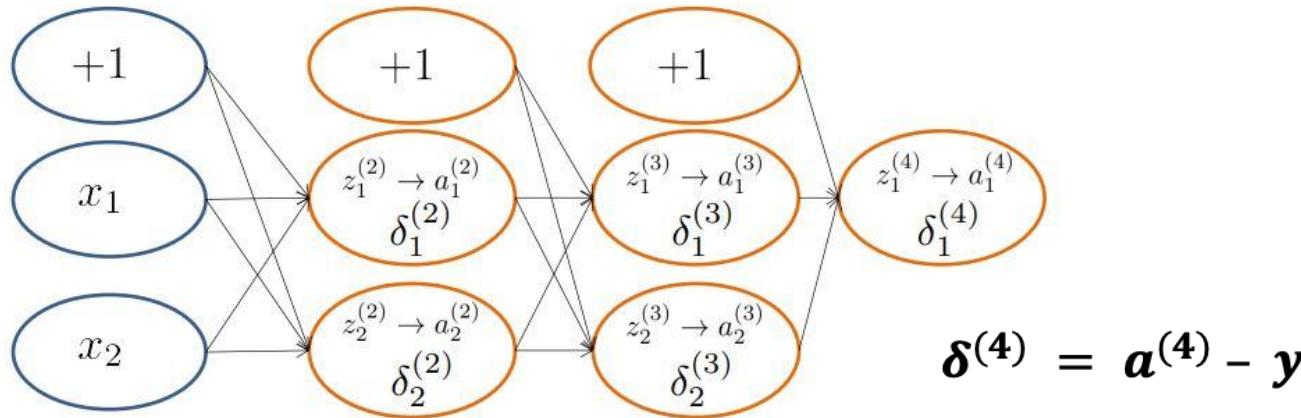
- 1: **while** stopping criteria not met **do**
- 2: Sample example $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ from training set
- 3: Compute gradient estimate:
- 4: $\hat{\mathbf{g}} \leftarrow +\nabla_{\theta} L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
- 5: Apply Update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$
- 6: **end while**

Batch and Stochastic Gradient Descent

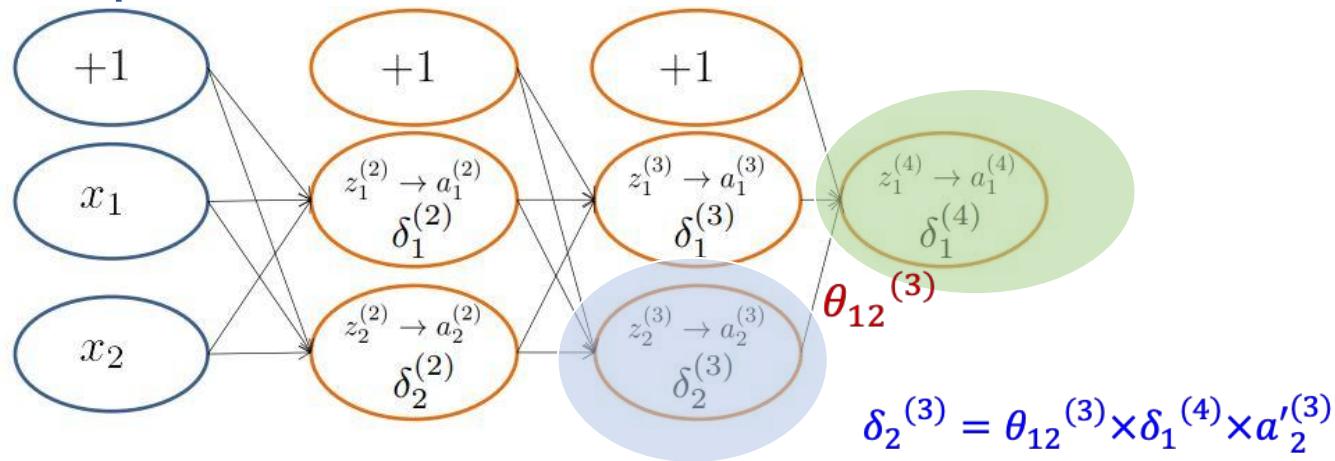


Extra Slides for BackProp Visualization

Backpropagation Visualization



Backpropagation

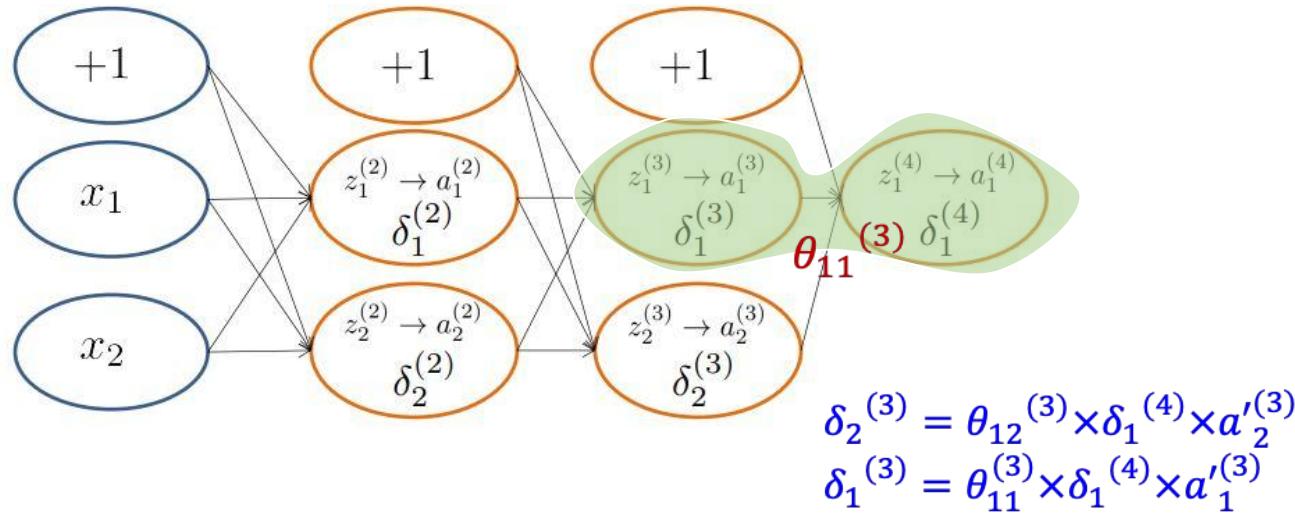


$\delta_j^{(l)}$ = “error” of node j in layer l

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

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

Backpropagation

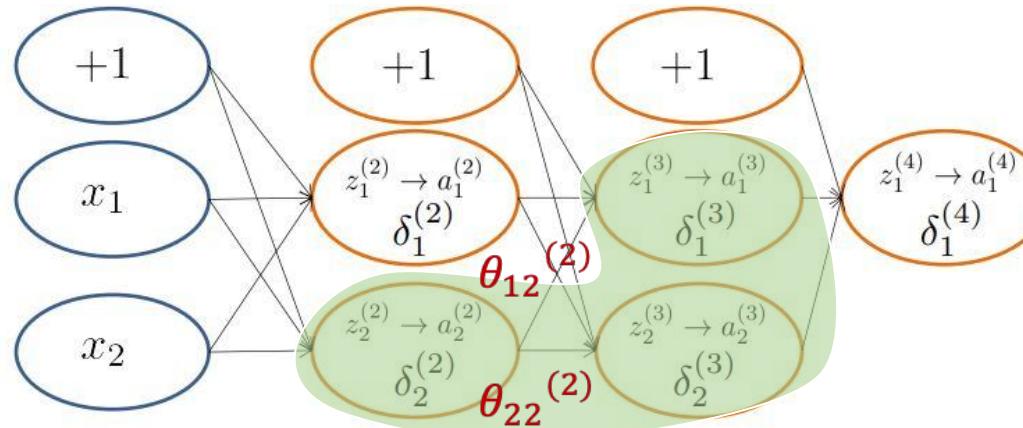


$\delta_j^{(l)}$ = “error” of node j in layer l

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

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

Backpropagation



$$\delta_2^{(2)} = \theta_{12}^{(2)} \times \delta_1^{(3)} \times a_2^{(2)} + \theta_{22}^{(2)} \times \delta_2^{(3)} \times a_2^{(2)}$$

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

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

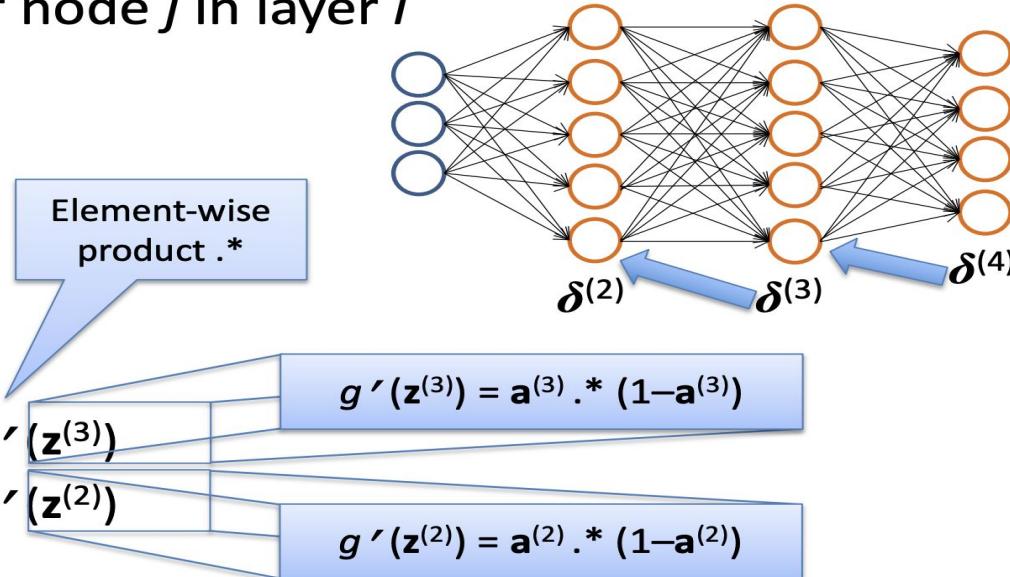
where $\text{cost}(\mathbf{x}_i) = y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))$

Gradient Computation

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

Backpropagation

- $\delta^{(4)} = a^{(4)} - y$
- $\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} .* g'(z^{(3)})$
- $\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} .* g'(z^{(2)})$
- (No $\delta^{(1)}$)



$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)} \quad (\text{ignoring } \lambda; \text{ if } \lambda = 0)$$