# DS 4400

# Machine Learning and Data Mining I
# Spring 2021

Alina Oprea

Associate Professor

Khoury College of Computer Science

Northeastern University

April 13 2021

# Announcements

- Project presentations
  - Tuesday, April 20, 11:45am-2:30pm
  - Thursday, April 22, 9am-12pm
  - 8 minutes per team (5 minute presentation + 3 min questions)

- Project report
  - Monday, April 26, at midnight
  - No late days!

# Outline

- Training Neural Networks
  - Backpropagation
  - Parameter Initialization
  - Derivation for feed-forward neural network for binary classification (sigmoid activation)
- Stochastic Gradient Descent
  - Gradient descent variants

# How to train Neural Networks?

- Backpropagation algorithm
- David Rumelhart, Geoffrey Hinton, Ronald Williams. "Learning representations by back-propagating errors". Nature. 323 (6088): 533–536. 1986
- Applicable to both FFNN and CNN
- Extension of Gradient Descent to multi-layer neural networks

# Reminder: Logistic Regression

$$J(\theta) = -\sum_{i=1}^{N} [y_i \log h_\theta(x_i) + (1 - y_i)\log(1 - h_\theta(x_i))]$$

- Cost of a single instance:

$$\text{cost}(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$

- Can re-write objective function as

$$J(\boldsymbol{\theta}) = \sum_{i=1}^{n} \text{cost}\left( h_\theta(x_i), y_i \right)$$

CROSS-ENTROPY LOSS
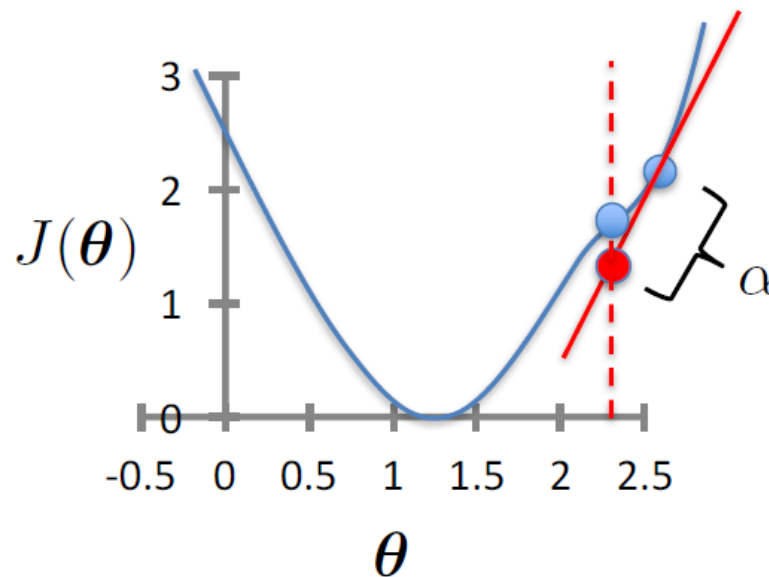
# Gradient Descent

- Initialize $\theta$  RANDOM          $\boldsymbol{\theta} = (W, b)$

- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

simultaneous update
for j = 0 … d

learning rate (small)
e.g., α = 0.05



$J(\boldsymbol{\theta})$

$\alpha$

$\theta$

- Converges for convex objective
- Could get stuck in local minimum for non-convex objectives

# Training Neural Networks

- Training data $x_1, y_1, \dots x_N, y_N$    *N training examples*

- One training example $x_i = (x_{i1}, \dots x_{id}), \text{label } y_i$

- One forward pass through the network
  - Compute prediction $\boxed{\hat{y}_i = h_\theta(x_i)}$    $\theta = [W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}]$

- Loss function

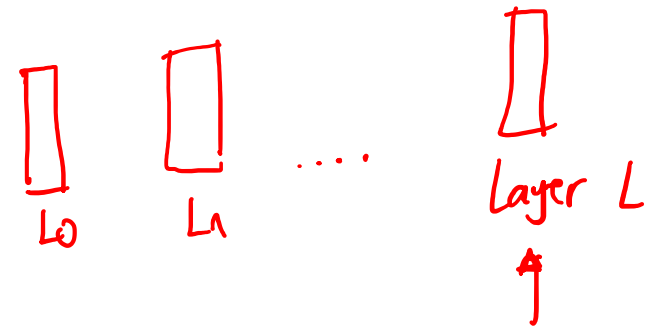$$L(y, \hat{y}) = -\left[ y \log \hat{y} + (1-y) \log (1-\hat{y}) \right] \quad \text{CROSS-ENTROPY LOSS}$$

$$J(\theta) = \sum_{i=1}^{N} L(y_i, \hat{y}_i)$$

# GD for Neural Networks

- ## Initialization
  - For all layers $\ell = 1, L$
    - Initialize $W^{[\ell]}, b^{[\ell]}$ RANDOM

- ## Backpropagation
  - Fix learning rate $\alpha$
  - Repeat
    - For all layers $\ell$ STARTING FROM LAST LAYER $L, \cdots, 1$

$$W^{[\ell]} \leftarrow W^{[\ell]} - \alpha \frac{\partial J(\theta)}{\partial W^{[\ell]}}$$

$$b^{[\ell]} \leftarrow b^{[\ell]} - \alpha \frac{\partial J(\theta)}{\partial b^{[\ell]}}$$

UNTIL STOPPING CONDITION

$L_0$    $L_1$   ....    Layer $L$

# GD for Neural Networks

- ## Initialization
  - For all layers $\ell$
    - Set $W^{[\ell]}, b^{[\ell]}$ at random

- ## Backpropagation
  - Fix learning rate $\alpha$
  - Repeat
    - For all layers $\ell$ (starting backwards)
    - $W^{[\ell]} = W^{[\ell]} - \alpha \sum_{i=1}^{N} \frac{\partial L(\hat{y}_i, y_i)}{\partial W^{[\ell]}}$
    - $b^{[\ell]} = b^{[\ell]} - \alpha \sum_{i=1}^{N} \frac{\partial L(\hat{y}_i, y_i)}{\partial b^{[\ell]}}$

This is expensive!

# Stochastic Gradient Descent

- ## Initialization
  - For all layers $\ell$
    - Set $W^{[\ell]}, b^{[\ell]}$ at random

- ## Backpropagation
  - Fix learning rate $\alpha$
  - Repeat
    - For all layers $\ell$ (starting backwards)
      - For all training examples $x_i, y_i$

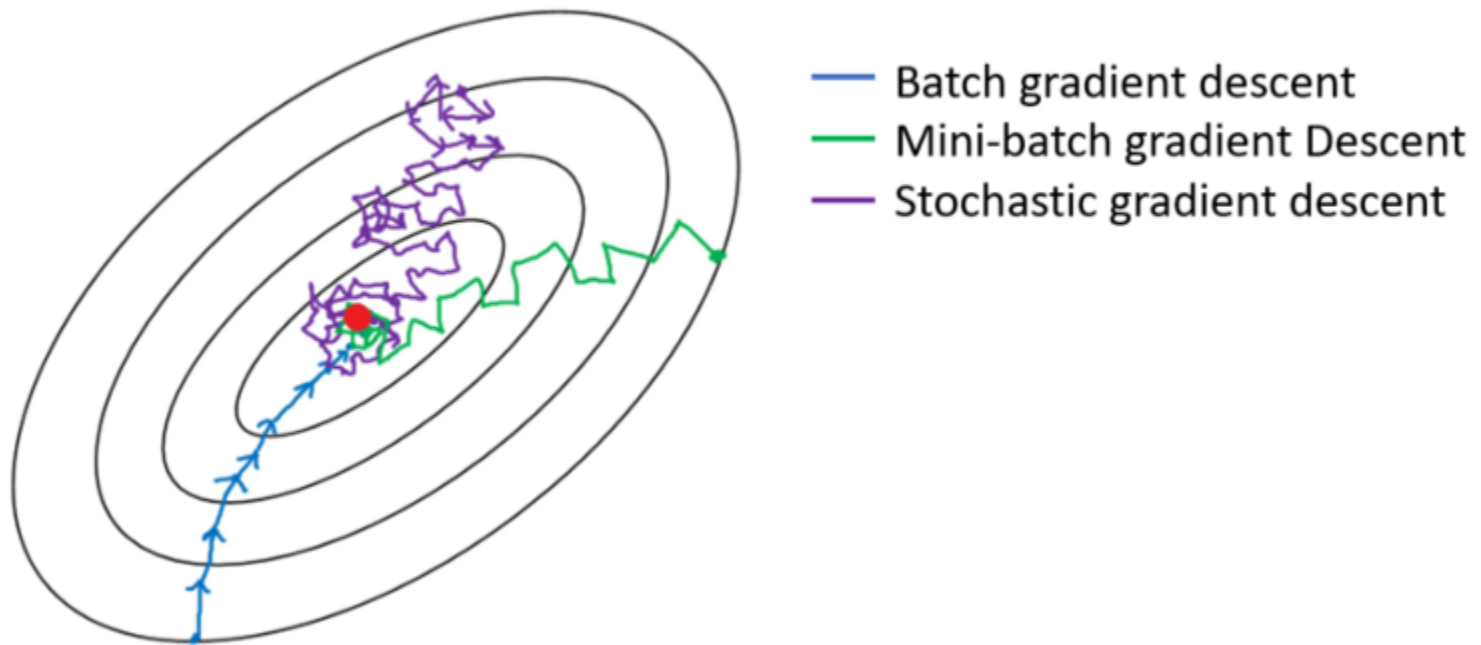$$W^{[\ell]} = W^{[\ell]} - \alpha \boxed{\frac{\partial L(\hat{y}_i, y_i)}{\partial W^{[\ell]}}}$$

$$b^{[\ell]} = b^{[\ell]} - \alpha \frac{\partial L(\hat{y}_i y_i)}{\partial b^{[\ell]}}$$

Incremental version of GD

# Online Perceptron

Let $\theta \leftarrow [0,0,\ldots,0]$
Repeat:
    Receive training example $(x_i, y_i)$
    If $y_i \theta^T x_i \leq 0$           // prediction is incorrect
        $\theta \leftarrow \theta + y_i\, x_i$
Until stopping condition

**Online learning** – the learning mode where the model update is performed each time a single observation is received

**Batch learning** – the learning mode where the model update is performed after observing the entire training set

# Mini-batch Gradient Descent

- Initialization
  - For all layers $\ell$
    - Set $W^{[\ell]}, b^{[\ell]}$ at random

- Backpropagation
  - Fix learning rate $\alpha$
  - Repeat

    *RANDOMIZED*

    - For all layers $\ell$ (starting backwards)
      - For all batches b of size B with training examples $x_{ib}, y_{ib}$

    $$W^{[\ell]} = W^{[\ell]} - \alpha \sum_{i=1}^{B} \frac{\partial L(\hat{y}_{ib}, y_{ib})}{\partial W^{[\ell]}}$$

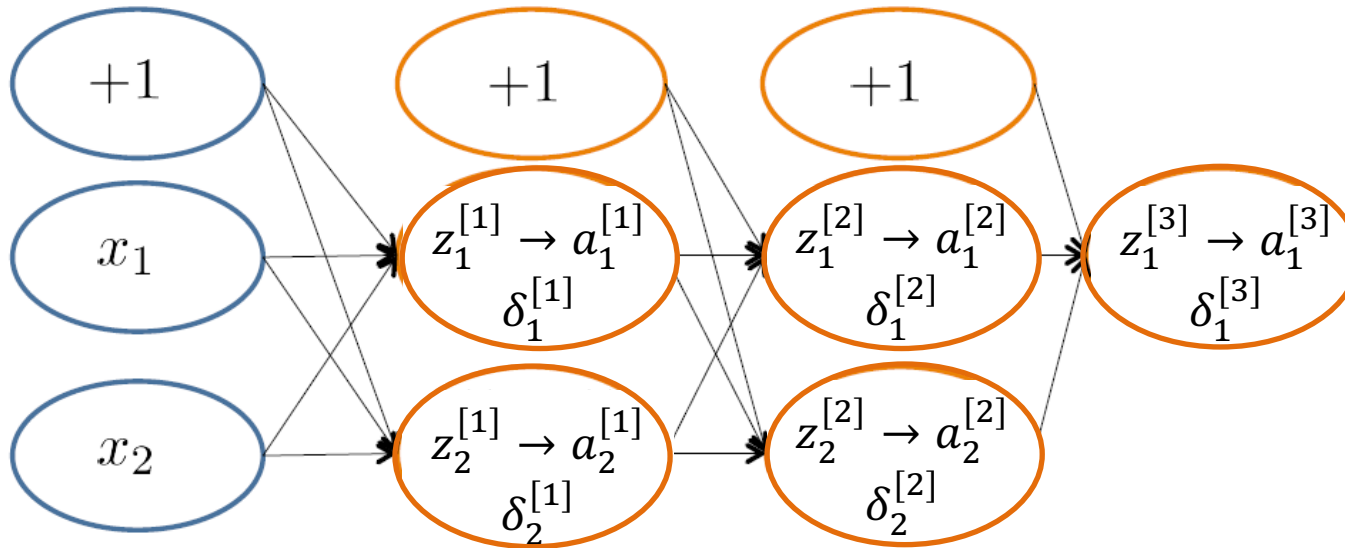    $$b^{[\ell]} = b^{[\ell]} - \alpha \sum_{i=1}^{B} \frac{\partial L(\hat{y}_{ib}, y_{ib})}{\partial b^{[\ell]}}$$

# Gradient Descent Variants

# Gradient Descent Variants



Batch gradient descent
Mini-batch gradient Descent
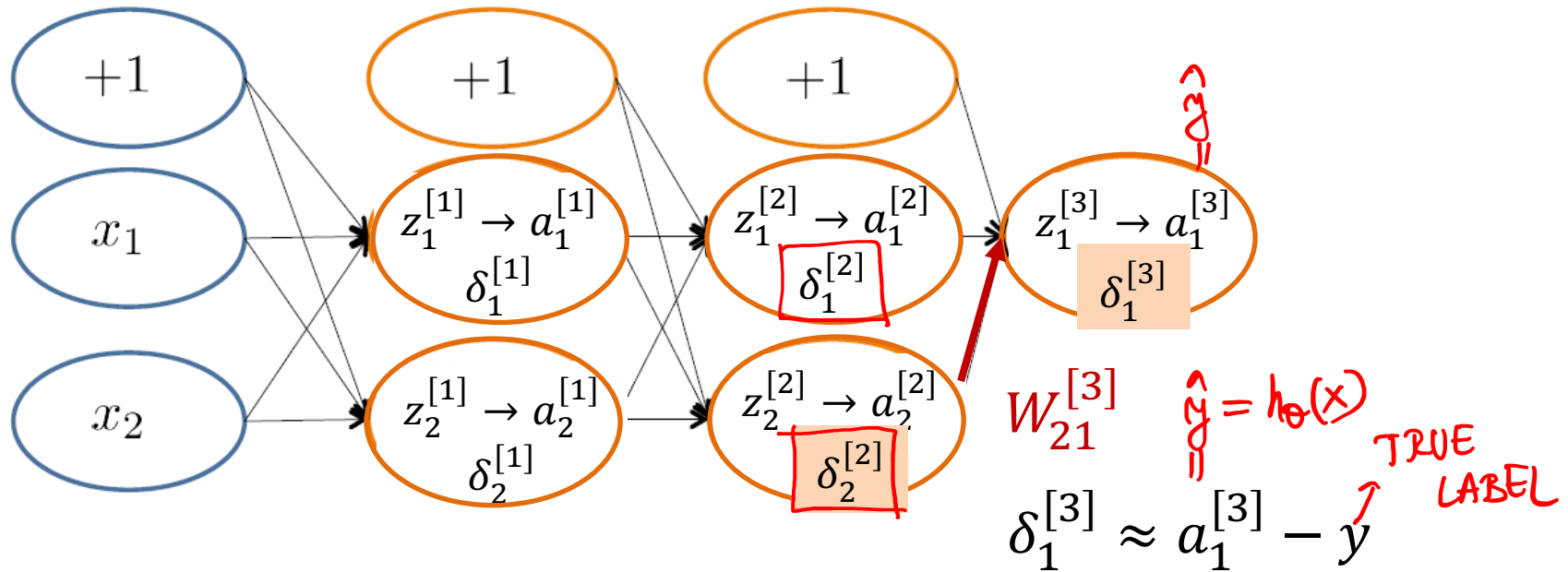Stochastic gradient descent

# Backpropagation Intuition



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

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

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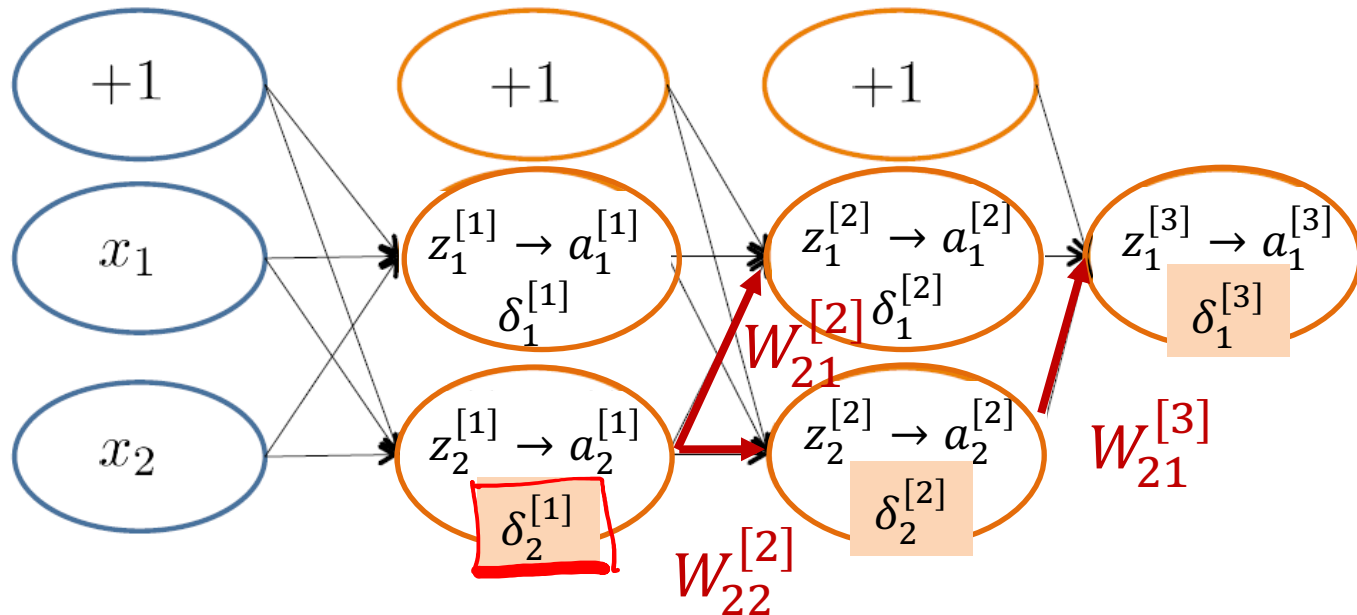
# Backpropagation Intuition



$$\delta_1^{[3]} \approx a_1^{[3]} - \hat{y}$$

$$\delta_2^{[2]} \approx \delta_1^{[3]} W_{21}^{[3]}$$

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

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

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

16

# Backpropagation Intuition



$$\delta_2^{[1]} \approx W_{21}^{[2]} \delta_1^{[2]} + W_{22}^{[2]} \delta_2^{[2]}$$

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

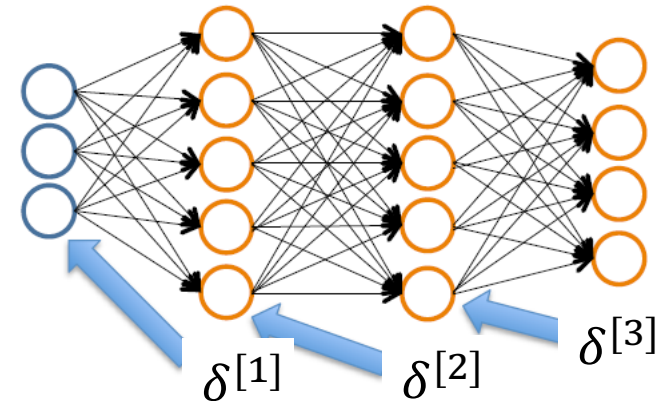Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

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

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

$$L(y, \hat{y}) = -[(1 - y)\log(1 - \hat{y}) + y\log\hat{y}]$$

## Definitions

$$- \quad z^{[\ell]} = W^{[\ell]} a^{[\ell-1]} + b^{[\ell]}, a^{[\ell]} = g(z^{[\ell]})$$

LINEAR · ACT

$$- \quad \delta^{[\ell]} = \frac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}}; \text{ Output } \hat{y} = a^{[L]} = g(z^{[L]})$$

DEF

$\delta^{[1]}$ $\delta^{[2]}$ $\delta^{[3]}$

- LAYER L: $\delta^{[L]} = \frac{\partial L(y,\hat{y})}{\partial z^{[L]}} = \frac{\partial L(y,\hat{y})}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{[L]}} = \frac{\partial L(y,\hat{y})}{\partial \hat{y}} \cdot g'(z^{[L]})$

- LAYER $\ell < L$: $\delta^{[\ell]} = \frac{\partial L(y,\hat{y})}{\partial z^{[\ell]}} = \frac{\partial L(y,\hat{y})}{\partial z^{[\ell+1]}} \cdot \frac{\partial z^{[\ell+1]}}{\partial a^{[\ell]}} \cdot \frac{\partial a^{[\ell]}}{\partial z^{[\ell]}}$

$z^{[\ell+1]} = W^{[\ell+1]} a^{[\ell]} + b^{[\ell+1]}$

$\delta^{[\ell+1]}$ $\quad W^{[\ell+1]}$ $\quad g'(z^{[\ell]})$

$= \delta^{[\ell+1]} W^{[\ell+1]} g'(z^{[\ell]})$

18

# Backpropagation

GD UPDATE

- $\frac{\partial L(y, \hat{y})}{\partial W^{[l]}} = \frac{\partial L(y, \hat{y})}{\partial z^{[l]}} \cdot \frac{\partial z^{[l]}}{\partial W^{[l]}} = \delta^{[l]} a^{[l-1]T}$
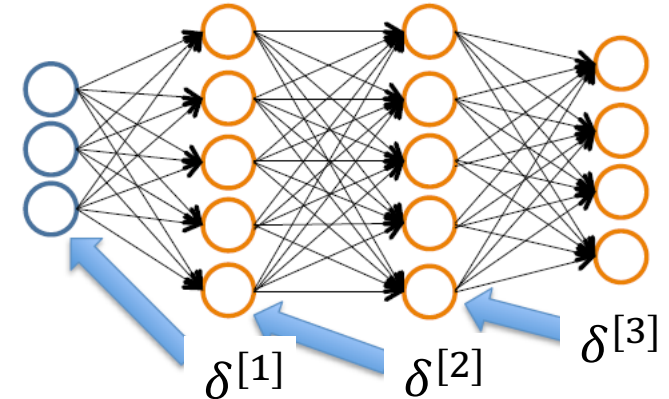
$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

with $\delta^{[l]}$ and $a^{[l-1]T}$ marked.

- $\frac{\partial L(y, \hat{y})}{\partial b^{[l]}} = \frac{\partial L(y, \hat{y})}{\partial z^{[l]}} \cdot \frac{\partial z^{[l]}}{\partial b^{[l]}} = \delta^{[l]}$

# Backpropagation

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

$$L(y, \hat{y}) = -[(1 - y)\log(1 - \hat{y}) + y\log\hat{y}]$$
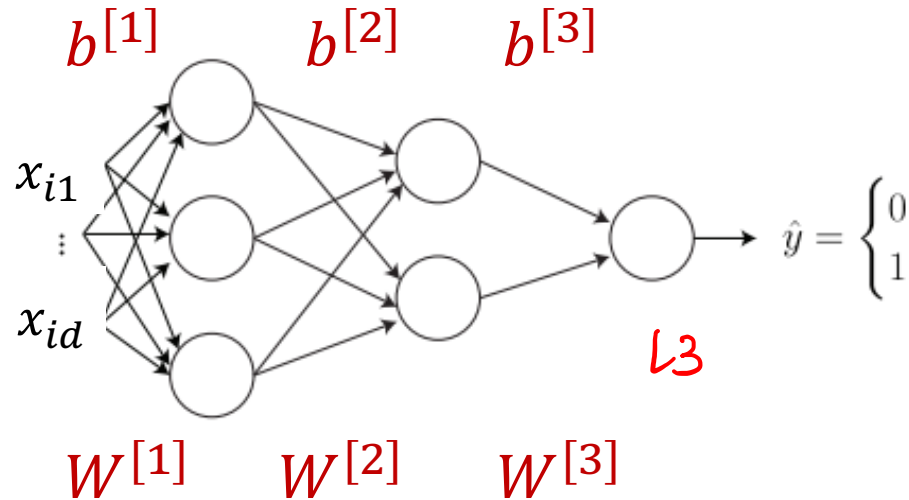
**Definitions**

- $z^{[\ell]} = W^{[\ell]} a^{[\ell-1]} + b^{[\ell]}, a^{[\ell]} = g(z^{[\ell]})$

- $\delta^{[\ell]} = \dfrac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}}$; Output $\hat{y} = a^{[L]} = g(z^{[L]})$

1. For last layer L: $\delta^{[L]} = \dfrac{\partial L(\hat{y}, y)}{\partial z^{[L]}} = \dfrac{\partial L(\hat{y}, y)}{\partial \hat{y}} \dfrac{\partial \hat{y}}{\partial z^{[L]}} = \dfrac{\partial L(\hat{y}, y)}{\partial \hat{y}} g'(z^{[L]})$

2. For layer $\ell$: $\delta^{[\ell]} = \dfrac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}} = \dfrac{\partial L(\hat{y}, y)}{\partial z^{[\ell+1]}} \dfrac{\partial z^{[\ell+1]}}{\partial a^{[\ell]}} \dfrac{\partial a^{[\ell]}}{\partial z^{[\ell]}} = \delta^{[\ell+1]} W^{[\ell+1]} g'(z^{[\ell]})$

3. Compute parameter gradients

- $\dfrac{\partial L(\hat{y}, y)}{\partial W^{[\ell]}} = \dfrac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}} \dfrac{\partial z^{[\ell]}}{\partial W^{[\ell]}} = \delta^{[\ell]} a^{[\ell-1]T}$

- $\dfrac{\partial L(\hat{y}, y)}{\partial b^{[\ell]}} = \dfrac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}} \dfrac{\partial z^{[\ell]}}{\partial b^{[\ell]}} = \delta^{[\ell]}$



$\delta^{[1]}$     $\delta^{[2]}$     $\delta^{[3]}$

# Example 2 Hidden Layers

Training data
Dimension d



$b^{[1]}$    $b^{[2]}$    $b^{[3]}$

$x_{i1}$

$\vdots$

$x_{id}$

$\hat{y} = \begin{cases} 0 \\ 1 \end{cases}$

L3

$W^{[1]}$    $W^{[2]}$    $W^{[3]}$

$$z^{[1]} = W^{[1]} x_i + b^{[1]}$$    LIN

$$a^{[1]} = g(z^{[1]})$$    ACT

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$    LIN

$$a^{[2]} = g(z^{[2]})$$    ACT

$$z^{[3]} = W^{[3]} a^{[2]} + b^{[3]}$$

$$\hat{y}^{(i)} = a^{[3]} = g(z^{[3]})$$    OUTPUT

$g = $ SIGMOID

$$\hat{y} = a^{[3]} = g(z^{[3]})$$

21

# Binary Classification Example

$$\delta^{[3]} = \frac{\partial l(y, \hat{y})}{\partial \hat{y}} \cdot g'(z^{[3]}) = \frac{\hat{y} - y}{(1-\hat{y})\hat{y}} \cdot \underbrace{g(z^{[3]})}_{a = \hat{y}} \underbrace{(1 - g(z^{[3]}))}_{1-\hat{y}} = \hat{y} - y$$

$$L(y, \hat{y}) = -\left[ (1-y)\log(1-\hat{y}) + y\log\hat{y} \right]$$

$$\frac{\partial l(y, \hat{y})}{\partial \hat{y}} = \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}} = \frac{\hat{y} - y\hat{y} - y + y\hat{y}}{(1-\hat{y})\hat{y}} = \frac{\hat{y} - y}{(1-\hat{y})\hat{y}}$$

$$g(z) = \frac{1}{1+e^{-z}} \quad, \quad g'(z) = g(z)(1-g(z))$$
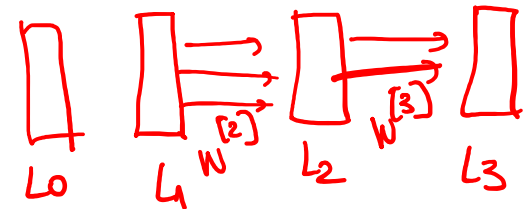
$$\delta^{[3]} = \hat{y} - y = a^{[3]} - y \qquad \longrightarrow \quad \text{ERROR} \quad \text{(ONLY FOR SIGMOID ACTIVATION)}$$

PREDICTION $= h_\theta(x)$

TRUE LABEL

# Binary Classification Example

$$\frac{\partial L(y, \hat{y})}{\partial W^{[3]}} = \delta^{[3]} a^{[2]T} = \left(a^{[3]} - y\right) a^{[2]T}$$

$$\frac{\partial L(y, \hat{y})}{\partial b^{[3]}} = \delta^{[3]} = a^{[3]} - y$$



$L_0 \quad L_1 \quad W^{[2]} \quad L_2 \quad W^{[3]} \quad L_3$

# Binary Classification Example

LAYER 2: $\delta^{[2]} = \delta^{[3]} \cdot W^{[3]} g'(z^{[2]})$

$$= (a^{[3]} - y) W^{[3]} \cdot g(z^{[2]})(1 - g(z^{[2]})) = (a^{[3]} - y) W^{[3]} a^{[2]}(1 - a^{[2]})$$

$$\overset{\shortparallel}{a^{[2]}}$$

$$\frac{\partial L(y, \hat{y})}{\partial W^{[2]}} = \delta^{[2]} a^{[1]T} = (a^{[3]} - y) W^{[3]} a^{[2]}(1 - a^{[2]}) a^{[1]T}$$

$$\frac{\partial L(y, \hat{y})}{\partial b^{[2]}} = \delta^{[2]} = (a^{[3]} - y) W^{[3]} a^{[2]}(1 - a^{[2]})$$

EXERCISE: COMPUTE $\delta^{[1]}$

$$\frac{\partial L(y, \hat{y})}{\partial W^{[1]}}, \quad \frac{\partial L(y, \hat{y})}{\partial b^{[1]}}$$

# Parameter Initialization

- How about we set all W and b to 0?

$$z^{[1]} = W^{[1]}x + b^{[1]} = [0, 0, \ldots, 0]$$

$$a^{[1]} = g(z^{[1]}) = \left[\frac{1}{2}, \ldots, \frac{1}{2}\right]$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} = [0, \ldots, 0]$$

$$a^{[2]} = \left[\frac{1}{2}, \ldots, \frac{1}{2}\right]$$

$$a^{[L]} = \left[\frac{1}{2}, \ldots, \frac{1}{2}\right]$$

$$(x, y) \rightarrow \left(\frac{1}{2}, \ldots, \frac{1}{2}\right)$$

INIT PARAMS AT RANDOM

# Training NN with Backpropagation

Given training set $(x_1, y_1), \ldots, (x_N, y_N)$
Initialize all parameters $W^{[\ell]}, b^{[\ell]}$ randomly, for all layers $\ell$
Loop

Set $\Delta_{ij}^{[l]} = 0$, for all layers $l$ and indices $i, j$     **EPOCH**    BATCH

For each training instance $(x_k, y_k)$:

    Compute $a^{[1]}, a^{[2]}, \ldots, a^{[L]}$ via forward propagation

    Compute errors $\delta^{[L]} = a^{[L]} - y_k, \delta^{[L-1]}, \ldots \delta^{[1]}$   BACKWARD

    Compute gradients $\Delta_{ij}^{[l]} = \Delta_{ij}^{[l]} + a_j^{[l-1]} \delta_i^{[l]}$   $\delta^{[l]} a^{[l-1]T}$

  Update weights via gradient step

- $W_{ij}^{[\ell]} = W_{ij}^{[\ell]} - \alpha \Delta_{ij}^{[\ell]}$   $\delta^{[l]} a^{[l-1]T}$

- Similar for $b_{ij}^{[\ell]}$   $b^{[l]} \leftarrow b^{[l]} - \alpha \delta^{[l]}$

Until weights converge or maximum number of epochs is reached

27

# Training Neural Networks

- Randomly initialize weights

- Implement forward propagation to get prediction $\hat{y}_i$ for any training instance $x_i$

- Compute loss function $L(\hat{y}_i, y_i)$

- Implement backpropagation to compute partial derivatives $\frac{\partial L(\hat{y}_i, y_i)}{\partial W^{[\ell]}}$ and $\frac{\partial L(\hat{y}_i, y_i)}{\partial b^{[\ell]}}$

- Use gradient descent with backpropagation to compute parameter values that optimize loss

- Can be applied to both feed-forward and convolutional nets

# Materials

- Stanford tutorial on training Multi-Layer Neural Networks
  - http://ufldl.stanford.edu/tutorial/supervised/MultiLayerNeuralNetworks/
- Notes on backpropagation by Andrew Ng
  - http://cs229.stanford.edu/notes-spring2019/backprop.pdf
- Deep learning notes by Andrew Ng
  - http://cs229.stanford.edu/notes2020spring/cs229-notes-deep_learning.pdf

# Representing Boolean Functions

**Simple example: AND**
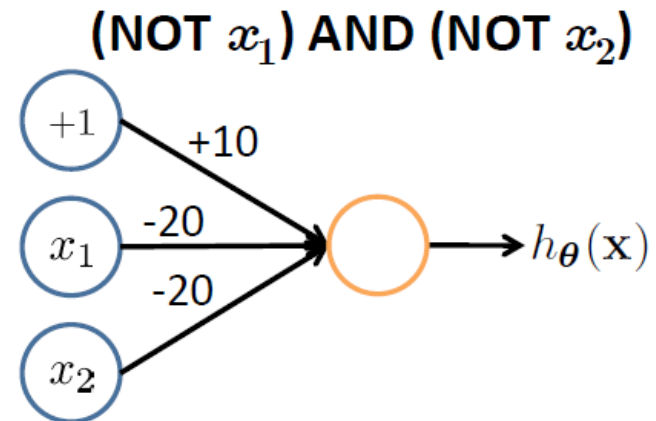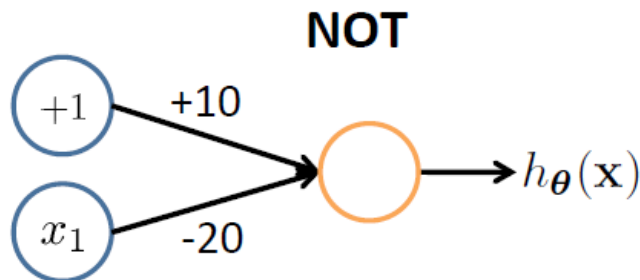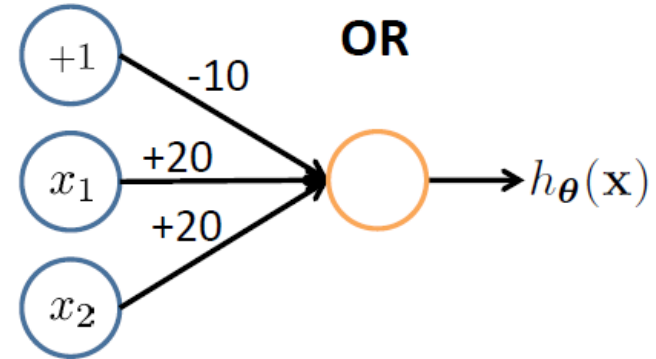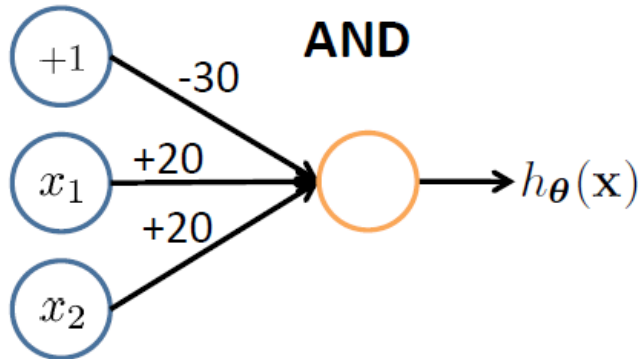
$x_1, x_2 \in \{0, 1\}$
$y = x_1 \text{ AND } x_2$

Activation function

$$g(z) = \begin{cases} 1, x \geq 0 \\ 0, x < 0 \end{cases}$$



$+1$ — ^

$x_1$ 0.5

$x_2$ 0.5

$\longrightarrow h_\theta(\mathbf{x})$

$h_\Theta(\mathbf{x}) = g(\ ? \ + \ ? \ x_1 + \ ? \ x_2)$

| $x_1$ | $x_2$ | $h_\Theta(\mathbf{x})$ |
|-------|-------|------------------------|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

30

# Representing Boolean Functions



**AND**

$+1$   $-30$
$x_1$   $+20$
$x_2$   $+20$
$\rightarrow h_{\boldsymbol{\theta}}(\mathbf{x})$

**OR**

$+1$   $-10$
$x_1$   $+20$
$x_2$   $+20$
$\rightarrow h_{\boldsymbol{\theta}}(\mathbf{x})$

**NOT**

$+1$   $+10$
$x_1$   $-20$
$\rightarrow h_{\boldsymbol{\theta}}(\mathbf{x})$

**(NOT $x_1$) AND (NOT $x_2$)**

$+1$   $+10$
$x_1$   $-20$
$x_2$   $-20$
$\rightarrow h_{\boldsymbol{\theta}}(\mathbf{x})$

# XOR

- Need at least one hidden layer to compute XOR!



NOT (XOR)

II        I
−        +        0

Non-linearly separable

+        −
III      IV

NOT (X1 XOR X2) = (X1 AND X2) OR (NOT X1 AND NOT X2))

NOT(XOR)

| | | NOT(XOR) |
|---|---|---|
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

# Combining Representations

# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
  - Andrew Moore
  - Yann LeCun
- Thanks!