

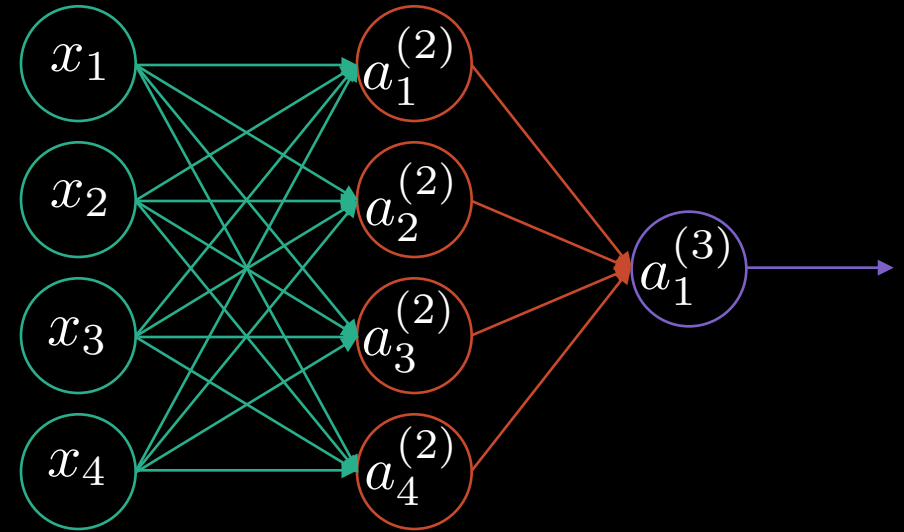
# Practical 2.1

Neural Networks – back propagation (training)

# Overview

- Review on forward pass
- Training data and labels
- Loss function
- Local error
  - Output error
  - Layer  $\ell$  error (Jacobian of composite functions)
- Parameters' gradient
- Back propagation
- Stochastic and (mini-) batch gradient descent

# Neural network



$$\mathbf{a}^{(1)} = \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{s_1}$$

$$\mathbf{a}^{(\ell+1)} = \sigma(\mathbf{z}^{(\ell+1)}) = \sigma(\Theta^{(\ell)} \hat{\mathbf{a}}^{(\ell)}), \ell = 1, 2, \dots, L-1 \quad \hat{\mathbf{a}} = \begin{bmatrix} +1 \\ \underline{\mathbf{a}} \end{bmatrix} \in \mathbb{R}^{s_{\ell+1}}$$

$$h_{\Theta}(\mathbf{x}) = \mathbf{a}^{(L)} \in \mathbb{R}^{s_L} = \mathbb{R}^k$$

$$\Theta^{(j)} \in \mathbb{R}^{s_{j+1} \times (s_j + 1)} \quad \Theta = \{\Theta^{(1)}, \Theta^{(2)}, \dots, \Theta^{(L-1)}\}$$

# Training data (I)

$$\mathbf{X} = \begin{bmatrix} \underline{x}^{(1)} \\ \underline{x}^{(2)} \\ \underline{x}^{(3)} \\ \vdots \\ \underline{x}^{(m)} \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} \underline{x}^{(1)} \\ \underline{x}^{(2)} \\ \underline{x}^{(3)} \\ \vdots \\ \underline{x}^{(m)} \end{bmatrix}} \right\} m \quad \mathbf{Y} = \begin{bmatrix} \underline{y}^{(1)} \\ \underline{y}^{(2)} \\ \underline{y}^{(3)} \\ \vdots \\ \underline{y}^{(m)} \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} \underline{y}^{(1)} \\ \underline{y}^{(2)} \\ \underline{y}^{(3)} \\ \vdots \\ \underline{y}^{(m)} \end{bmatrix}} \right\} m$$

$\underbrace{\hspace{10em}}_n \qquad \underbrace{\hspace{10em}}_K \quad h_{\theta}(z) \in \mathbb{R}^K$

# Training data (II)

$$\begin{array}{c}
 \mathbf{X} \\
 \begin{array}{cccc}
 \underline{x_1} & \underline{x_2} & \dots & \underline{x_n} \\
 \hline
 x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\
 x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\
 \vdots & \vdots & & \vdots \\
 x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)}
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 \mathbf{Y} \\
 \begin{array}{cccc}
 \underline{y_1} & \underline{y_2} & \dots & \underline{y_k} \\
 \hline
 y_1^{(1)} & y_2^{(1)} & \dots & y_k^{(1)} \\
 y_1^{(2)} & y_2^{(2)} & \dots & y_k^{(2)} \\
 \vdots & \vdots & & \vdots \\
 y_1^{(m)} & y_2^{(m)} & \dots & y_k^{(m)}
 \end{array}
 \end{array}$$

if  $K=1 \Rightarrow [\mathbf{Y}] \equiv \underline{y_1}$

# Loss function

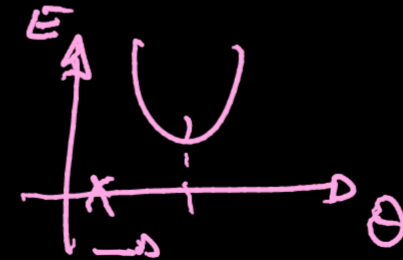
$$\mathcal{L}(\Theta) = \frac{1}{m} \sum_{i=1}^m E^{(i)}$$

$$E = E(h_{\Theta}(x)) \quad \begin{array}{l} x, y \text{ are fixed} \\ [\Theta] \text{ can vary} \end{array}$$

$$E = E(h_{\Theta}(x)) = \frac{1}{2} \| \underline{y}^{(i)} - \underbrace{h_{\Theta}(x^{(i)})}_{\underline{a}^{(L)}} \|^2 = \frac{1}{2} \sum_k (y_k - a_k^{(L)})^2$$

$$\frac{\partial E}{\partial \Theta_{ij}^{(\ell)}}$$

$$\underline{\Theta} \rightarrow \underline{\Theta} - \eta \frac{\partial E}{\partial \underline{\Theta}}$$



# Local error

$$\underline{z_i^{(\ell)}} \rightarrow x_i^{(e)} + \Delta x_i^{(e)}$$

$$\underline{E} \rightarrow E + \frac{\partial E}{\partial x_i^{(e)}} \cdot \Delta x_i^{(e)}$$

$$\left| \frac{\partial E}{\partial z_i^{(\ell)}} \right| \gg 0 \Rightarrow \Delta x_i^{(e)} \text{ can influence } E$$

$$\simeq 0 \Rightarrow \text{near optimal}$$

$$\delta_i^{(\ell)} := \frac{\partial E}{\partial x_i^{(e)}}$$

error @  $i$ -th neuron  
@ layer  $\ell$

# Output error

$$E = \frac{1}{2} \sum_k (y_k - a_k^{(L)})^2$$

$$\Rightarrow \frac{1}{2} \cdot 2 (a_i^{(L)} - y_i)$$

$$\delta_i^{(L)} = \frac{\partial E}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} = \frac{\partial E}{\underbrace{\partial a_i^{(L)}}_{a_i^{(L)} - y_i}}$$

$$\underbrace{\sigma'(z_i^{(L)})}_{a_i^{(L)}(1-a_i^{(L)})}$$

$$\sigma(z) = (1 + \exp(-z))^{-1}$$

$$\begin{aligned} \sigma'(z) &= + (1 + \exp(-z))^{-2} \cdot \exp(-z)(+1) \\ &= (1 + \exp(-z))^{-2} \cdot [(1 + \exp(-z)) - 1] \\ &= (1 + \exp(-z))^{-1} - (1 + \exp(-z))^{-2} = \\ &= \sigma(z) - \sigma^2(z) = \sigma(z)(1 - \sigma(z)) \\ &= a(1-a) \end{aligned}$$

$$\delta^{(L)} = \underbrace{\underline{a}^{(L)} - \underline{y}^{(L)}}_{\nabla_a E} \odot \underbrace{\underline{a}^{(L)} \odot (1 - \underline{a}^{(L)})}_{\sigma'(z^{(L)})}$$



# Layer $\ell$ error (I)

$$\delta^{(\ell)} = \delta^{(e)} \left( \delta^{(e+1)} \right)$$

# Jacobian of composite function

$$g: \mathbb{R} \rightarrow \mathbb{R}^n \text{ diff. } \mu_0, \quad f: \mathbb{R}^n \rightarrow \mathbb{R} \text{ diff. } \underline{x}^0 = g(\mu_0)$$

$$h = f \circ g: \mathbb{R} \rightarrow \mathbb{R} \text{ der. in } \mu_0 \Rightarrow$$

$$\begin{aligned} h'(\mu_0) &= (\nabla f)(\underline{x}^0) (\nabla g)(\mu_0) = \\ &= \left[ \frac{\partial f}{\partial x_1}(\underline{x}^0), \dots, \frac{\partial f}{\partial x_n}(\underline{x}^0) \right] \begin{bmatrix} g'_1(\mu_0) \\ \vdots \\ g'_n(\mu_0) \end{bmatrix} = \\ &= \langle \nabla f(\underline{x}^0), g'(\mu_0) \rangle \end{aligned}$$

# Layer $\ell$ error (II)

$$\underline{\delta^{(\ell)}} = \delta^{(\ell)} (\underline{\delta^{(\ell+1)}})$$

$$\delta_i^{(\ell)} = \frac{\partial E}{\partial z_i^{(\ell)}} \quad \delta_j^{(\ell+1)} = \frac{\partial E}{\partial z_j^{(\ell+1)}}$$

$\delta_i^{(\ell)} = \sum_j \overbrace{\frac{\partial E}{\partial z_j^{(\ell+1)}}}^{\delta_j^{(\ell+1)}} \underbrace{\frac{\partial z_j^{(\ell+1)}}{\partial z_i^{(\ell)}}}_{\text{pink box}}$

## Layer $\ell$ error (III)

$$\frac{\partial z_j^{(\ell+1)}}{\partial z_i^{(\ell)}} = \Theta_{ji}^{(e)} \sigma'(z_i^{(e)})$$

$$z_j^{(\ell+1)} = \sum_k \Theta_{jk}^{(e)} z_k^{(e)} = \sum_k \Theta_{jk}^{(e)} \sigma(z_k^{(e)})$$

$$\delta_i^{(\ell)} = \sum_j \delta_j^{(\ell+1)} \frac{\partial z_j^{(\ell+1)}}{\partial z_i^{(\ell)}} = \sum_j \Theta_{ji}^{(e)} \delta_j^{(\ell+1)} \sigma'(z_i^{(e)})$$

## Layer $\ell$ error (IV)

$$\delta_i^{(\ell)} = \sum_j \Theta_{ji}^{(\ell)} \delta_j^{(\ell+1)} \sigma'(z_i^{(\ell)})$$

$$\boldsymbol{\delta}^{(\ell)} = [(\boldsymbol{\Theta}^{(\ell)})^\top \boldsymbol{\delta}^{(\ell+1)}] \odot \boldsymbol{\sigma}'(\boldsymbol{z}^{(\ell)})$$

# Parameters' gradient

$$\frac{\partial E}{\partial \Theta_{ij}^{(l)}} = \underbrace{\frac{\partial E}{\partial z_i^{(e+1)}}}_{\substack{\text{from back} \\ \delta_i^{(e+1)}}} \underbrace{\frac{\partial z_i^{(e+1)}}{\partial \Theta_{ij}^{(e)}}}_{\substack{a_j^{(e)} \\ \text{forward pass}}}, \forall i, j$$

$$\frac{\partial E}{\partial \Theta^{(l)}} = \underline{a}^{(e)} (\underline{\delta}^{(e+1)})^T \Rightarrow \boxed{\quad}$$

# Back propagation

i)  $\mathbf{a}^{(1)} = \mathbf{x}, \hat{\mathbf{a}}^{(1)} = \begin{bmatrix} +1 \\ \mathbf{a}^{(1)} \end{bmatrix}$

ii)  $\mathbf{z}^{(\ell+1)} = \mathbf{\Theta}^{(\ell)} \hat{\mathbf{a}}^{(\ell)}, \mathbf{a}^{(\ell)} = \sigma(\mathbf{z}^{(\ell)}), \hat{\mathbf{a}}^{(\ell)} = \begin{bmatrix} +1 \\ \mathbf{a}^{(\ell)} \end{bmatrix}$

iii)  $\boldsymbol{\delta}^{(L)} = \nabla_h E \odot \sigma'(\mathbf{z}^{(L)})$

iv)  $\boldsymbol{\delta}^{(\ell)} = [(\mathbf{\Theta}^{(\ell)})^\top \boldsymbol{\delta}^{(\ell+1)}] \odot \sigma'(\mathbf{z}^{(\ell)})$

v)  $\frac{\partial E}{\partial \mathbf{\Theta}^{(\ell)}} = \mathbf{a}^{(\ell)} (\boldsymbol{\delta}^{(\ell+1)})^\top$

# Weight update

## GRADIENT DESCENT

- Stochastic SGD

$$\Theta^{(\ell)} \rightarrow \Theta^{(e)} - \eta \underline{a}^{(e)} \underline{\delta}^{(e+1)}$$

- Batch (or mini-batch)

$$\Theta^{(\ell)} \rightarrow \Theta^{(e)} - \eta \frac{1}{m} \sum_{i=1}^m (\underline{a}^{(e)} \underline{\delta}^{(e+1)})^{(i)}$$