

# Introduction to Neural Networks

DATA 607 — Session 5 — 11/03/2019

depth of network (number of layers):  $d$   
 number of neurons in layer  $\ell$ :  $p_\ell$   
 activation (output) of neuron  $k$  in layer  $\ell$ :  $a_k^{[\ell]}$   
 bias of neuron  $k$  in layer  $\ell$ :  $b_k^{[\ell]}$   
 weight of the connection between neuron  $j$   
 in layer  $\ell$  and neuron  $i$  in layer  $\ell + 1$ :  $w_{ij}^{[\ell]}$   
 activation function:  $h$

$$a_i^{[\ell+1]} = h\left(z_i^{[\ell+1]}\right), \quad \text{where} \quad z_i^{[\ell+1]} = b_i^{[\ell]} + \sum_{j=1}^{p_\ell} w_{ij}^{[\ell]} a_j^{[\ell]}$$

Vectorize:

$$\mathbf{z}^{[\ell]} = \begin{bmatrix} z_1^{[\ell]} \\ \vdots \\ z_{p_\ell}^{[\ell]} \end{bmatrix} \in \mathbb{R}^{p_\ell}, \quad \mathbf{a}^{[\ell]} = \begin{bmatrix} a_1^{[\ell]} \\ \vdots \\ a_{p_\ell}^{[\ell]} \end{bmatrix} \in \mathbb{R}^{p_\ell},$$

$$\mathbf{b}^{[\ell]} = \begin{bmatrix} b_1^{[\ell]} \\ \vdots \\ b_{p_\ell}^{[\ell]} \end{bmatrix} \in \mathbb{R}^{p_\ell}, \quad \mathbf{W}^{[\ell]} = \begin{bmatrix} w_{11}^{[\ell]} & \cdots & w_{1p_\ell}^{[\ell]} \\ \vdots & \ddots & \vdots \\ w_{p_{\ell+1}1}^{[\ell]} & \cdots & w_{1p_\ell}^{[\ell]} \end{bmatrix} \in \mathbb{R}^{p_{\ell+1} \times p_\ell}$$

Apply  $h$  componentwise:

$$\mathbf{a}_i^{[\ell+1]} = h\left(\mathbf{z}^{[\ell+1]}\right) = h\left(\mathbf{b}^{[\ell]} + \mathbf{W}^{[\ell]} \mathbf{a}^{[\ell]}\right)$$

Useful intermediate quantity:

$$\mathbf{b}^{[\ell]} + \mathbf{W}^{[\ell]} \mathbf{a}^{[\ell]}$$

The process of computing  $\hat{\mathbf{y}}$  from  $\mathbf{x}$ , given  $\mathbf{b}^{[\ell]}$  and  $\mathbf{W}^{[\ell]}$ ,  $\ell = 1, \dots, d$ , is called **forward propagation** of data.

A **loss function**,  $L(\hat{\mathbf{y}}, \mathbf{y})$ , assesses a penalty based on the error in approximating  $\mathbf{y}$  by  $\hat{\mathbf{y}}$ .

The total loss associated to a training set

$$T = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$$

is called the **cost** of  $T$ :

$$C(T) = \sum_{i=1}^n L(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

We adjust the variables  $\mathbf{b}^{[\ell]}$  and  $\mathbf{W}^{[\ell]}$  based on the derivatives

$$\frac{\partial C}{\partial b_i^{[\ell]}} \quad \text{and} \quad \frac{\partial C}{\partial w_{ij}^{[\ell]}}$$

- Softmax regression
- Gradient descent