

CH3.4 Neural Networks

→ Add more layers between the input and output layers. This makes your model DEEPER leading to the term Deep Learning.

→ If we just pass linear functions to linear functions the final result is again a linear function. To overcome this, output of each layer is evaluated at a specific nonlinear function called "activation function".

→ We can add as many layers we want. Every layer, just like input and output layers, is composed of nodes called "neuron". Every neuron is connected to the output of all neurons in the previous layer and the input of all neurons in the next layer. This is inspired by the actual brain neurons.

→ So we have a network of neurons. That's why we call this structure as Neural Network.

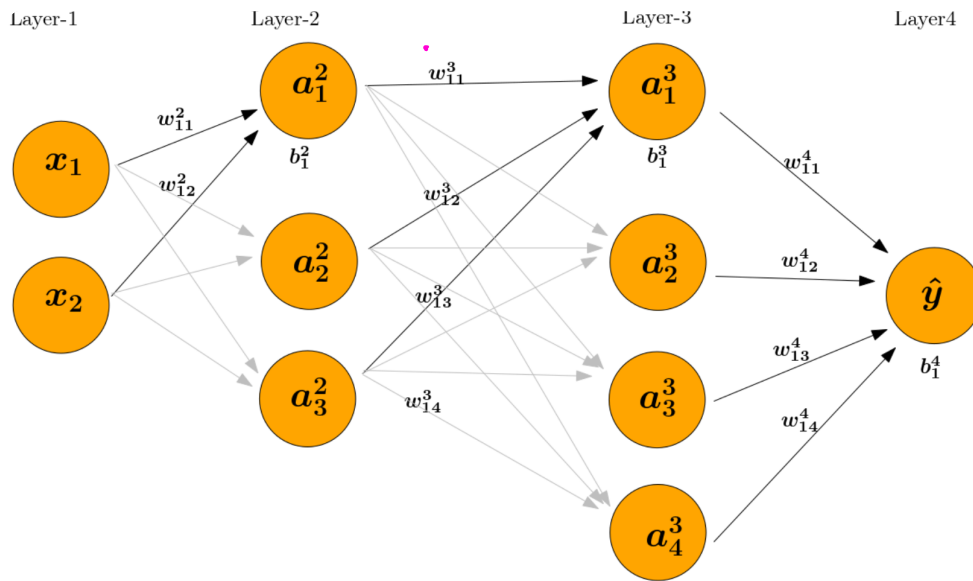
→ "I know a guy know who knows a guy"

→ Almost all neural network models have two components.

(1) **Forward Pass**: get the input and pass it through the layers all the way to the last layer and compute the cost (loss)

(2) **Backpropagation**: compute the derivative of the cost function with respect to ALL parameters and update them using gradient descent method

① Forward Pass



w_{jk}^l → weight from k^{th} neuron in layer- $(l-1)$ to the j^{th} neuron in layer- l
 b_j^l → bias of the j^{th} neuron of layer- l
 a_j^l → activation of j^{th} neuron in the l^{th} layer

Ex: w_{31}^2 : weight from 3rd neuron in layer-1 to the 1st neuron in layer-2
 b_1^2 : bias of the 1st neuron in layer-2

From Layer-1 to Layer-2

$$\begin{aligned} z_1^2 &= w_{11}^2 x_1 + w_{12}^2 x_2 + b_1^2 \\ z_2^2 &= w_{21}^2 x_1 + w_{22}^2 x_2 + b_2^2 \\ z_3^2 &= w_{31}^2 x_1 + w_{32}^2 x_2 + b_3^2 \end{aligned} \Rightarrow \begin{matrix} \begin{bmatrix} z_1^2 \\ z_2^2 \\ z_3^2 \end{bmatrix} = \begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1^2 \\ b_2^2 \\ b_3^2 \end{bmatrix} \\ \underline{z}^{(2)} \quad \quad \underline{w}^{(2)} \quad \quad a^{(1)} \quad \quad \underline{b}^{(2)} \end{matrix}$$

Then apply nonlinear activation to obtain:

$$\begin{bmatrix} a_1^2 \\ a_2^2 \\ a_3^2 \end{bmatrix} = \sigma \left(\begin{bmatrix} z_1^2 \\ z_2^2 \\ z_3^2 \end{bmatrix} \right) \Rightarrow a^{(2)} = \sigma(z^{(2)})$$

$$\underline{z}^{(2)} = \underline{w}^{(2)} a^{(1)} + \underline{b}^{(2)} \Rightarrow a^{(2)} = \sigma(\underline{z}^{(2)})$$

From Layer-2 to Layer-3

$$\begin{aligned} z_1^3 &= w_{11}^3 a_1^2 + w_{12}^3 a_2^2 + w_{13}^3 a_3^2 + b_1^3 \\ z_2^3 &= w_{21}^3 a_1^2 + w_{22}^3 a_2^2 + w_{23}^3 a_3^2 + b_2^3 \\ z_3^3 &= w_{31}^3 a_1^2 + w_{32}^3 a_2^2 + w_{33}^3 a_3^2 + b_3^3 \\ z_4^3 &= w_{41}^3 a_1^2 + w_{42}^3 a_2^2 + w_{43}^3 a_3^2 + b_4^3 \end{aligned} \Rightarrow \begin{matrix} \begin{bmatrix} z_1^3 \\ z_2^3 \\ z_3^3 \\ z_4^3 \end{bmatrix} = \begin{bmatrix} w_{11}^3 & w_{12}^3 & w_{13}^3 \\ w_{21}^3 & w_{22}^3 & w_{23}^3 \\ w_{31}^3 & w_{32}^3 & w_{33}^3 \\ w_{41}^3 & w_{42}^3 & w_{43}^3 \end{bmatrix} \begin{bmatrix} a_1^2 \\ a_2^2 \\ a_3^2 \end{bmatrix} + \begin{bmatrix} b_1^3 \\ b_2^3 \\ b_3^3 \\ b_4^3 \end{bmatrix} \\ \underline{z}^{(3)} \quad \quad \underline{w}^{(3)} \quad \quad \underline{a}^{(2)} \quad \quad \underline{b}^{(3)} \end{matrix}$$

Apply nonlinear activation

$$\begin{bmatrix} a_1^3 \\ a_2^3 \\ a_3^3 \\ a_4^3 \end{bmatrix} = \sigma \begin{pmatrix} z_1^3 \\ z_2^3 \\ z_3^3 \\ z_4^3 \end{pmatrix} \Rightarrow a^{(3)} = \sigma(z^{(3)})$$

$$z^{(3)} = W^{(3)} a^{(2)} + b^{(3)} \Rightarrow a^{(3)} = \sigma(z^{(3)})$$

From Layer-3 to Layer-4

$$z_1^{(4)} = w_{11}^4 a_1^3 + w_{12}^4 a_2^3 + w_{13}^4 a_3^3 + w_{14}^4 a_4^3 + b_1^4$$

$$\underset{z^4}{z_1^{(4)}} = \underset{W^4}{\begin{bmatrix} w_{11}^4 & w_{12}^4 & w_{13}^4 & w_{14}^4 \end{bmatrix}} \underset{a^3}{\begin{bmatrix} a_1^3 \\ a_2^3 \\ a_3^3 \\ a_4^3 \end{bmatrix}} + \underset{b^4}{\begin{bmatrix} b_1^4 \end{bmatrix}}$$

Apply activation $a^4 = \sigma(z^4)$

$$z^{(4)} = W^{(4)} a^{(3)} + b^{(4)} \Rightarrow a^{(4)} = \sigma(z^{(4)})$$

Summary

$$z^{(2)} = \sum a^{(1)} + b^{(2)} \Rightarrow a^{(2)} = \sigma(z^{(2)})$$

$$z^{(3)} = \sum a^{(2)} + b^{(3)} \Rightarrow a^{(3)} = \sigma(z^{(3)})$$

$$z^{(4)} = \sum a^{(3)} + b^{(4)} \Rightarrow a^{(4)} = \sigma(z^{(4)})$$

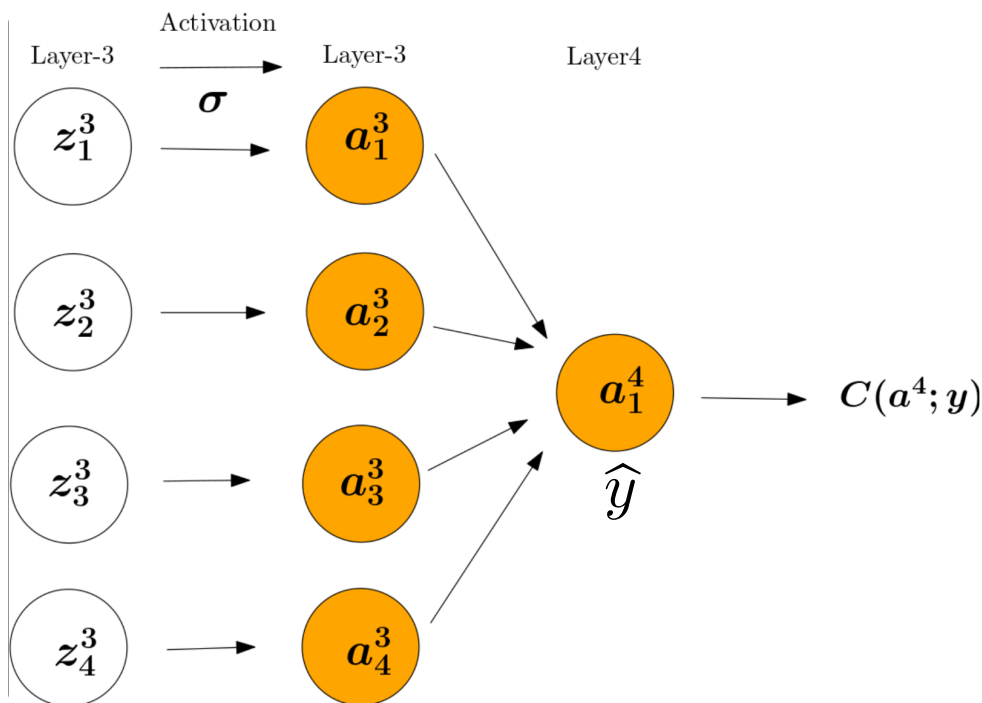
Forward Pass

$$z^l = \sum a^{l-1} + b^l, \quad a^l = \sigma(z^l), \quad l=2,3,\dots,L$$

- In any deep learning publication, you should see pretty much the same form.
- Just like multiple linear regression, we will now compute the derivative of the cost function with respect to all weights w_{jk}^l and biases b_j^l in all layers

② Backpropagation

- Let's consider the picture in the very last layer. We don't need to specify a certain form for the cost function C and the nonlinear activation σ at this point.



- Convince yourself that the cost function C is a function of all weights w_{jk}^l and all biases b_j^l in the network. The ultimate goal is to find an expression for

$$\frac{\partial C}{\partial w_{jk}^l} \quad \text{and} \quad \frac{\partial C}{\partial b_j^l} \quad l=1, 2, \dots, L$$

Let's define

$$\delta_1^{(4)} = \frac{\partial C}{\partial z_1^{(4)}} \rightarrow \text{derivative of cost function with respect to the output of } j^{\text{th}} \text{ neuron in last layer.}$$

In this case we have only 1 neuron. $j=1$

Using chain rule

$$a_1^{(4)} = \sigma(z_1^{(4)}) \Rightarrow \delta_1^{(4)} = \frac{\partial C}{\partial a_1^{(4)}} \frac{\partial a_1^{(4)}}{\partial z_1^{(4)}} = \frac{\partial C}{\partial a_1^{(4)}} \sigma'(z_1^{(4)})$$

In general, this expression is a vector since we have more than one output neuron. We can write

$$\delta^{(4)} = \bar{\nabla}_a C \odot \sigma'(z^{(4)})$$

\odot : elementwise multiplication. $\bar{\nabla}_a C = [\partial C / \partial a^{(4)}]$

Note: Once we know $z_1^{(4)}$, we can compute $\delta_1^{(4)}$, the derivative of cost function with respect to the last layer. Now we will propagate this error to the previous layers

Formula: the actual derivation of the formulas below can be found in my lecture notes I will post

$$\delta^l = (W^{l+1})^T \delta^{l+1} \odot \sigma'(z^l), l = L-1, L-2, \dots, 2$$

Then we have

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

$$\frac{\partial C}{\partial b_j^l} = 1 \cdot \delta_j^l$$

component-wise

\Rightarrow

$$\frac{\partial C}{\partial W^l} = \delta^l (a^{l-1})^T \leftarrow \text{matrix}$$

$$\frac{\partial C}{\partial b^l} = \delta^l \leftarrow \text{vector}$$

layer-wise

Now we can perform gradient descent "layer-wise"

$$\delta^l = (W^{l+1})^T \delta^{l+1} \odot \sigma'(z^l)$$

$$W^{l+1} = W^l - \tau \frac{\partial C}{\partial W^l} = W^l - \eta \cdot \delta^l (a^{l-1})^T \leftarrow \text{matrix}$$

$$b^{l+1} = b^l - \tau \frac{\partial C}{\partial b^l} = b^l - \tau \delta^l \leftarrow \text{vector}$$

-

