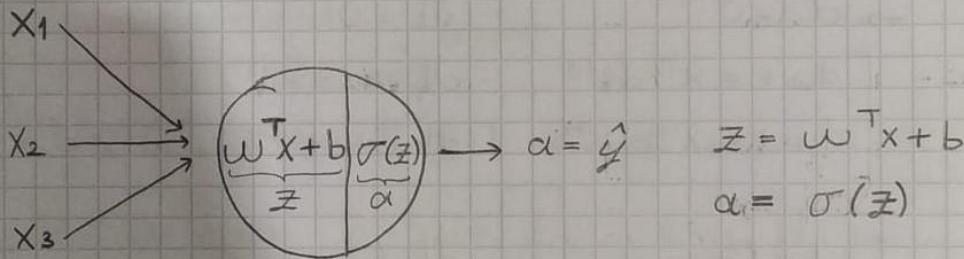
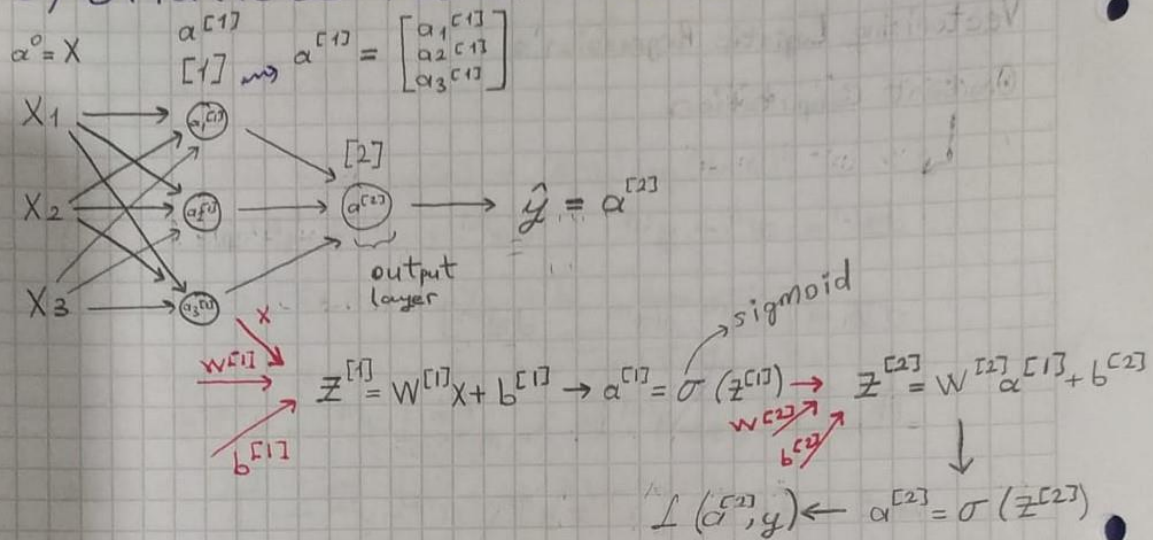
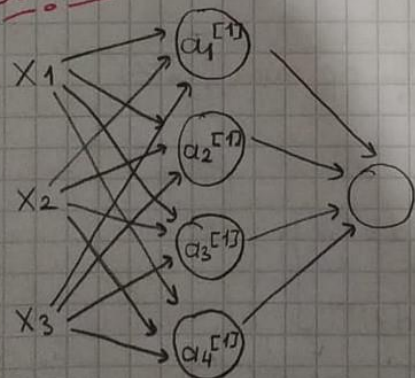


3) Shallow Neural Network



Gösterin



$$z_1^{[1]} = W_1^{[1]T} X + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = W_2^{[1]T} X + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = W_3^{[1]T} X + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = W_4^{[1]T} X + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

$$\begin{bmatrix} W_1^{[1]T} \\ W_2^{[1]T} \\ W_3^{[1]T} \\ W_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} W_1^{[1]T} X + b_1^{[1]} \\ W_2^{[1]T} X + b_2^{[1]} \\ W_3^{[1]T} X + b_3^{[1]} \\ W_4^{[1]T} X + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

Hepsi single Training example için.

Vectorizing across Multiple Training Examples:

$$X \longrightarrow a^{[2]} = \hat{y}$$

$$X^{(1)} \longrightarrow a^{[2]}(1) = \hat{y}^{(1)}$$

$$X^{(2)} \longrightarrow a^{[2]}(2) = \hat{y}^{(2)}$$

⋮

$$X^{(m)} \longrightarrow a^{[2]}(m) = \hat{y}^{(m)}$$

ör.

$$a^{[2]}(i)$$

Refers to training example i.

layer 2.

Different Training Examples

$$X = \begin{bmatrix} X^{(1)} & X^{(2)} & \dots & X^{(m)} \end{bmatrix}$$

↓ Different input features

$$Z = \begin{bmatrix} z^{[1]}(1) & \dots \\ \vdots & \\ z^{[n]}(1) & \dots \end{bmatrix}$$

↓ Hidden layers

Different Training Examples

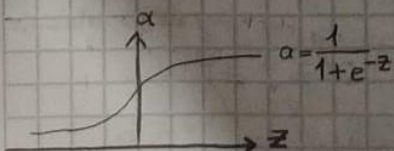
$$A = \begin{bmatrix} a^{[1]}(1) & \dots \\ \vdots & \\ a^{[n]}(1) & \dots \end{bmatrix}$$

↓ Hidden layers

Different training examples

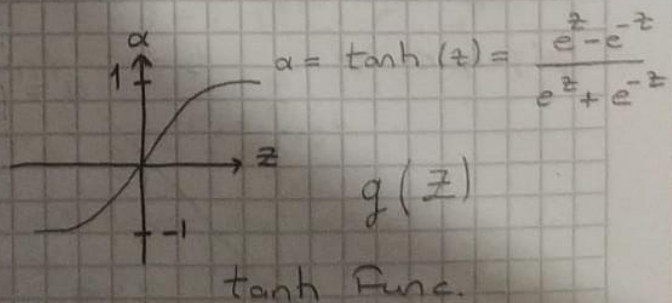
Activation Functions:

Sigmoid Function'den farklı olarak, başka non linear fonksiyonlarda kullanılabilir.



Sigmoid Func.

$$\sigma(z)$$



tanh Func.

tanh, değerleri 0'a center'ladığı için Sigmoid'e göre daha iyidir. tanh'da -1 ile 1 arasında değişir değerler.

→ output layer'da "Sigmoid Func" kullanmak gerekir, hidden'da tanh kullansak bile. Çünkü 0 ile 1 arasında olmasına isteriz Output'un.

* Activation functions can be different for different layers

* Hem tanh hem Sigmoid fonksiyonun z 'nin çok küçük veya çok büyük olduğu durumlarda slope of fonksiyon ends up being close to 0. Bu da Gradient Descent'i yavaşlatır.

→ Örneğin

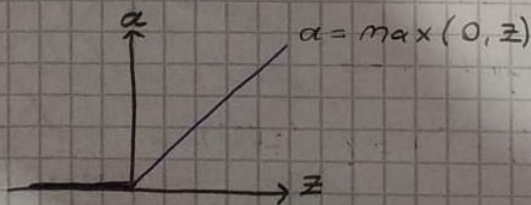
$$W(i) = W(i) - \alpha \frac{dW(i)}{dz}$$

$$\frac{dW(i)}{dz} = \frac{\partial L}{\partial W(i)} = X \cdot dz$$

$dz \rightarrow$ çok küçük veya büyük olduğu durumda Gradient Descent çok yavaşlar.

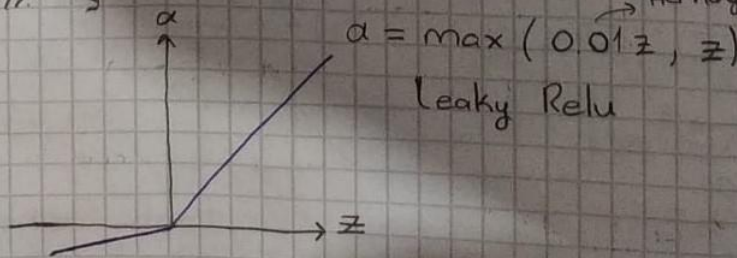
Bu sorunların üstesinden gelmek için;

"Relu" fonksiyonu kullanılır.



Rule of Thumb.
Hidden Layer için Relu kullan;
Output layer için Sigmoid kullan.

Leaky Relu diye bir aktiv. Func var. Relu'dan daha iyi çalışır.



Herhangi değer olabilir.

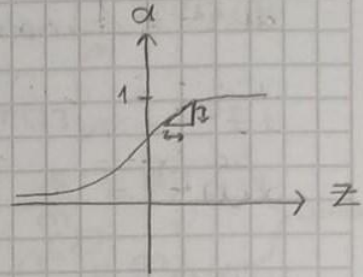
örnek

Örneğin sigmoid func. derivative'sini almak isteyelim.

$$g(z) = \frac{1}{1+e^{-z}}$$

$$\frac{d}{dz} g(z) = \text{slope of } g(x) \text{ at } z.$$

$$= \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right)$$
$$= g(z) (1 - g(z))$$



örneğin Relu func. derivative almak isteyelim:

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Şimdi de Leaky Relu derivative alalım.

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Gradient Descent for Neural Networks (one hidden layer)

Forward Propagation:

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]}) = \sigma(Z^{[2]})$$

↳ sigmoid (output layer)

Backpropagation:

$$dZ^{[2]} = A^{[2]} - Y$$

$$Y = [y^{(1)}, \dots, y^{(m)}]$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$$

↳ element wise product

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims}=\text{True})$$

Summary of Gradient Descent: (one hidden layer)

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims}=\text{True})$$

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims}=\text{True})$$

! Network'e başlarken, w ve b random initialize edilmelidir. (log reg'deki gibi 0'a set edilmemelidir.)

! Başlangıç değerleri 0'a set edilmemelidir, ancak küçük değerlere set edilmelidir.

$$\text{loss func} = \dots$$
$$J(\cdot) = \frac{1}{m} \sum_{i=1}^n L(\hat{y}_i, y_i)$$