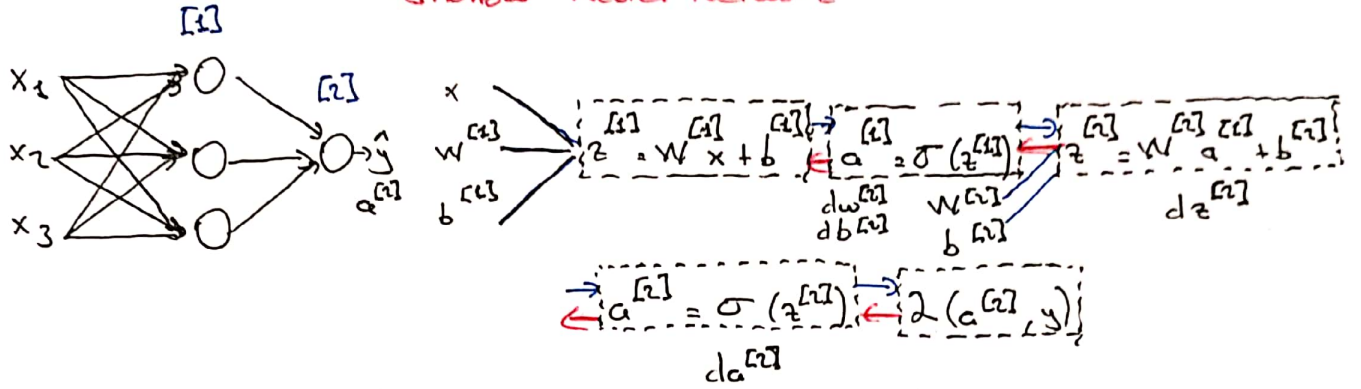
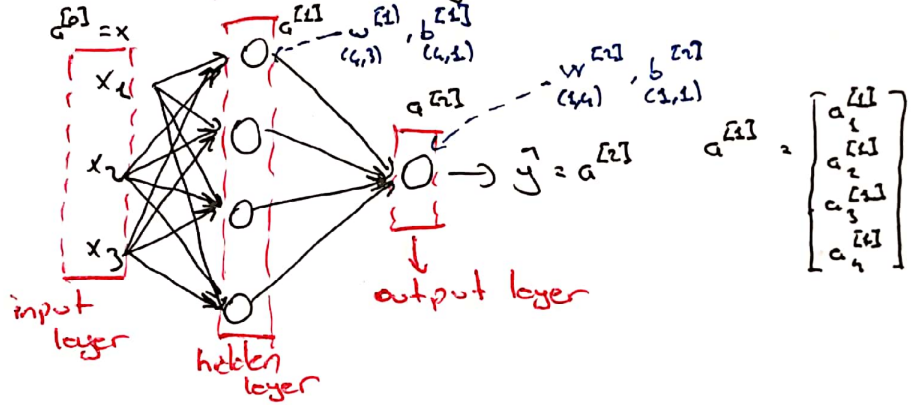


Shallow Neural Network



Sinir ağında her bir nöron hem doğrusal bir fonksiyon hesaplar hemde bir aktivasyon fonksiyonu hesaplar.

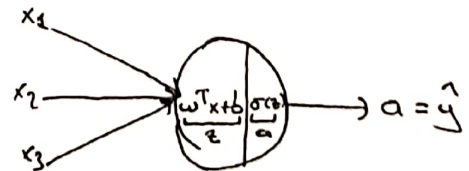
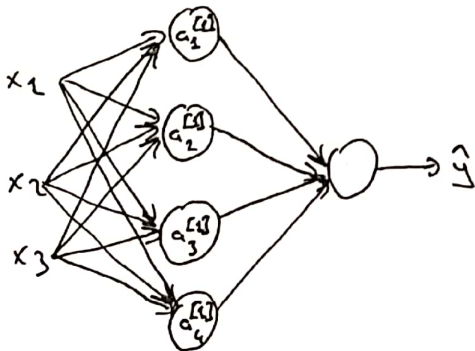
Katmanların gösterimi köşeli parantez ile gösterilmektedir.



Input layer, özellik girdilerimiz. Output layer, çıktı verimiz. Bu iki sütün eğitim örneklerimizde bulunurken hidden layer girdilerimiz bulunmaz. Bu bir nevi neden gizli olarak adlandırıldığını açıklar.

Yukarıdaki örnek sinir ağı 2 katmanlıdır. Çünkü sinir ağına giriş katmanları sayılmaz.

Computing a Neural Network Output



$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, & a_1^{[1]} &= \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, & a_2^{[1]} &= \sigma(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, & a_3^{[1]} &= \sigma(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, & a_4^{[1]} &= \sigma(z_4^{[1]}) \end{aligned}$$

Bir sinir ağında bir katmandaki nöronların nasıl hesaplandığını gördük. Gelin bunun vektörizasyonunu yapalım.

$$z^{[1]} = \underbrace{\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix}}_{W^{[1]}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}}_{b^{[1]}} = \begin{bmatrix} w_1^{[1]T}x + b_1^{[1]} \\ w_2^{[1]T}x + b_2^{[1]} \\ w_3^{[1]T}x + b_3^{[1]} \\ w_4^{[1]T}x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$\begin{aligned} 1 \text{ layer} \Rightarrow z^{[1]} &= W^{[1]}x + b^{[1]} \\ a^{[1]} &= \sigma(z^{[1]}) \end{aligned}$$

$$\begin{aligned} 2 \text{ layer} \Rightarrow z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

Yukarıdaki dört denklemleri kullanarak gizli katmanlı sinir ağını eğitebilirsiniz. Logistic regresyondan çokda büyük bir farkı yoktur.

Multiple Examples Vectorizing

En başta tek bir örnek üzerinde tek katmanlı gizli bir sinir ağının nasıl vektörize edildiğini görmüştük. Şimdi çoklu örnekte nasıl vektörize edilip eğiteceğimize bakalacağız.

$$\begin{aligned} x &\longrightarrow a^{[1]} = \hat{y} \\ x^{(1)} &\longrightarrow a^{1} = \hat{y}^{(1)} \\ x^{(2)} &\longrightarrow a^{[1](2)} = \hat{y}^{(2)} \\ \vdots & \\ x^{(m)} &\longrightarrow a^{[1](m)} = \hat{y}^{(m)} \end{aligned}$$

example:
layer 2

for $i=1$ to m ,

$$z^{[1]}(i) = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1]}(i) = \sigma(z^{[1]}(i))$$

$$z^{[2]}(i) = W^{[2]}a^{[1]}(i) + b^{[2]}$$

$$a^{[2]}(i) = \sigma(z^{[2]}(i))$$

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

$$z^{[1]} = \begin{bmatrix} | & | & & | \\ z^{[1]}(1) & z^{[1]}(2) & \dots & z^{[1]}(m) \\ | & | & & | \end{bmatrix}$$

train examples

$$A^{[1]} = \begin{bmatrix} | & | & & | \\ a^{[1]}(1) & a^{[1]}(2) & \dots & a^{[1]}(m) \\ | & | & & | \end{bmatrix}$$

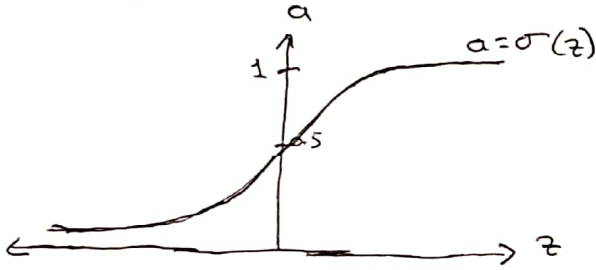
hidden units

$$\begin{aligned} z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= \sigma(z^{[1]}) \\ z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

Activation Functions

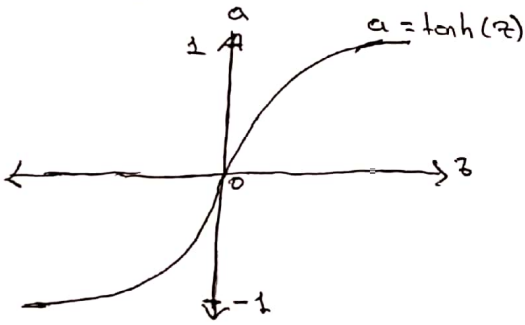
Halihazırda çokça kullandığımız sigmoid fonksiyonu bir aktivasyon fonksiyonudur. Ancak bazı aktivasyon fonksiyonları daha iyi sonuçlar verebiliyor.

1. Sigmoid Function



$$\sigma = \frac{1}{1 + e^{-z}}$$

2. Tangent Function (Hiperbolik, Tanjant)



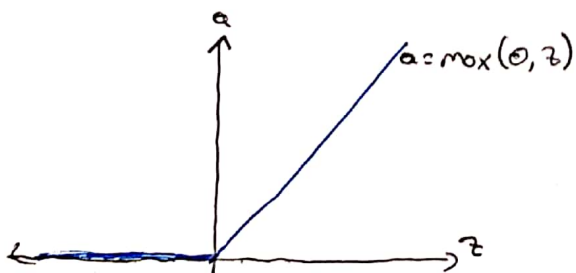
$$a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

- Sigmoid fonksiyonundan daha iyi çalışır.
- Sigmoid fonksiyonunun matematiksel olarak kaydırılmış halidir.
- Verilerini 0 merkeze ortalayacağından (0) bir sonraki katman için öğrenmeyi kolaylaştırır.
- Tangent aktivasyon fonksiyonu her halde

sigmoid fonksiyonun çok üstünde bir verimliliği vardır. Tek istisna ikili bir sınıflandırma yaparsanız değer 0'dan 1 arasında olmalıdır. Ancak tangent bize $-1 \leq a \leq 1$ arasında sonuç döndürür. Bunun için sinir ağımızın son çıkış katmanında sigmoid fonksiyon kullanabilirsiniz.

Not Hem sigmoid fonksiyonunun hemde tanh fonksiyonunun dezavantajlarından biri, z çok büyük veya çok küçükse, bu fonksiyonun gradyan veya türevi (eğim) 0'a yakın olur. Bu da gradyan inişini yavaşlatabilir.

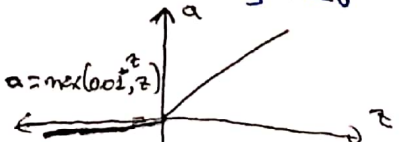
3. ReLU (Rectified Linear Unit)



- Pozitif olduğu sürece (z) lineer 0 'dan Negatif veya 0 ise türevi 0 'dır. ~~0~~ 0'dan itibaren tanımlanır.
- Tek dezavantajı z negatif olduğunda türevi 0 olmasıdır. Bunu ~~çözmek için~~ ^{pozitif bir eğime} inmesine neden olur. Bu duruma Leaky ReLU (sıradan ReLU) deniyor.

• İşlem gücünü oldukça hafifletir.

Leaky ReLU $a = \max(0.01, z)$



Niçin Aktivasyon Fonksiyonlarına İhtiyaç Duyuyoruz?

$$a^{[1]} = z^{[1]} = w^{[1]} x + b^{[1]}$$

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

$$\begin{aligned} a^{[2]} &= w^{[2]} (w^{[1]} x + b^{[1]}) + b^{[2]} \\ &= \underbrace{(w^{[2]} w^{[1]})}_{w'} x + \underbrace{(w^{[2]} b^{[1]} + b^{[2]})}_{b'} \\ &= w' x + b' \end{aligned}$$

Eğer aktivasyon fonksiyonlarını kullanmazsak birçok gizli katmandaki nöronun hesaplanması yine bize bir lineer fonksiyon verir.

Lineer bir gizli katman bir işe yaramaz. Çünkü iki lineer fonksiyonun birleşimi kendisine yarı bir lineer fonksiyona eşittir.

Not Lineer fonksiyon regresyon problemlerinde kullanılabilir.

Daha derin ağırlara ilerlerken onemli fonksiyonlar elde edebilmek için lineer fonksiyonu non-linear şekle çevirerek lineer olmayan problemlere yanıt verebiliriz.

Aktivasyon Fonksiyonlarının Türevi

Yapay sinir ağı için geri yayılımı (back propagation) uygularken, aktivasyon fonksiyonlarının türevlerini hesaplayabilmemiz gerekir.

1. Sigmoid Function

$$\begin{aligned} g(z) &= \frac{1}{1+e^{-z}} \\ \frac{d}{dz} g(z) &= -\frac{\frac{d}{dz} [e^{-z} + 1]}{(e^{-z} + 1)^2} \quad \text{formül} \\ &= -\frac{e^{-z}}{(1 + e^{-z})^2} \\ &= \left(\frac{1}{1+e^{-z}}\right) \cdot \left(1 - \frac{1}{1+e^{-z}}\right) \end{aligned}$$

$$a(1-a) = g(z) \cdot (1 - g(z))$$

2. Tanh Function

$$\begin{aligned} g(z) &= \frac{e^z - e^{-z}}{e^z + e^{-z}} \\ g'(z) &= \frac{(e^z + e^{-z}) \cdot (e^z - e^{-z}) - (e^z - e^{-z}) \cdot (e^z + e^{-z})}{(e^z + e^{-z})^2} \\ &= \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} \\ &= 1 - \frac{(e^z - e^{-z})^2}{(e^z + e^{-z})^2} \end{aligned}$$

$$= 1 - (g(z))^2 = 1 - a^2$$

3. ReLU ve Leaky ReLU

$$\begin{aligned} g(z) &= \max(0, z) \\ g'(z) &= \begin{cases} 0, & z < 0 \\ 1, & z \geq 0 \end{cases} \end{aligned}$$

$$\begin{aligned} g(z) &= \max(0.01z, z) \\ g'(z) &= \begin{cases} 0.01 & z < 0 \\ 1 & z \geq 0 \end{cases} \end{aligned}$$

Gradient Descent for Neural Networks

Parameters: $W^{[1]}$, $b^{[1]}$, $W^{[2]}$, $b^{[2]}$ $n_x = n^{[0]}$, $n^{[1]}$, $n^{[2]} = 1$
 $(n^{[1]}, n^{[0]})$ $(n^{[1]}, 1)$ $(n^{[2]}, n^{[1]})$, $(n^{[2]}, 1)$ katmandeki değişim sayısı

$$\text{Cost Function: } J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, \hat{y}_i)$$

Gradient Descent:

Repeat {

Compute predicts $(\hat{y}^{(i)}, i = 1 \dots m)$

$$d\omega^{[1]} = \frac{d\mathcal{F}}{d\omega^{[1]}} \quad , \quad db^{[1]} = \frac{d\mathcal{F}}{db^{[1]}}$$

$$W^{[1]} = W^{[1]} - \alpha dW^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha d b^{[1]}$$

$$dw^{[2]} = \frac{d\mathcal{F}}{dw^{[2]}} \quad , \quad db^{[2]} = \frac{d\mathcal{F}}{db^{[2]}}$$

$$\psi^{[2]} = \psi^{[2]} - \alpha dW^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha d b^{[2]}$$

3

Forward Propagation

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

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

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

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

Backpropagation

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

$$d_{\omega}^{[22]} = \frac{1}{3} d_{\alpha}^{[21]} A^{[13]T}$$

$$d1b^{[2]} = \frac{1}{3} \text{np.sum}(d2^{[2]}, \text{axis} = 1, \text{keepdims} = \text{True})$$

$$d\tau^{[1]} = \underbrace{\omega^{[1]T} d\tau^{[1]}}_{(A^{[1]}, m)} * \underbrace{g^{[1]}(\tau^{[1]})}_{(n^{[1]}, m)}$$

$$d\omega^{[1]} = \frac{1}{3} dz^{[1]} x^T$$

$$d6^{[1]} = \frac{1}{3} \cdot \text{np.sum}(d7^{[1]}, \text{axis}=2, \text{keepdims} = \text{True})$$

Backpropagation Sergei



$$dz = a - y \quad \text{or} \quad dz = \frac{dL(a, y)}{da} = -y \log a - (1-y) \log(1-a)$$

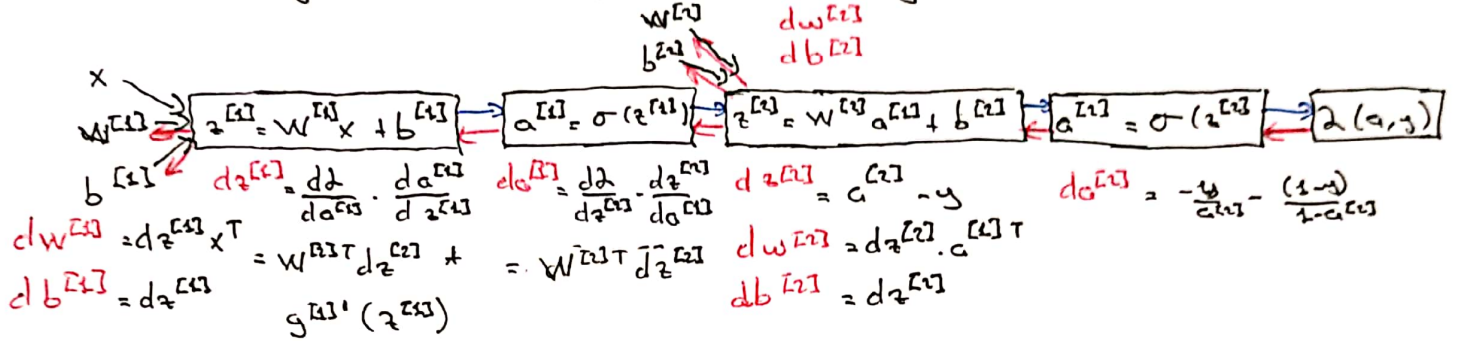
$$dz = da \cdot g'(z) = \frac{dh}{da} \cdot \frac{da}{dz}$$

$$dw = dz \cdot x$$

$$db = dz$$

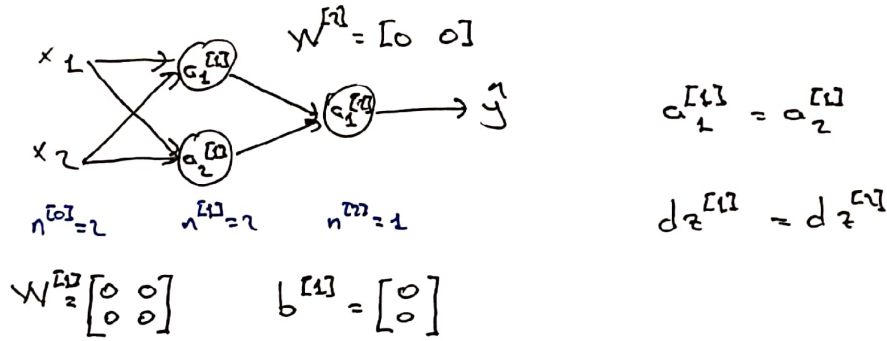
$$= \frac{d}{dt} \cdot \frac{d}{dz} g(z) = \frac{1-y}{y} + \frac{1-y}{1-a}$$

Lojistik regresyonda back propagation tek seferde gerçekleşir. Ancak sinir ağlarında katman seviyesi arttıkça işlem miktarıda artar. Gelin 2 katmanlı bir sinir ağı back propagation'ını hesaplayalım.



Parametre Ağırlıklarının Random Ayarlanması

Lojistik regresyonda ağırlıkların sıfıra ayarlanması ile gradient çakışır ancak sinir ağında uygulanan gradyan inişi çok fazla işe yarar.



Eğer ağırlık parametrelerini aynı değerler ile başlatırsanız. Her iki nöronda birbirinin simetrisi olacaktır. Nöronlar aynı değerleri göleyeceğinden $z^{[1]}$ ve $z^{[2]}$ türederide aynı çıkar.

$$dw^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad W^{[1]} = W^{[1]} - \alpha dw$$

$dw^{[1]}$ 'in türevi yukarıdaki gibi aynı iki satırdan oluşacak ve parametreleri güncellediğimizde benzerlik devam edecek ve nöronlar iterasyonlar geçse de simetrikliğini koruyacak. Bu durum sinir ağını amaçsız bir durum haline getiriyor. Bunun çözümü parametreleri rastgele başlatmaktır. Ağırlıkları rastgele başlattığımızı da bias değerimizi 0 başlatma sorun test etmiyor.

$$W^{[1]} = np.random.randn((2, 2)) * 0.01$$

$$b = np.zeros((2, 1))$$

~~Not~~ Parametre ağırlıklarını 0.01 gibi küçük bir sayı ile çarparak küçültmemizin nedeni eğer çok büyük bir sayı ile çarparsak sigmoid fonksiyonuna sokacağımız "z" çok büyük veya çok küçük olur bu da ağırlığın 0'a yaklaşmasından gradyan inişini yavaşlatır. Bu durum sigmoid ve tanh aktivasyon fonk. larında geçerlidir.