



# Deep Learning

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir



<https://www.aparat.com/mehran.safayani>



[https://github.com/safayani/deep\\_learning\\_course](https://github.com/safayani/deep_learning_course)



Department of Electrical and computer engineering, Isfahan university of technology, Isfahan, Iran

# **Basics of Neural Network Programming**

Binary Classification

# Binary classification

		Blue				
	Green		255	134	93	22
Red		255	134	202	22	2
		255	231	42	22	4
		123	94	83	2	192
		34	44	187	92	34
		34	76	232	124	94
		67	83	194	202	

$$\vec{x} = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 254 \\ 253 \\ 250 \\ 220 \end{bmatrix} \quad 64 \times 64 \times 3 = \underbrace{12288}_{n=n_x}$$

$\vec{x} \rightarrow \text{model} \rightarrow \hat{y}$

- Notation

$$(\vec{x}, y) \quad x \in R^{n_x}, y \in \{0,1\}$$

# Binary classification

- $m$  training example:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

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

$n_x$

$m$

$$X \in \mathbb{R}$$

$$X.\text{shape} = (n_x, m)$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$

# Logistic Regression

- Given  $x$ , output  $\hat{y} = P(y=1|x)$   $0 \leq \hat{y} \leq 1$

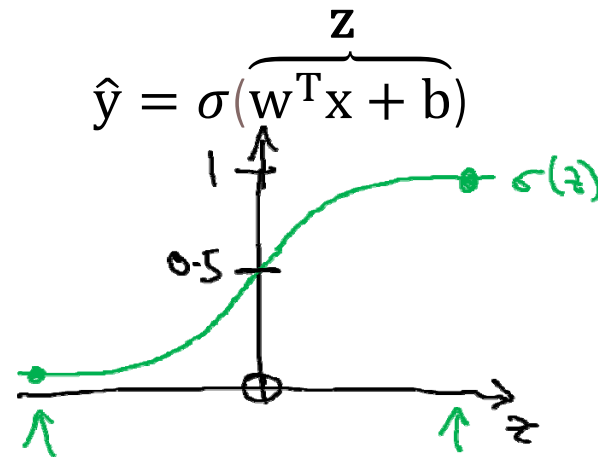
$x \in \mathbb{R}^{n_x}$  parameters:  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$

$$\hat{y} = w^T x + b$$

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

if  $z$  large  $\sigma(z) \approx 1$

if  $z$  large negative  $\sigma(z) \approx 0$



# Logistic Regression

- $\hat{y} = \sigma(\underbrace{w^T x}_z + b)$

$$x_0 = 1, x \in \mathbb{R}^{n_x+1}$$

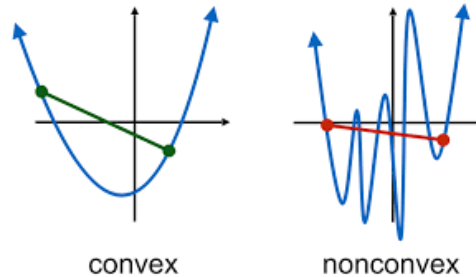
$$\hat{y} = w^T x$$

$$W = \begin{bmatrix} b = w_0 \\ w_1 \\ w_2 \\ w_3 \\ \vdots \\ \vdots \\ \vdots \\ w_{nx} \end{bmatrix} \begin{matrix} \left. \vphantom{\begin{bmatrix} b = w_0 \\ w_1 \\ w_2 \\ w_3 \\ \vdots \\ \vdots \\ \vdots \\ w_{nx} \end{bmatrix}} \right\} b \\ \left. \vphantom{\begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ \vdots \\ \vdots \\ w_{nx} \end{bmatrix}} \right\} w \end{matrix}$$

# Logistic Regression cost function

- Loss (error) function:  $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2 = \frac{1}{2} (\sigma(w^T x + b) - y)^2$  SE: Square Error

- Non-convex graph:



<https://mlstory.org/optimization.html>

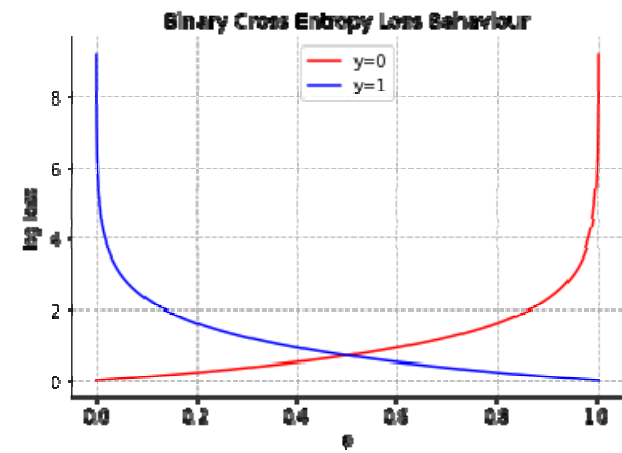
- In mathematics, a real-valued function is called **convex** if the line segment between any two distinct points on the graph of the function lies above the graph between the two points.
- Equivalently, a function is convex if its epigraph (the set of points on or above the graph of the function) is a convex set.

# Cross Entropy

- $L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$

if  $y=1$  :  $L(\hat{y}, y) = -\log \hat{y}$

if  $y=0$  :  $L(\hat{y}, y) = -\log(1-\hat{y})$



<https://datamonje.com/classification-loss-functions/>

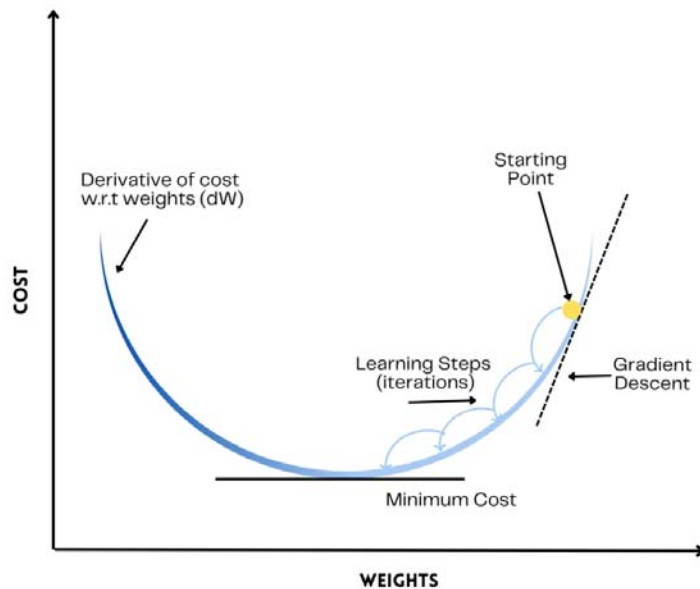
- Cost function: 
$$J(w,b) = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)})$$
$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$



# Basics of Neural Network Programming

## Gradient Descent

# Gradient Descent



- 1)  $\alpha > 0$

Repeat{

$$w = w - \alpha \frac{dJ(w)}{dw}$$

}until convergence

$$w = w - \alpha dw$$

$$w^* = (x^T x)^{-1} x^T y$$

$$y = (x - \frac{n^3}{1})^2$$

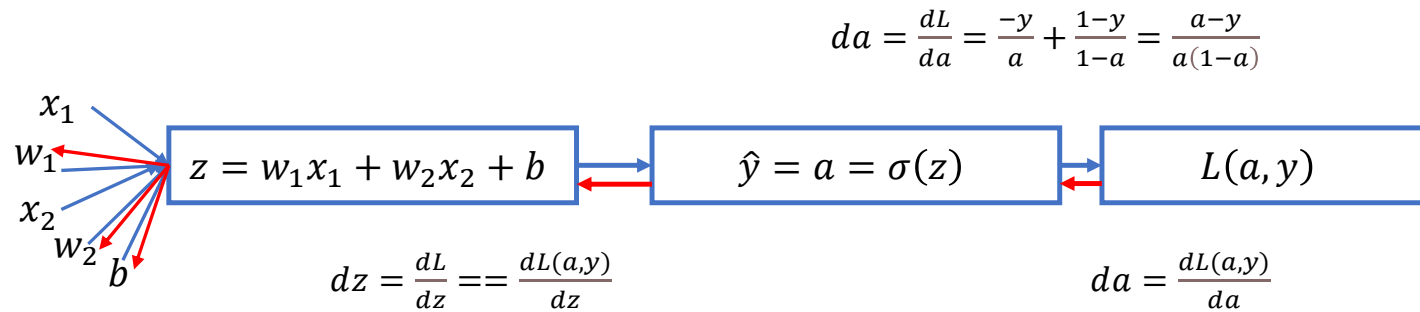
$$\frac{dy}{dx} = 2(x - 1) = 0$$

$$x = 1$$

- $z = w^T x + b$
- $\hat{y} = a = \sigma(z)$
- $L(a, y) = -(y \log a + (1 - y) \log(1 - a))$

# Gradient Descent

## Computational Graph



$$dz = a - y_{a(1-a)}$$

$$dz = \frac{dL}{da} \times \frac{\vec{da}}{dz} = \frac{dL}{dz}$$

$$a = \sigma(z)$$

$$\frac{da}{dz} = \sigma'(z) = \underbrace{\sigma(z)}_a \underbrace{(1 - \sigma(z))}_{1-a}$$

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

$$\frac{dL}{da} = da = \frac{a-y}{a(1-a)} \times a(1-a) = a - y$$

$$dw_1 = \frac{dL}{dw_1} = \frac{dL}{dz} \times \underbrace{\frac{dz}{dw_1}}_{x_1} = x_1 dz$$

$$dw_2 = x_2 dz$$

$$db = \frac{dL}{db} = \frac{dL}{dz} \times \underbrace{\frac{dz}{db}}_1 = dz$$

# Gradient Descent

$$\bullet \left\{ \begin{array}{l} w_1 = w_1 - \alpha dw_1 \\ w_2 = w_2 - \alpha dw_2 \\ b = b - \alpha db \end{array} \right\} \text{بروزرسانی همزمان}$$

$$\bullet \left\{ \begin{array}{l} w_1temp = w_1 - \alpha dw_1 \\ w_2temp = w_2 - \alpha dw_2 \\ btemp = b - \alpha db \end{array} \right.$$

$$\bullet \left\{ \begin{array}{l} w_1 = w_1temp \\ w_2 = w_2temp \\ b = btemp \end{array} \right.$$

$$\bullet J(w,b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

$$\bullet a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(wx^{(i)} + b)$$

$$\bullet dw_1^{(i)} \quad dw_2^{(i)} \quad db^{(i)}$$

$$\bullet \underbrace{dJ(w,b)}_{dw_1} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{dL(a^{(i)}, y^{(i)})}{dw_1}}_{dw_1^{(i)}}$$

# Logistic regression on $m$ examples

$J = 0;$       $dw_1 = 0;$   $dw_2 = 0;$   $db = 0;$

$w_1 \leftarrow \text{random}$       $w_2 \leftarrow \text{random}$       $b \leftarrow \text{random}$

**Repeat**{

**For**      $i=1$      **to**      $m$

$$z^{(i)} = w^T x^{(i)} + b$$

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

$$J += [y^{(i)} \text{Log} a^{(i)} + (1 - y^{(i)}) \text{Log}(1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J /= m;$$

$$dw_1 /= m;$$

$$dw_2 /= m;$$

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

**} until convergence**

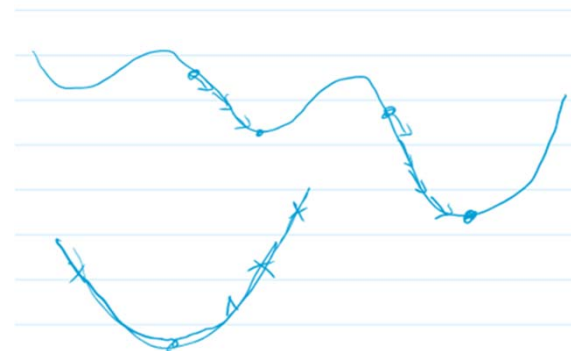
$$d\theta = \begin{bmatrix} dw_1 \\ dw_2 \\ db \end{bmatrix}$$

$$\|d\theta\| \leq \varepsilon = 10^{-4}$$

$$\theta^t = \begin{bmatrix} w_1^t \\ w_2^t \\ b^t \end{bmatrix}$$

$$\theta^{t+1} = \begin{bmatrix} w_1^{t+1} \\ w_2^{t+1} \\ b^{t+1} \end{bmatrix}$$

$$\|\theta^{t+1} - \theta^t\|_2 \leq \varepsilon$$



# Basics of Neural Network Programming

Vectorizing Logistic Regression's Gradient Computation

# Vectorizing Logistic Regression

- $Z=0;$

*For*  $i$  *in*  $\text{range}(n_x)$   
 $z += w[i] * x[i]$   
 $z += b$

روش اول

- $Z=0;$

$z = \text{np.dot}(w, x) + b$

روش دوم  
SIMD  
GPU

# Vectorizing Logistic Regression

$$\bullet \begin{array}{lll} z^{(1)} = w^T x^{(1)} + b & z^{(2)} = w^T x^{(2)} + b & z^{(3)} = w^T x^{(3)} + b \\ a^{(1)} = \sigma(z^{(1)}) & a^{(2)} = \sigma(z^{(2)}) & a^{(3)} = \sigma(z^{(3)}) \end{array}$$

$$\bullet X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \in R^{n_x \times m} \quad \underbrace{[w_1 \dots w_{n_x}]}_{W^T} \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$$\bullet Z = \begin{bmatrix} \underbrace{z^{(1)}}_{w^T x^{(1)} + b} & \underbrace{z^{(2)}}_{w^T x^{(2)} + b} & \underbrace{z^{(3)}}_{w^T x^{(3)} + b} & \dots & \underbrace{z^{(m)}}_{w^T x^{(m)} + b} \end{bmatrix} = W^T X + [b \ b \ \dots \ b]_{1 \times m}$$

$$\bullet Z = \underbrace{np \cdot dot(W \cdot T, X)}_{1 \times m} + \underbrace{b}_{(1,1)} \text{ "broadcasting" } [b, b \dots, b]_{1 \times m}$$



# Vectorizing Logistic Regression

- $A = [a^{(1)}, a^{(2)}, \dots, a^{(m)}] = \sigma(\underbrace{Z}_{1 \times m})$

- $dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)}$

- $dZ = [dz^{(1)} \ dzz^{(2)} \ \dots \ dz^{(m)}]_{1 \times m}$

- $A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] \quad Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$

- $dZ = A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots \ a^{(m)} - y^{(m)}]$

- $\begin{bmatrix} dw_1 \\ dw_2 \\ dw_3 \\ \vdots \\ dw_{nx} \end{bmatrix}_{nx \times 1} \quad \begin{cases} dw=0 \\ dw+=\widetilde{x^1} \overbrace{dz^1}^{عدد \ a^1-y^1 \text{ بردار}} \\ dw+=x^2 dz^2 \\ dw+=xmdzm \\ dw/=m \end{cases} \quad \begin{cases} db=0 \\ db+=\overbrace{dz^1}^{عدد \ a^1-y^1} \\ db+=dz^2 \\ db+=dzm \end{cases} \quad db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} = \frac{1}{m} np.sum(dZ)$

$$dz = [dz^1 \ dzz^2 \ \dots \ dz^m]_{1 \times m}$$

$$\boxed{dW = \frac{1}{m} X dZ^T} = \frac{1}{m} \begin{bmatrix} x^{(1)}x^{(2)} & \dots & \dots & \dots & x^{(m)} \end{bmatrix}_{nx \times m}$$

$$dw = \frac{1}{m} [x^{(1)}dz^1 + x^{(2)}dz^2 + \dots + x^{(m)}dz^m]$$

$$\begin{bmatrix} dz^1 \\ dz^2 \\ dz^3 \\ \vdots \\ dz^m \end{bmatrix}_{m \times 1}$$

# Vectorizing Logistic Regression

- $w_1, w_2, b \leftarrow \text{random}$

For  $\text{iter in range}(1000)$

1 epoch

$$Z = W^T x + b = \text{np} \cdot \text{dot}(W.T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

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

$$db = \frac{1}{m} \text{np.sum}(dZ)$$

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

$$a = \text{np} \cdot \text{random} \cdot \text{randn}(5,1)$$

$$a.\text{shape} = (5,1)$$

# Logistic Regression Cost function

- $\hat{y} = \sigma(w^T x + b) \quad 0 < \sigma(z) = \frac{1}{1+e^{-z}} < 1$
- $\hat{y} = p(y = 1|x)$
- $\left. \begin{array}{l} \text{if } y = 1 : p(y|x) = \hat{y} \\ \text{if } y = 0 : p(y|x) = 1 - \hat{y} \end{array} \right\} p(y|x)$
- $\boxed{p(y|x) = \hat{y}^y + (1 - \hat{y})^{1-y}}$  *distribution? Bernoulli*
- $\text{if } y = 1 : p(y|x) = \hat{y}$
- $\text{if } y = 0 : p(y|x) = 1 - \hat{y}$
- $\log p(y|x) = \log[\hat{y}^y + (1 - \hat{y})^{1-y}] = y \log \hat{y} + (1 - y) \log(1 - \hat{y})$   
 $= -L(\hat{y}, y)$  *Max Likelihood*

# Logistic Regression Cost function

- $\log P(\text{labels in trainingset}) = \log \prod_{i=1}^m P(y^i | x^i)$
- $\log P(\dots) = \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}) = -\sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$
- Cost function 
$$\underbrace{J(w, b)}_{\text{minimize}} = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

# Neural Networks

