# Deep Learning
## Logistic Regression Classifier

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir

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

https://github.com/safayani/deep_learning_course

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

# Supervised Learning

| Input(x) | Output (y) | Application |
|---|---|---|
| Home features | Price | Real Estate |
| Ad, user info | Click on ad? (0/1) | Online Advertising |
| Image | Object (1,…,1000) | Photo tagging |
| Audio | Text transcript | Speech recognition |
| English | Chinese | Machine translation |
| Image, Radar info | Position of other cars | Autonomous driving |

# Basics of Neural Network Programming

Binary Classification

# Binary Classification



$\longrightarrow$  1 (cat) vs 0 (non cat)

Blue

Green

Red

| 255 | 134 | 93 | 22 |
|-----|-----|-----|-----|
| 255 | 134 | 202 | 22 | 2 |
| 255 | 231 | 42 | 22 | 4 | 30 |
| 123 | 94 | 83 | 2 | 192 | 124 |
| 34 | 44 | 187 | 92 | 34 | 142 |
| 34 | 76 | 232 | 124 | 94 |
| 67 | 83 | 194 | 202 |

# Binary classification



$$\vec{x} = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 254 \\ 253 \\ 250 \\ 220 \end{bmatrix}$$

$64 \times 64 \times 3 = \underbrace{12288}_{n_x}$

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

- Notation

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

# Binary classification

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

- $X = \begin{bmatrix} x^{(1)} & x^{(2)} & ... ... ... & x^{(m)} \end{bmatrix} n_x$

$m$

$X \in \mathbb{R}^{n_x \times m}$

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

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

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

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

# Logistic Regression

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

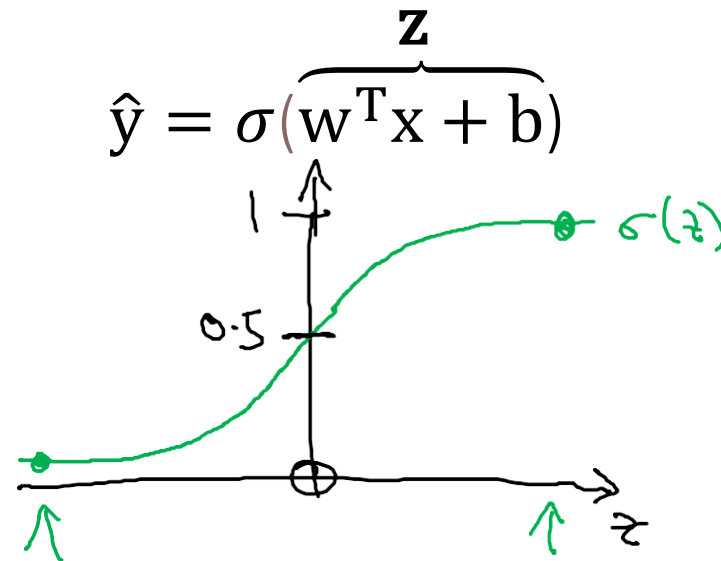  $x \in \mathbb{R}^{n_x}$      parameters: $w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$

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

$$\hat{y} = \sigma(\overbrace{w^T x + b}^{z})$$

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

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

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

# Logistic Regression

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

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

$$W = \begin{bmatrix} b = w_0 \\ w_1 \\ w_2 \\ w_3 \\ . \\ . \\ . \\ w_{n_x} \end{bmatrix} \begin{array}{l} \} \text{ b} \\ \\ \\ \\ \} \text{ w} \\ \\ \\ \end{array}$$

# 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

- What is the problem?
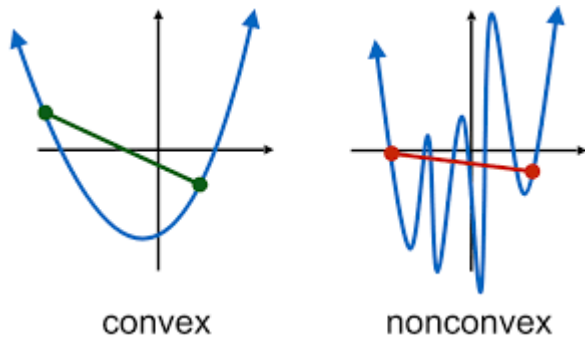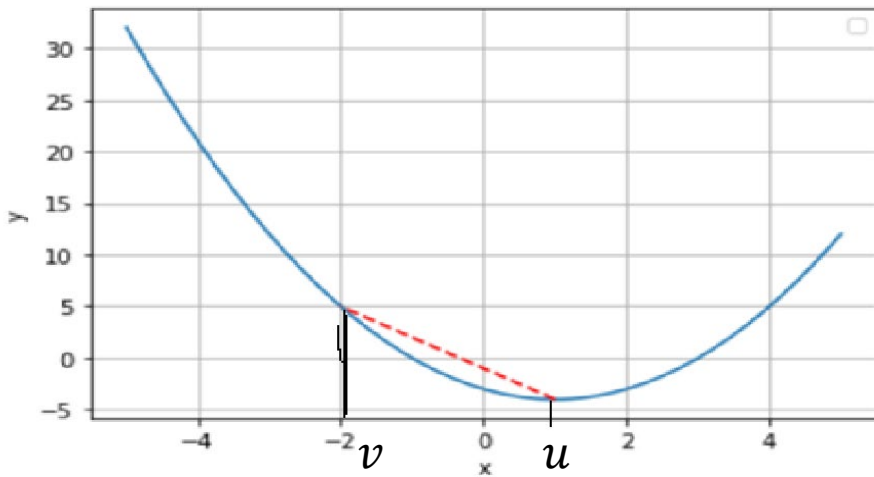
- .

# Convexity



Function h(u) with u∈ X is convex if for any u, v ∈ X and for any $0 \le \lambda \le 1$ we have:

**h($\lambda$u +(1- $\lambda$)v) $\le$ $\lambda$ h(u) + (1- $\lambda$) h(v)**
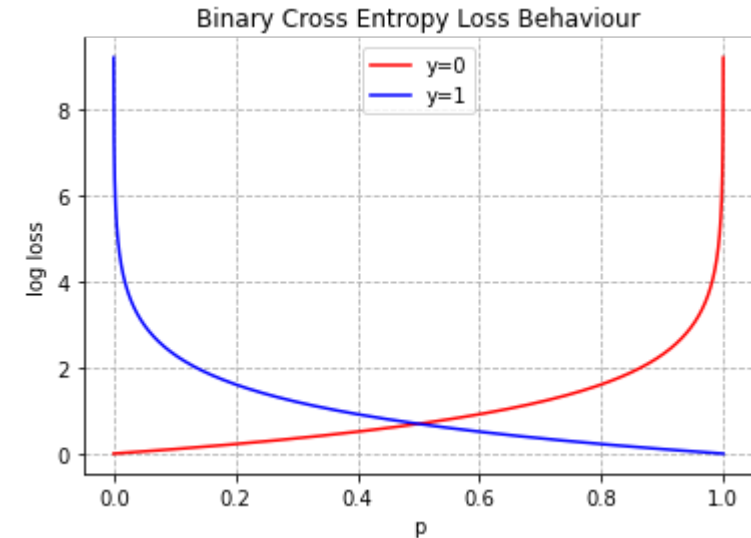
برای توابع محدب هر بهینه محلی یک بهینه سراسری است.



convex        nonconvex

https://mlstory.org/optimization.html

# Cross Entropy

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

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

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



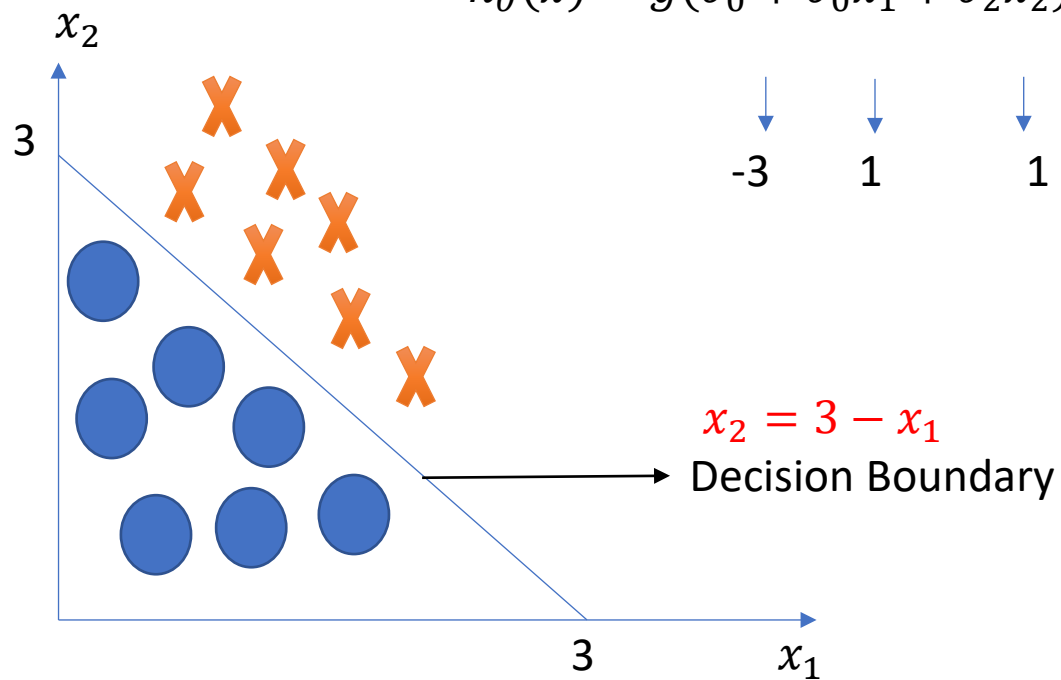https://datamonje.com/classification-loss-functions/

- Cost function: $\qquad J(w,b) = \dfrac{1}{m} \sum_{i=1}^{m} L\left(y^{(i)}, \hat{y}^{(i)}\right)$

$$= -\dfrac{1}{m} \sum_{i=1}^{m} [y^{(i)}\log \hat{y}^{(i)} + (1-y^{(i)})\log(1-\hat{y}^{(i)})]$$

11

# Decision Boundary



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$
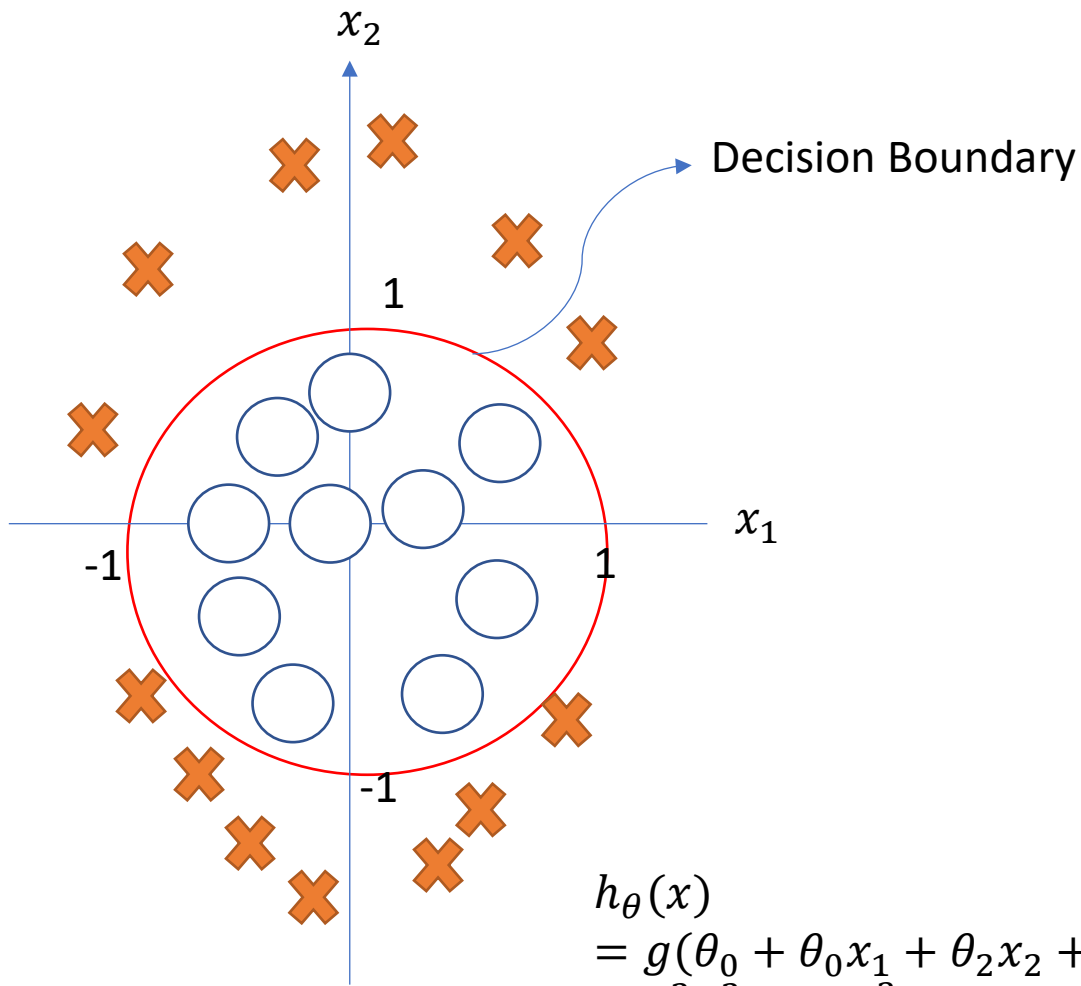
$$h_\theta(x) = g(\theta_0 + \theta_0 x_1 + \theta_2 x_2)$$

-3    1    1

$x_2 = 3 - x_1$

Decision Boundary

Predict    $y = 1, if \underbrace{-3 + x_1 + x_2 \geq 0}_{\theta^T X}$

$x_1 + x_2 \geq 3$

$x_2$

3

3    $x_1$

# Non-Linear Decision Boundaries



Decision Boundary

$$-1 \quad\quad 0 \quad\quad\quad 0$$
$$h_\theta(x) = g(\theta_0 + \theta_0 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

Predict $\quad y = 1, if \underbrace{-3 + x_1^2 + x_2^2 \geq 0}_{\theta^T X}$

$$x_1{}^2 + x_2{}^2 \geq 1$$

$$h_\theta(x) = g(\theta_0 + \theta_0 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 + \dots)$$

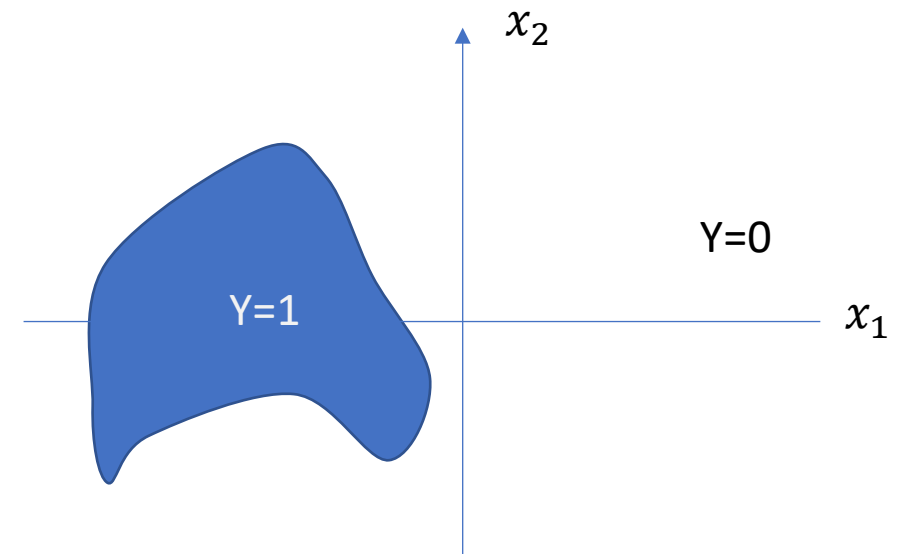Y=0

Y=1

# Gradient Descent

# Cost Function

**Minimize** $J(b, w_1)$

$b, w_1$

If $J(w_1) = (w_1 - 2)^2$

$\dfrac{dJ(w_1)}{dw_1} = 0$ $\Longrightarrow$ $\dfrac{dJ(w_1)}{dw_1} = 2(w_1 - 2) = 0$ $\Longrightarrow$ $w_1 = 2$

# Gradient Descent

**Minimize** $J(\boldsymbol{b}, \boldsymbol{w_1})$

$$\boldsymbol{b}, \boldsymbol{w_1}$$

**Minimize** $J(\boldsymbol{b}, \boldsymbol{w_1}, \ldots, \boldsymbol{w_n})$

$$\boldsymbol{b}, \boldsymbol{w_1}, \ldots, \boldsymbol{w_n}$$

Repeat until convergence:

b= $\boldsymbol{w_0}$

For j=0,...,n

$$w_j = w_j - \alpha \frac{dJ(\boldsymbol{b}, \boldsymbol{w_1}, \ldots, \boldsymbol{w_n})}{d\boldsymbol{w_j}}$$

$\alpha$ **is learning rate**

**Updating all** $\boldsymbol{w_j}$ $\boldsymbol{Simultaneous}$**ly**

Convergence condition:
$$\|W^{t+1} - W^t\|_2 \leq \varepsilon$$

# Gradient Descent

Correct form

$$\text{temp0} = b - \alpha \frac{dJ(b, w_1)}{db}$$

$$\text{temp1} = w_1 - \alpha \frac{dJ(b, w_1)}{dw_1}$$

بروزرسانی همزمان

$b = \text{temp0}$

$w_1 = \text{temp1}$

✔

Incorrect form

$$b = b - \alpha \frac{dJ(b, w_1)}{db}$$

$$w_1 = w_1 - \alpha \frac{dJ(b, w_1)}{dw_1}$$

✖

# Gradient Descent



خطوط مماس نشان داده شده دارای شیب یا مشتق مثبت هستند.
درنتیجه:

$$\frac{dJ(w^1)}{dw^1} > 0 \ , \ \boldsymbol{\alpha} > 0 \implies \boldsymbol{\alpha}\frac{dJ(w^1)}{dw^1} > 0$$

$$\implies w^2 = w^1 - \alpha dw^1$$

$w$ کوچکتر میشود و به سمت چپ حرکت میکنیم.

# Gradient Descent
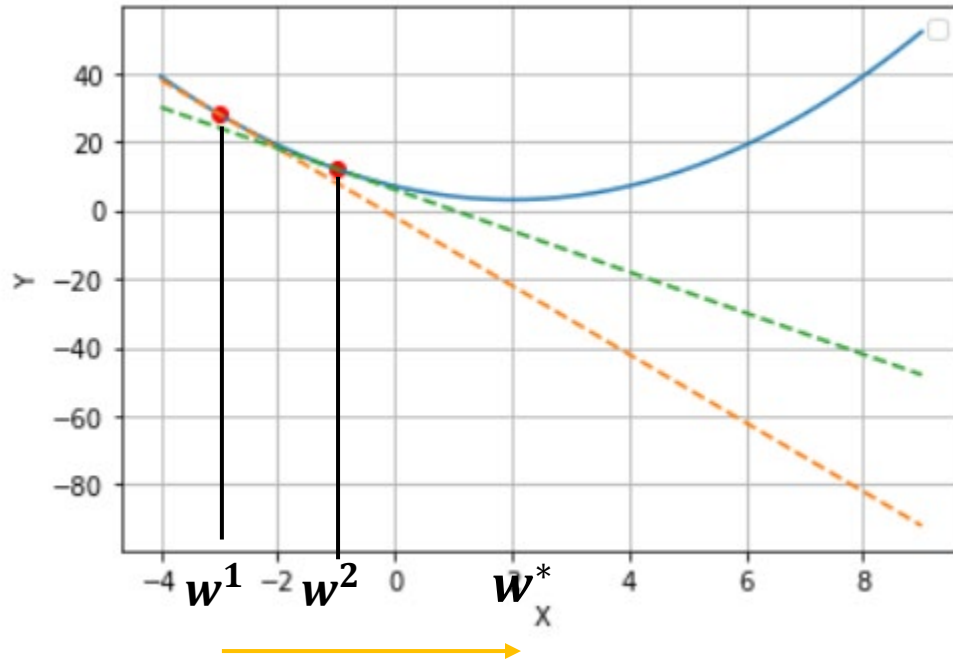


خطوط مماس نشان داده شده دارای شیب یا مشتق مثبت هستند. درنتیجه:

$$\frac{dJ(w^1)}{dw^1} < 0 , \ \alpha > 0 \implies \alpha \frac{dJ(w^1)}{dw^1} < 0$$

$$\implies w^2 = w^1 - \alpha dw^1$$

$w$ بزرگتر میشود و به سمت راست حرکت میکنیم.

# Choosing Learning Rate



$\alpha$ is too large

$\alpha$ is small

# Gradient Descent Weakness

# Gradient Descent



- 1) $\alpha > 0$

  Repeat{

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

  }until convergence

  $$w = w - \alpha dw$$

  $$z = w^T x + b$$

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

# Gradient Descent

$$L(a, y) = -(y \log a + (1-y) \log(1-a))$$

## Computational Graph

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

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

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

$x_1$
$w_1$
$x_2$
$w_2$
$b$

| $z = w_1 x_1 + w_2 x_2 + b$ | $\hat{y} = a = \sigma(z)$ | $L(a, y)$ |

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

$$dw_1 = \frac{dL}{dw_1} = \frac{dL}{\underbrace{dz}_{dz}} \times \frac{dz}{\underbrace{dw_1}_{x_1}} = x_1 \, dz \qquad dw_2 = x_2 \, dz \qquad db = \frac{dL}{db} = \frac{dL}{dz} \times \overbrace{\frac{dz}{db}}^{1} = dz$$

# Gradient Descent

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

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

- $dw_j = \frac{1}{m} \sum_{i=1}^{m} \frac{d L\left(a^{(i)}, y^{(i)}\right)}{dw_j}$

# Logistic regression on $m$ examples

$J = 0; \quad dw_1 = 0; \qquad dw_2 = 0; \qquad db = 0;$

$w_1 \leftarrow \text{ran}dom \quad w_2 \leftarrow \text{ran}dom \quad b \leftarrow \text{ran}dom$

*Repeat{*

    *For   i=1   to   m*

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

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

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

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

      $dw_1 \mathrel{+}= x_1^{(i)} dz^{(i)} \qquad dw_2 \mathrel{+}= x_2^{(i)} dz^{(i)} \qquad db \mathrel{+}= dz^{(i)}$

$J /\mathrel{=} m; \qquad dw_1 /\mathrel{=} m; \qquad\qquad dw_2 /\mathrel{=} m; \qquad\qquad db /\mathrel{=} m;$

$w_1 = w_1 - \alpha\, dw_1 \quad w_2 = w_2 - \alpha\, dw_2 \qquad b = b - \alpha\, db$

   *} until convergence*

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

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

$$dw = \begin{bmatrix} dw_1 \\ dw_2 \\ db \end{bmatrix}$$

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

## What's wrong with the code?

# Logistic regression on $m$ examples

$w_1 \leftarrow random \qquad w_2 \leftarrow random \qquad b \leftarrow random$

**Repeat{**

$J = 0; \quad dw_1 = 0; \qquad dw_2 = 0; \qquad db = 0;$

$\qquad For \qquad i=1 \qquad to \qquad m$

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

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

$\qquad\qquad J \mathrel{+}= \left[ y^{(i)} Log a^{(i)} + (1 - y^{(i)}) Log(1 - a^{(i)}) \right]$

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

$\qquad\qquad dw_1 \mathrel{+}= x_1^{(i)} dz^{(i)} \qquad dw_2 \mathrel{+}= x_2^{(i)} dz^{(i)} \qquad db \mathrel{+}= dz^{(i)}$

$J \mathrel{/}= m; \qquad dw_1 \mathrel{/}= m; \qquad dw_2 \mathrel{/}= m; \qquad db \mathrel{/}= m;$

$w_1 = w_1 - \alpha \, dw_1 \qquad w_2 = w_2 - \alpha \, dw_2 \qquad b = b - \alpha \, db$

**} until convergence**

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

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

$$dw = \begin{bmatrix} dw_1 \\ dw_2 \\ db \end{bmatrix}$$

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

# Vectorizing Logistic Regression

- $z=0;$

  $For\ i\ in\ range(n\_x)$

  $\quad z\ +=\ w[i] * x[i]$

  $\quad z\ +=\ b$

  First method

- $z=0;$

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

  Second method

  **SIMD**

  **GPU**

# Vectorizing Logistic Regression

$$z^{(1)} = w^T x^{(1)} + b \quad z^{(2)} = w^T x^{(2)} + b \qquad\qquad\qquad z^{(m)} = w^T x^{(m)} + b$$

$$a^{(1)} = \sigma\left(z^{(1)}\right) \qquad a^{(2)} = \sigma\left(z^{(2)}\right) \qquad\qquad\qquad a^{(m)} = \sigma\left(z^{(m)}\right)$$

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

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

- $Z = \underbrace{np \cdot dot(w \cdot T, X)}_{1 \times m} + \underbrace{b}_{(1,1)} \quad ''broadcasting'' \quad [b, b \ldots, b]_{1 \times m}$

# Vectorizing Logistic Regression

- $A = \left[a^{(1)}, a^{(2)}, \ldots, a^{(m)}\right] = \sigma(\underset{1\times m}{Z})$

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

- $dZ = \left[dz^{(1)}\ dz^{(2)}\ \ldots\ dz^{(m)}\right]_{1\times m}$

- $A = \left[a^{(1)}\ a^{(2)}\ \ldots\ a^{(m)}\right] \qquad\qquad Y = \left[y^{(1)}\ y^{(2)}\ \ldots\ y^{(m)}\right] \qquad 1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ \vdots \\ \vdots \\ \vdots \\ 1 \end{bmatrix}_{m\times 1}$

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

- $\begin{bmatrix} dw_1 \\ dw_2 \\ dw_3 \\ \vdots \\ dw_{nx} \end{bmatrix}_{nx\times 1}$

$\begin{cases} \bullet\ \textit{dw=0} \\ \qquad\qquad \overset{\text{عدد}}{\phantom{x}} \\ \qquad\qquad \overset{\text{بردار } a^1-y^1}{\phantom{x}} \\ \textit{dw+=}\ \widetilde{x^1}\ \widetilde{dz^1} \\ \textit{dw+=x²dz²} \\ \textit{dw+=x}^m dz^m \\ \textit{dw/=m} \end{cases}$
$\begin{cases} \textit{db=0} \\ \qquad \overset{\text{عدد}}{\phantom{x}} \\ \qquad \overset{a1-y1}{\phantom{x}} \\ \textit{db+=}\ \widetilde{dz1} \\ \textit{db+=dz2} \\ \textit{db+=dz}^m \\ \textit{db/=m} \end{cases}$
$db = \frac{1}{m}\sum_{i=1}^{m} dz^{(i)} = \frac{1}{m} np\,.\,sum(dz) \qquad db = \frac{1}{m} dz \times 1$

$dz = \left[dz^1\ dz^2\ \ldots\ dz^m\right]_{1\times m}$

$dw = \frac{1}{m}\left[x^{(1)}dz^1 + x^{(2)}dz^2 + \cdots + x^{(m)}dz^m\right] \qquad\qquad \boxed{dw = \frac{1}{m} X\,dz^T} = \frac{1}{m}\begin{bmatrix} x^{(1)}x^{(2)} & \ldots\ldots\ldots & x^{(m)} \end{bmatrix} \times \begin{bmatrix} dz^1 \\ dz^2 \\ dz^3 \\ \vdots \\ \vdots \\ dz^m \end{bmatrix}$

# Vectorizing Logistic Regression

- $w_1, w_2, b \leftarrow random$

*For iter in range(1000)*

**1 epoch**

$$Z = W^T X + b = np \cdot dot(w.T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dw = \frac{1}{m} X \, dz^T$$

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

$$w = w - \alpha \, dw$$

$$b = b - \alpha \, db$$

$$w = np \cdot random \cdot \text{rand}n(n_x, 1)$$

# Logistic Regression Cost function

- $\hat{y} = \sigma(w^T x + b)$        $0 < \sigma(z) = \frac{1}{1+e^{-z}} < 1$

- $\hat{y} = p(y = 1 | x)$

- $if \quad y = 1 \quad : \quad p(y|x) = \hat{y}$

- $if \quad y = 0 \quad : \quad p(y|x) = 1 - \hat{y}$    $p(y|x)$

- $\boxed{p(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}}$     *distribution? Bernoulli*

- $if \quad y = 1 \quad : \quad p(y|x) = \hat{y}$

- $if \quad y = 0 \quad : \quad p(y|x) = 1 - \hat{y}$

- $\log p(y|x) = \log[\hat{y}^y (1 - \hat{y})^{1-y}] = y \log \hat{y} + (1 - y) \log(1 - \hat{y})$

$$= -L(\hat{y}, y) \quad Max \ Likelihood$$

# Logistic Regression Cost function

- $\log P(\ labels\ in\ trainingset\ ) = \log \prod_{i=1}^{m} P(y^i \mid x^i)$

- $\log P(\cdots) = \sum_{i=1}^{m} \log P(y^{(i)} \mid x^{(i)}) = -\sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)})$

- Cost function $\qquad \underbrace{J(w, b)}_{\textbf{\textit{minimize}}} = \frac{1}{m} \sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)})$

# Neural Networks