# 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

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



→ 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



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

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

- Notation

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

# Binary classification

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

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

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

X.**shape** $= (n_x, m)$

$Y = \begin{bmatrix} y^{(1)}, y^{(2)}, ..., y^{(m)} \end{bmatrix}$

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

Y.**shape** $= (1, m)$

# Logistic Regression

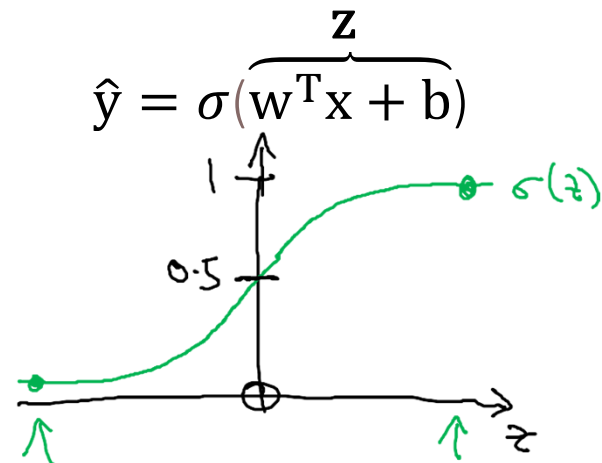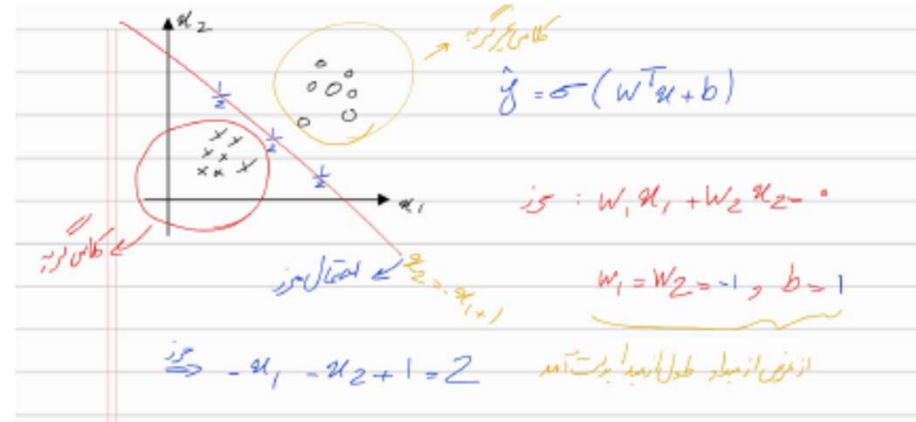• Given x , output  $\hat{y} = P(y=1|x)$      $0 \le \hat{y} \le 1$

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

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

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

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

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

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

# Logistic Regression

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

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

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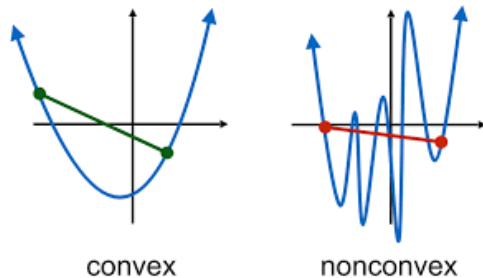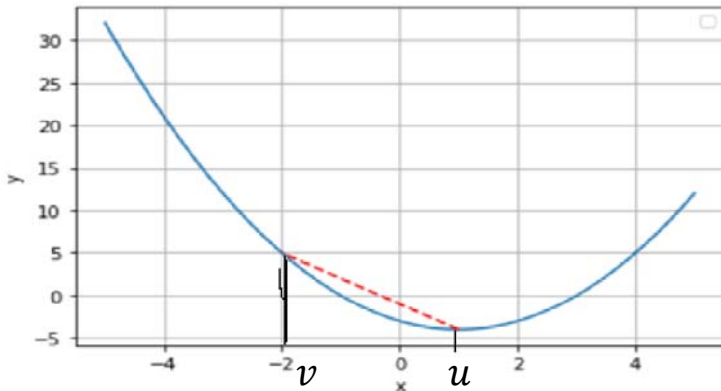
$\}$ b

$\}$ w

# 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



https://mlstory.org/optimization.html

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

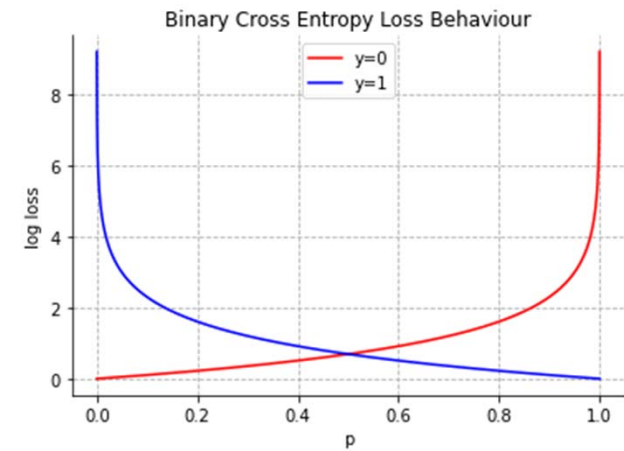**h($\lambda$u +(1- $\lambda$)v) $\leq$ $\lambda$ h(u) + (1- $\lambda$) h(v)**

# Cross Entropy

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

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

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



Binary Cross Entropy Loss Behaviour

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

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

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

# Gradient Descent

# Cost Function

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

$b, w_1$

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

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

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

$w_1 = 2$

# Gradient Descent

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

Minimize $J(b, w_1, \ldots, w_n)$
$b, w_1, \ldots, w_n$

Repeat until convergence:

For j=0,…,n

$b = w_0$

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

$\alpha$ **is learning rate**

**Updating all $w_j$ $Simultaneous$ly**

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

# Gradient Descent

**Correct form**

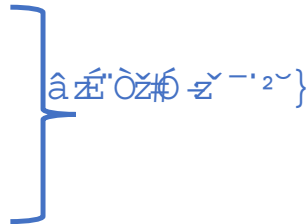$\text{temp0} = b - \alpha \dfrac{dJ(b, w_1)}{db}$

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

$b$ = temp0

$w_1$ = temp1

✔

**Incorrect form**

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

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

❌

# Gradient Descent



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

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

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

$w$ 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.

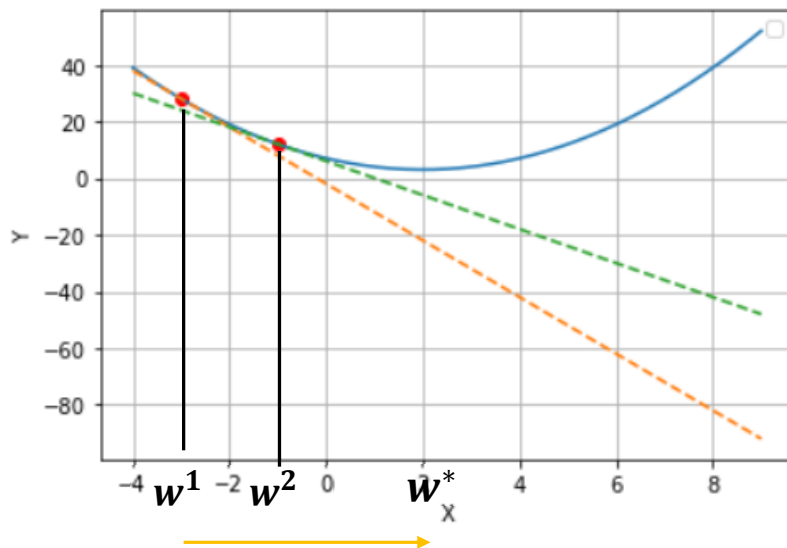# Choosing Learning Rate



$\alpha$ is too large

$\alpha$ is small

# Gradient Descent Weakness

# Gradient Descent



- 1) $\alpha > 0$

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

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

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

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

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

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

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

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

# Gradient Descent

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

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

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

# Logistic regression on $m$ examples

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

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

**Repeat{**

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

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

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

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

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

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

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

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

$\qquad$ **} 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

$w_1 \leftarrow random \quad w_2 \leftarrow random \quad b \leftarrow random$

**Repeat{**

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

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

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

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

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

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

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

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

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

$\quad$ **} until convergence**

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

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

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

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

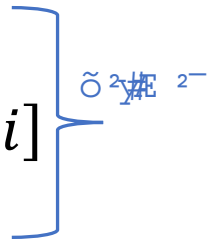# Basics of Neural Network Programming

## Vectorizing Logistic Regression's Gradient Computation
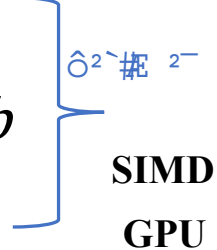
# Vectorizing Logistic Regression

- $z=0;$

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

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

  $\quad z\ +=\ b$


- $z=0;$

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

  **SIMD**

  **GPU**

# Vectorizing Logistic Regression

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

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

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

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

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

# Vectorizing Logistic Regression

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

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

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

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

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

$$1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ \vdots \\ \vdots \\ 1 \end{bmatrix}_{m\times 1}$$

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

$\bullet\ dw=0 \qquad\qquad\qquad\qquad db=0$

$\left\{ \begin{array}{l} \overset{\text{عدد}}{\underset{\text{بردار}}{}} \overset{a^1 - y^1}{} \\ dw\mathrel{+}= \overset{\sim}{x^1}\ \overset{\sim}{dz^1} \\ dw\mathrel{+}=x^2 dz^2 \\ dw\mathrel{+}=x^m dz^m \\ dw/=m \end{array}\right.$

$\left\{ \begin{array}{l} \overset{\text{عدد}}{a1 - y1} \\ db\mathrel{+}= \overset{\sim}{dz1} \\ db\mathrel{+}=dz2 \\ db\mathrel{+}=dz^m \end{array}\right.$

$db/=m$

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

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

$db = \frac{1}{m} dz \times 1$

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

$\boxed{dw = \frac{1}{m} X\, dz^T} = \frac{1}{m}\begin{bmatrix} x^{(1)} x^{(2)} & \dots\dots\dots & x^{(m)} \end{bmatrix} \times \begin{bmatrix} dz^1 \\ dz^2 \\ dz^3 \\ \vdots \\ \vdots \\ dz^m \end{bmatrix}$

# Vectorizing Logistic Regression

- $w_1, w_2, b \leftarrow random$

*For iter in range(1000)*

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

$$A = \sigma(Z)$$

$$dZ = A - Y$$

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

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

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

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

1 epoch

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

# Logistic Regression Cost function

- $\hat{y} = \sigma(w^T x + b)$  $\qquad 0 < \sigma(z) = \frac{1}{1+e^{-z}} < 1$
- $\hat{y} = p(y = 1|x)$
- $if \quad y = 1 \quad : \quad p(y|x) = \hat{y}$  $\quad \Big\}\, p(y|x)$
- $if \quad y = 0 \quad : \quad p(y|x) = 1 - \hat{y}$

- $\boxed{p(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}}$  $\qquad$ *distribution? Bernoulli*
- $if \quad y = 1 \quad : \quad p(y|x) = \hat{y}$
- $if \quad y = 0 \quad : \quad p(y|x) = 1 - \hat{y}$
- $\log p(y|x) = \log[\hat{y}^y (1 - \hat{y})^{1-y}] = y \log \hat{y} + (1 - y) \log(1 - \hat{y})$
$$= -L(\hat{y}, y) \; \textit{Max Likelihood}$$

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

# Neural Networks