

Linear Regression



$m, c?$

Single feature

(x)
hours

(y)
marks

$$y = \theta_1 x + \theta_0$$

weight θ_1 → bias

$\theta_1, \theta_0?$

$$J(\theta)_{\text{loss}} = \frac{1}{m} \sum_{i=1}^m (\hat{y} - y)^2$$

(MSE)

$h_{\theta}(x) = \theta_1 x + \theta_0$

→ Random θ_0, θ_1

→ How good θ_0, θ_1 is?

→ Update θ_0, θ_1

GRADIENT DESCENT



$$\theta = \theta - \eta \left(\frac{\partial J(\theta)}{\partial \theta} \right) \nabla J(\theta)$$

$$\theta_0 = \theta_0 - \eta \cdot \frac{2}{m} \sum_{i=1}^m (\theta_1 x + \theta_0 - y)$$

$$\theta_1 = \theta_1 - \eta \cdot \frac{2}{m} \sum_{i=1}^m (\theta_1 x + \theta_0 - y) x^i$$

Multiple features

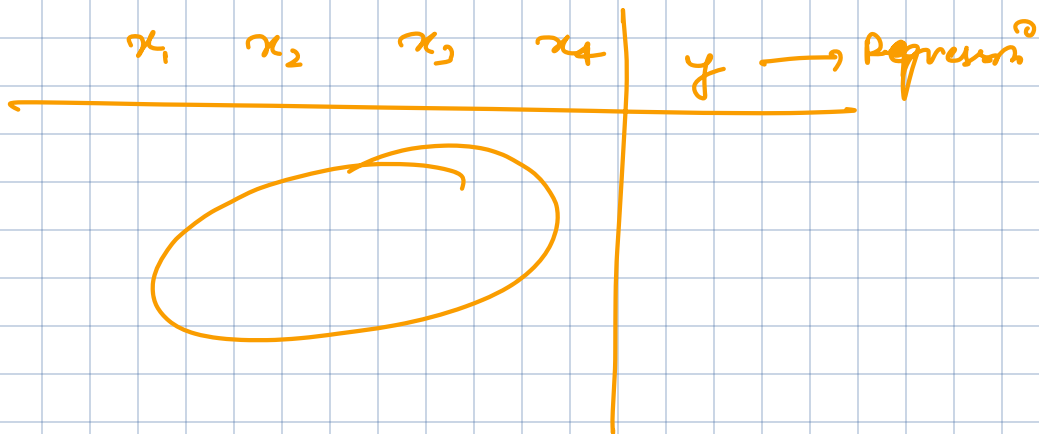
$$\begin{bmatrix} \text{---} x^1 \text{---} \\ \text{---} x^2 \text{---} \end{bmatrix}$$

→

$$\begin{bmatrix} x^1_1 & x^1_2 & x^1_3 & \dots & x^1_n \\ x^2_1 & x^2_2 & x^2_3 & \dots & x^2_n \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\begin{bmatrix} x_1^m & x_2^m & x_3^m & \dots & x_n^m \end{bmatrix}$$



$$h_{\theta}(x) = y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

bio
weights

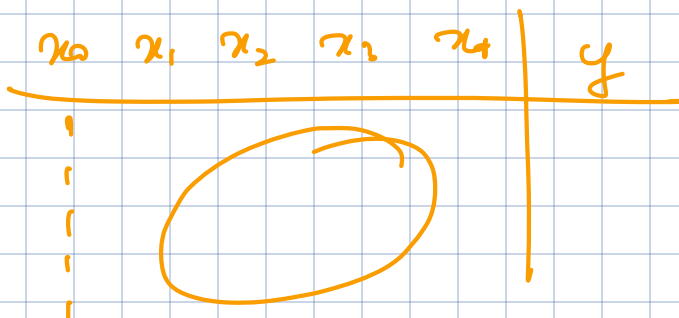
$$h_{\theta}(x) = \theta_0 + \sum_{i=1}^n \theta_i x_i$$

$$h_{\theta}(x) = \theta_0 x_0 + \sum_{i=1}^n \theta_i x_i \quad n=1$$

$$h_{\theta}(x) = \sum_{i=0}^n \theta_i x_i$$

feature 1

$$\left. \begin{array}{l} \theta_0 x_0 + \theta_1 x_1 \\ \underline{1} \\ \theta_0 + \theta_1 x_1 \end{array} \right\} \text{single feature}$$



$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\Theta^T x$$

$$[\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n] \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\hat{y} = h_{\Theta}(x) = \sum_{i=0}^n \theta_i x_i = \Theta^T x$$

Loss f(x):

MSE

$$J(\Theta) = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$$= \frac{1}{m} \sum_{i=1}^m (h_{\Theta}(x^{(i)}) - y^{(i)})^2$$

TASK:

$$\frac{\partial J(\Theta)}{\partial \Theta} ?$$

$$\nabla_{\Theta} J(\Theta)$$

$$\frac{\partial}{\partial \theta_j} J(\Theta) = \frac{\partial}{\partial \theta_j} \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$$= \frac{\partial}{\partial \theta_j} \frac{1}{m} \sum_{i=1}^m (h_{\Theta}(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} (\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_j x_j^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)})^2$$

$$= \frac{1}{m} \sum_{i=1}^m 2 (\underbrace{\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)}}_{\text{Hypothesis}} - y^{(i)}) x_j^{(i)}$$

✓ → Hypothesis
 ✓ → Error / Loss
 → Gradient

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{n} \sum_{i=1}^n 2 (\log(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Final Gradient Update Rule

$$\theta_j = \theta_j - \eta \cdot \frac{1}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$$

$$\hat{y}^{(i)} = \log(x^{(i)}) = \sum_{i=0}^n \theta_i x_i = \theta^T \cdot x$$

$[\theta_0 \ \theta_1 \ \dots \ \theta_n]$ Randomly

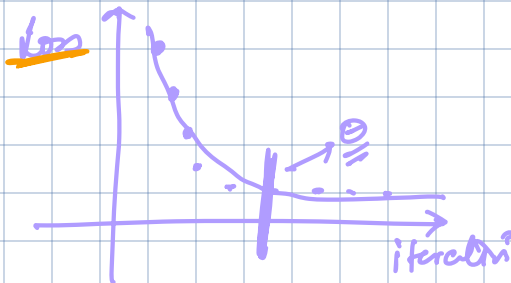
do
{

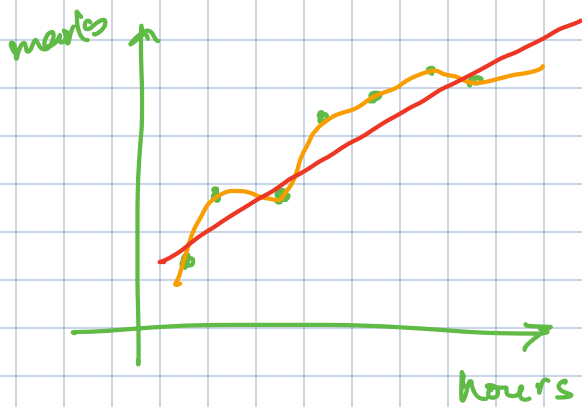
loss fnⁿ (MSE) \rightarrow how good our θ 's are

update θ

} while (converge)

(i) 2500 times





Interpolation

OVERFITTING

- not generalizable for test data points

- has info only about training data points.

ML Algo

y given
Supervised

↓
Unsupervised

Reinforcement

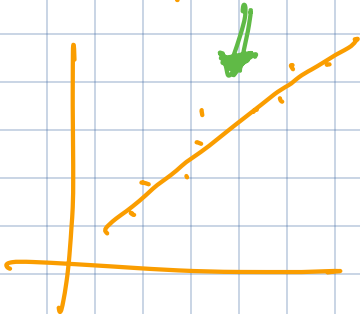
y continuous

Regression

(Linear Regression)

Classification & discrete

(Logistic Regression)



Logistic Regression

↳ Classification Algo

Eg:

→ wt, ht → Dog, Cat

→ Spain or not Spain

→ Animal → Cat
→ Dog

Binary Classification

Training Data:

$$\{x^{(i)}, y^{(i)}\}$$

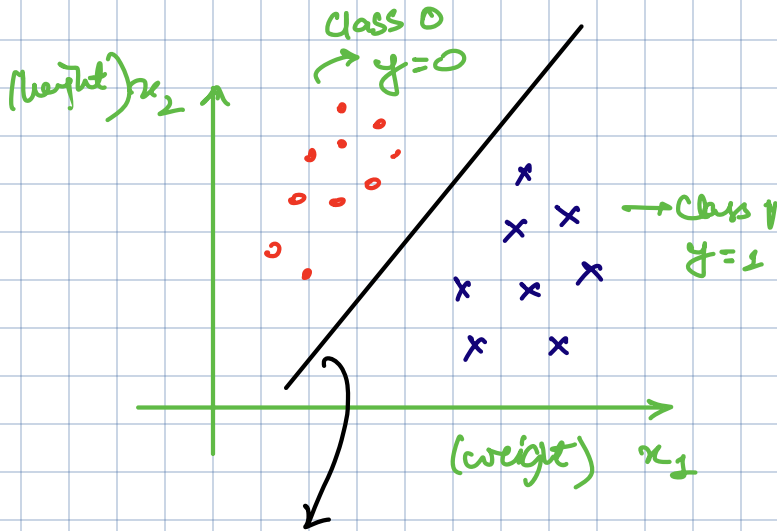
Classification:

y should be a discrete value.

$x \in \mathbb{R}^n \rightarrow x$ has n features
& all features are real nos.

$$y \in \{0, 1\}$$

\swarrow cat \searrow dog



x : cat : 1
 0 : dog : 0

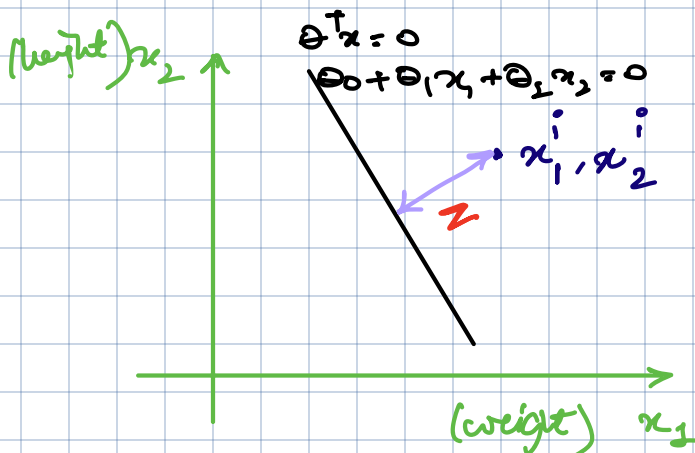
$$y = mx + c$$

$$ax + by + c = 0$$

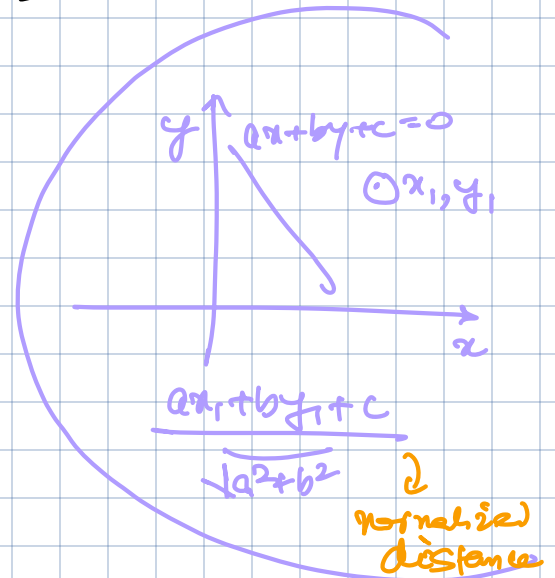
$$\underbrace{\theta_0 + \theta_1 x_1 + \theta_2 x_2}_{\theta^T x} = 0$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \quad \theta^T x = \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \quad x_0 = 1$$

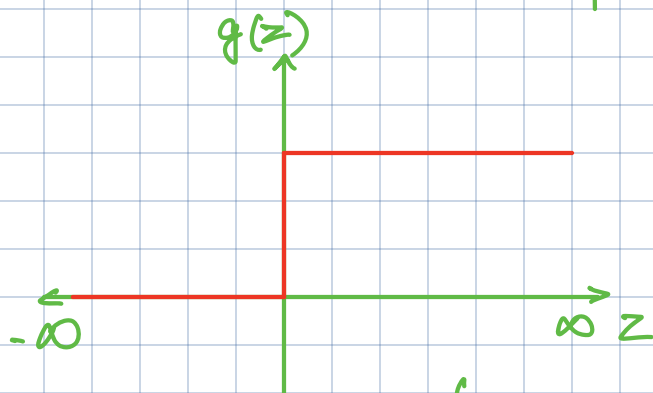
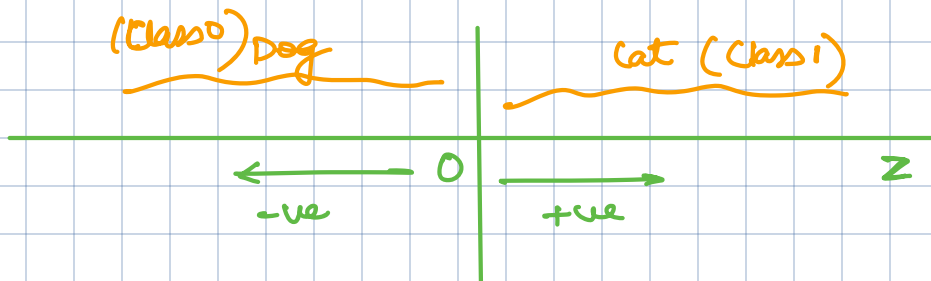
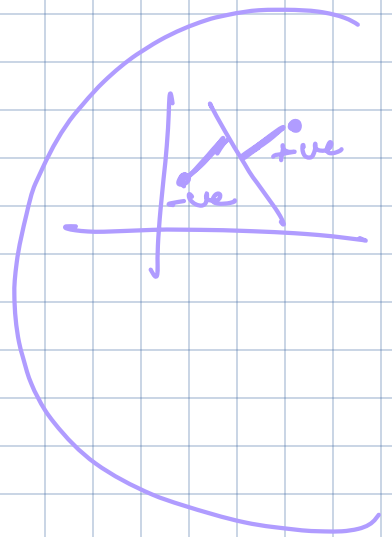
$$= \theta_0 + \theta_1 x_1 + \theta_2 x_2$$



$\theta_0 + \theta_1 x_1^i + \theta_2 x_2^i \rightarrow$ unnormalized distance b/w x^i & line $\theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0$



$$\underbrace{[\theta_0 \ \theta_1 \ \theta_2]}_{\theta^T} \underbrace{\begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}}_{x^{(i)}} = \underline{\underline{\theta^T x^{(i)}}}$$

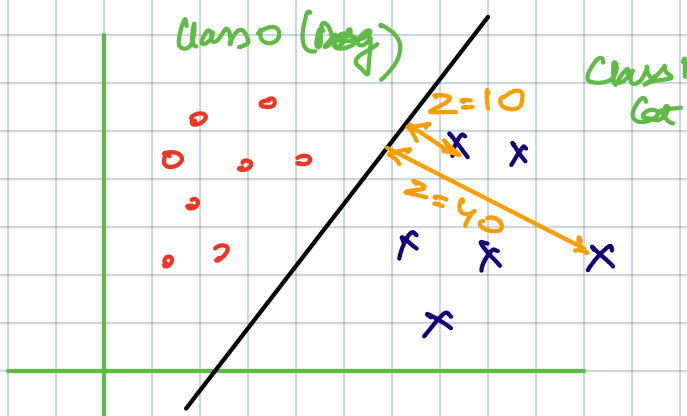


$$g(z) = 1 \quad \text{if } z \geq 0$$

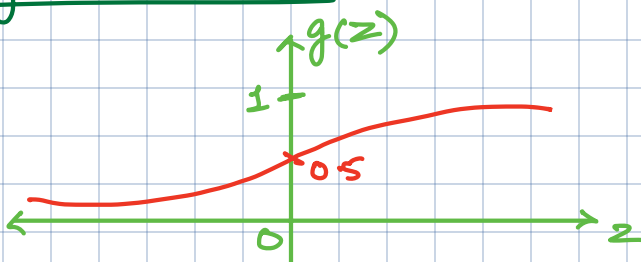
$$g(z) = 0 \quad \text{if } z < 0$$

PROBLEM?

Not able to distinguish b/w points which are close to line & which are far.



Sigmoid function



$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{sigmoid fun}$$

$$z = \infty$$

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

$$z = 0$$

$$g(z) = \frac{1}{1 + \frac{1}{e^0}} = \frac{1}{2} = 0.5$$

$$z = -\infty$$

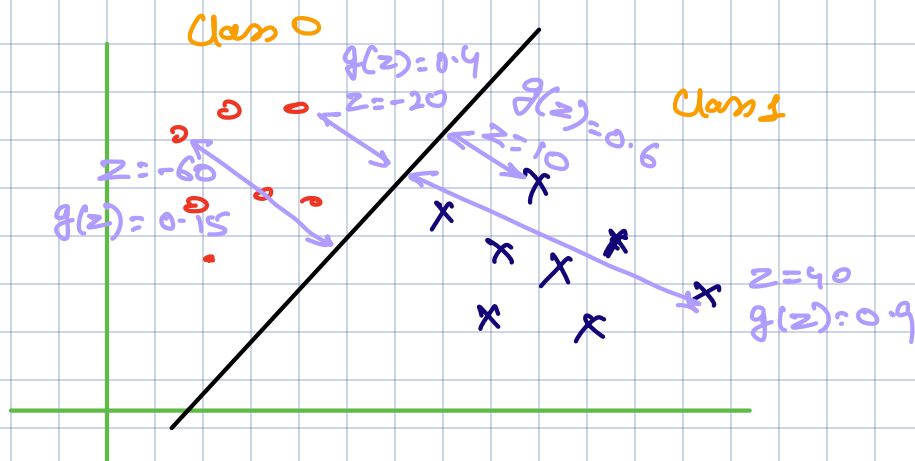
$$g(z) = \frac{1}{1 + \frac{1}{e^{-\infty}}} = \frac{1}{1 + \infty} = 0$$

$$h_0(x) = g(\Theta^T x) = \frac{1}{1 + e^{-\Theta^T x}}$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{where } z = \Theta^T x$$

value b/w 0 & 1

Probability / Confidence with which
you can say the point belongs
to class 1



$$g(z) = 0.4$$

40% sure points
belong to class 1

60% sure point
belongs to class 0

point lies on the line $g(z) = 0.5$

50% sure belongs to Class 0
Class 1

$$h_{\theta}(x) = g(z) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

Probability point \in Class 1

$$\hat{y} = 1 \text{ if } h_{\theta}(x) \geq 0.5$$

$$= 0 \text{ if } h_{\theta}(x) < 0.5$$

Loss funⁿ:

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2$$

Linear hypothesis

$$h_{\theta}(x) = \theta^T x$$

$$= \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

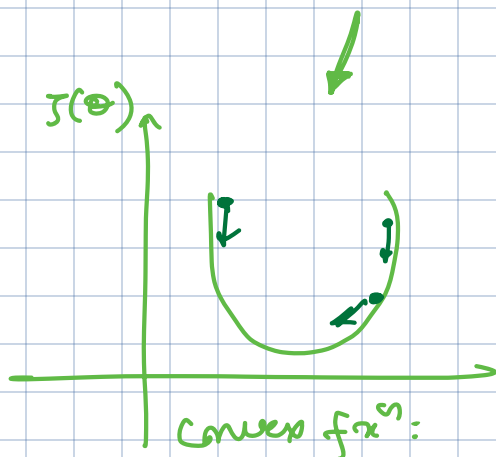
Logistic hypothesis

$$\frac{1}{1 + e^{-\theta^T x}}$$

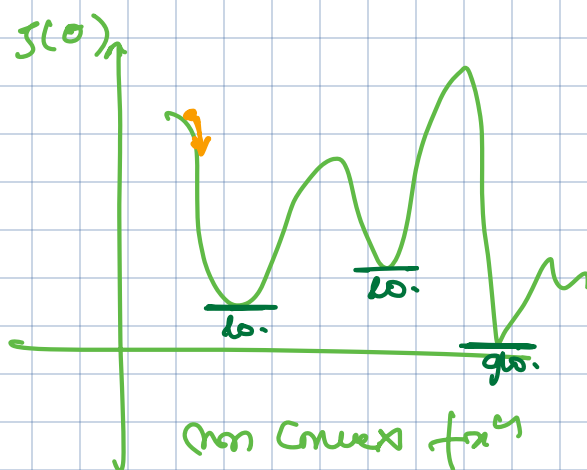
→ Hypothesis ✓

→ Error

→ Gradient



Local minima = Global minima



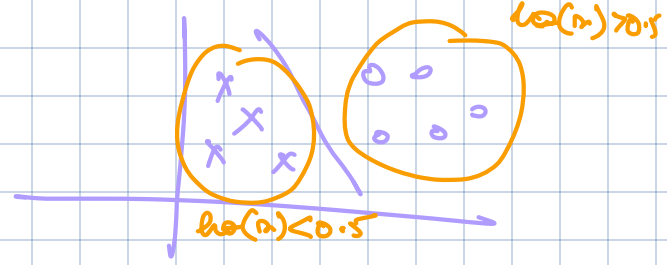
Binary Cross Entropy

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Prob. that point belongs to Class 1

$$P(y=1|x; \theta) = h_{\theta}(x)$$

$$P(y=0|x; \theta) = 1 - h_{\theta}(x)$$



Probability Mass function

$$P(y|x; \theta) = [h_{\theta}(x)]^y [1 - h_{\theta}(x)]^{1-y}$$

$$\begin{array}{l} y=1 \swarrow \\ h_{\theta}(x) \\ y=0 \searrow \\ 1 - h_{\theta}(x) \end{array}$$

0.6 60% Class 1
40% Class 0

Bernoulli Distribution

Likelihood

$$P(y^{(1)} y^{(2)} \dots y^{(m)} | x^{(1)} x^{(2)} \dots x^{(m)}; \theta)$$

$$= \underline{P(y^{(1)} | x^{(1)}; \theta)} \cdot \underline{P(y^{(2)} | x^{(2)}; \theta)} \cdot \dots \cdot \underline{P(y^{(m)} | x^{(m)}; \theta)}$$

$$= \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta)$$

likelihood of the data y .

$$LL(\theta) = \prod_{i=1}^m [h_{\theta}(x)]^y [1 - h_{\theta}(x)]^{1-y}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

loss minimize
LL maximize