

Logistic Regression: classification method.

Prerequisite:

- * binary classification algorithm
- * linear algorithm

- * what is ml?
- * what is classification?
- * what is training & testing dataset?
- * what is x, y in training data
- * what is sigmoid?
- * Gradient descent?

I assume, you have some knowledge on this and starting

training data \rightarrow data used to train the model.

x \rightarrow features - N dimensional array

y \rightarrow label/output/Target value - single dimension array.
class (0 or 1 when binary classification)

Sigmoid \rightarrow will convert real value to probabilities.
(0-1)

need to find the best value for the weights & bias which best fits the line.

Cost function \rightarrow say how close your predicted value to the original target value

Common formula $\Rightarrow \frac{1}{2} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2$

predicted via original
hypothesis value

parameter estimation:

try to minimize the weight & bias

here: for samples labelled "1":

Estimate $(\hat{w})^T x = \hat{\beta}$ such that $\hat{p}(x)$ is as close to 1 as possible

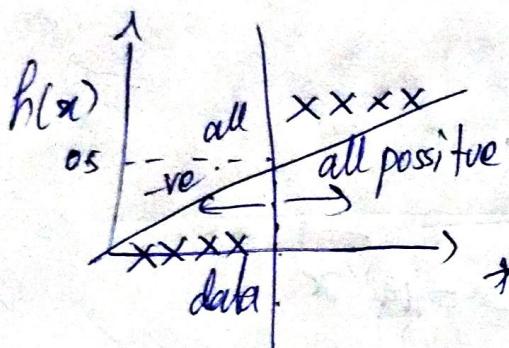
for " " " 0 " :

" " " "

(or)

$1 - \hat{p}(x)$ is as close to 0 as possible

$\hat{p}(x)$ is as close to 0 as possible



* we need something to change the value b/w 0-1

* so we are using sigmoid fun to change the o/p to b/w 0-1

* then $\hat{p}(x)$ is the probability

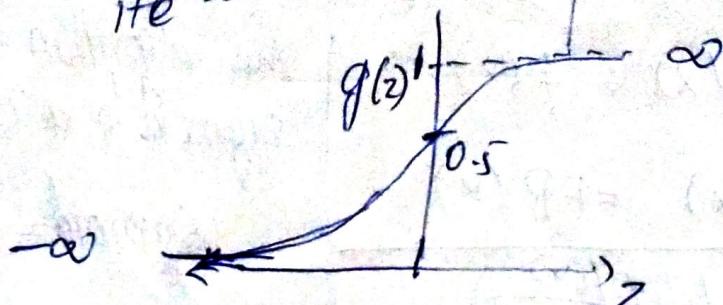
If $\hat{p}(x) \geq 0.5$ then 1
 < 0.5 then 0.

hypothesis representation:

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

| sigmoid (or)
logistic function.

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



so what is the meaning of it... $h_{\theta}(x)$

let's say $0/p \ 1 \Rightarrow$ cancer cell is malignant

$0 =$ cancer cell is not malignant (dangerous)

something
if features $x = [x_{11} \ x_{12} \dots]$
 \downarrow $x_{21} \ x_{22} \dots$

$h_{\theta}(x) = 0.7$ // let's say my hypothesis giving

$0/p$ as 0.7

then I can say that patient having 70% change to being malignant.

$h_{\theta}(x) = P(Y=1/x; \theta)$ = $P(\text{output}/\text{input})$
↳ parameterized by given.

so as per statistics: $P(1) + P(0) = 1$

so $P(Y=0/x; \theta) + P(Y=1/x; \theta) = 1$

then $P(Y=0/x; \theta) = 1 - P(Y=1/x; \theta)$

this inform will help you later.

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq 0.5 \\ 0 & \text{if } h_{\theta}(x) < 0.5 \end{cases}$$

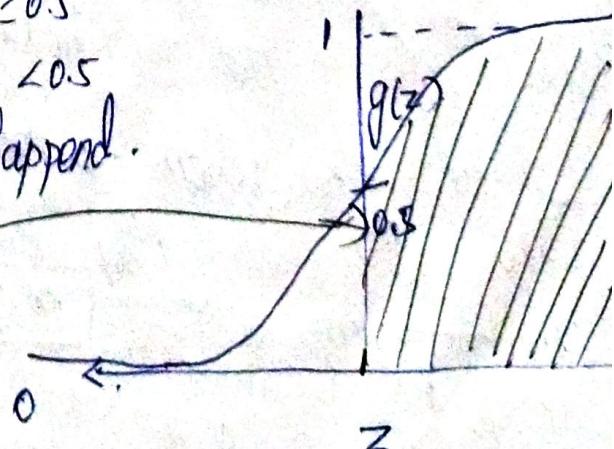
so we can find when this will happen.

so, $g(z) \geq 0.5$ when $z \geq 0$

so, $h_{\theta}(x) = g(\theta^T x) \geq 0.5$

when

$$\theta^T x \geq 0$$



$y = \begin{cases} 1 & \text{if } h_0(x) \geq 0.5 \Rightarrow \text{this will happen } g(\theta^T x) \geq 0.5 \\ 0 & \text{if } h_0(x) < 0.5 \Rightarrow \text{this will happen } g(\theta^T x) < 0.5 \end{cases}$
 when $\theta^T x \geq 0$
 when $\theta^T x < 0$

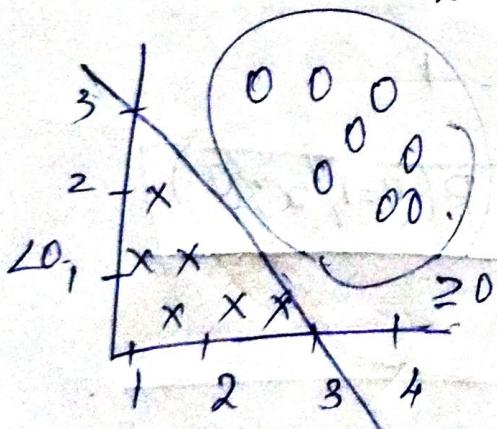
simply $\Rightarrow \theta^T x \geq 0 \quad \text{o/p} = 1$
 $\theta^T x < 0 \quad \text{o/p} = 0$

so what i need to find here:

I need to find the value of $\theta \Rightarrow$ which correctly classifies things.

(or) in other hand
 θ which minimizes the cost function

so lets start from [linear example]



here $h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$

find the value for θ , which correctly classifies this.

so let take $\theta_0 = -3 \quad \theta_1 = 1 \quad \theta_2 = 1$

$h_0(x) = g(-3 + x_1 + x_2)$

so when $-3 + x_1 + x_2 \geq 0$ o/p is 1

other hand $x_1 + x_2 \geq 3$. o/p is 1

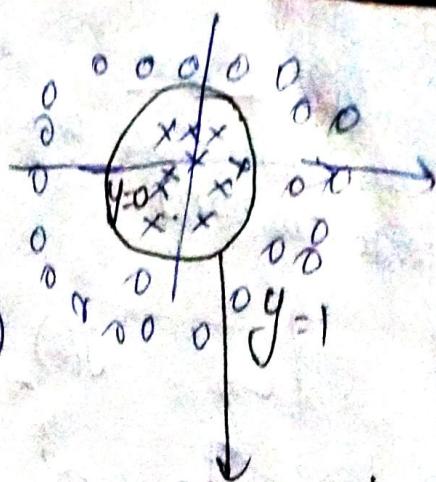
let's plot this

Non-linear example:

need to add non-linearity to the hypothesis so.

$$h(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\text{lets say } \theta_0 = -1, \theta_1 = 0, \theta_2 = 0, \\ \theta_3 = 1, \theta_4 = 1.$$



decision boundary.

$$\frac{-1 + x_1^2 + x_2^2 \geq 0 \text{ then } y=1}{x_1^2 + x_2^2 \geq 1}$$

equation of circle with radius 1 centered to 2010

Cost function: → used to fit the parameters (weights, θ)

Training set: $\{(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)\}$

$$m \# \text{training data} \quad x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_m \end{bmatrix} \quad x_0 = 1 \quad y \in \{0, 1\}$$

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

Cost function:

$$\text{linear regression} = J(\theta) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{1}{2}(h_\theta(x^i) - y^i)^2}_{\text{cost}(h_\theta(x^i), y^i)}$$

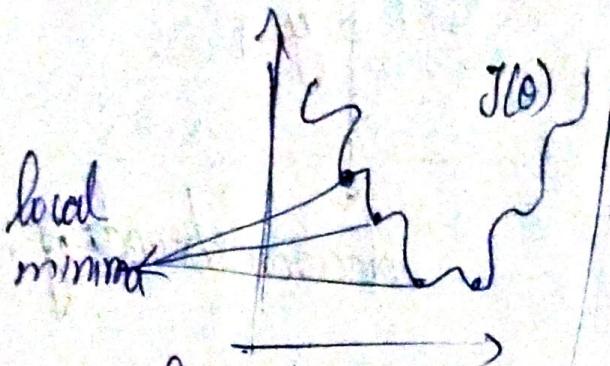
$$\text{so } \text{cost}(h_\theta(x^i), y^i) = \frac{1}{2}(h_\theta(x^i) - y^i)^2$$

diff b/w predicted from original value.

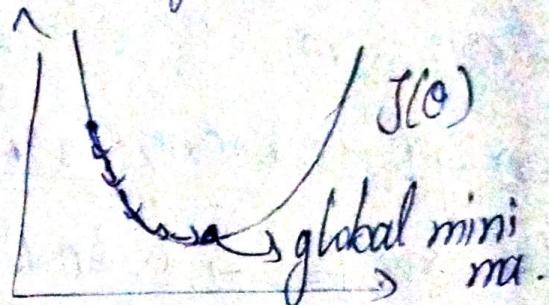
but this $J(\theta)$ will won't work for logistic regression.

because log.Reg is a non-linear function so the cost function may ~~be convex~~ non-convex → we need convex cost function

so, it has non-convex
more than one
local minima



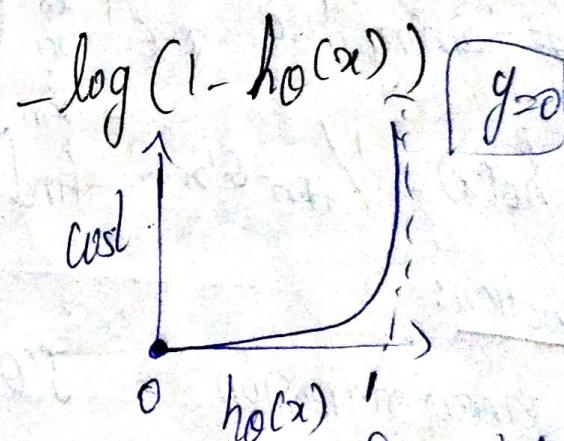
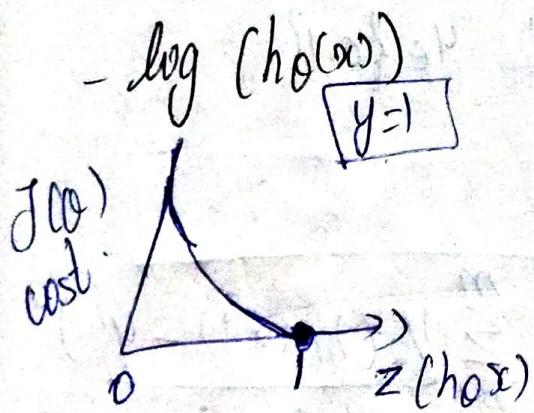
convex
one local minima
(global minima)



so, it has always converged
to global minima.

so we are using some complex function here.

$$\text{cost}(h_0(x), y) = \begin{cases} -\log(h_0(x)) & \text{if } y=1 \\ -\log(1-h_0(x)) & \text{if } y=0. \end{cases}$$



so if $h_0(x) = 1$, then cost is 0

$h_0(x) > 1$ then cost ↑.

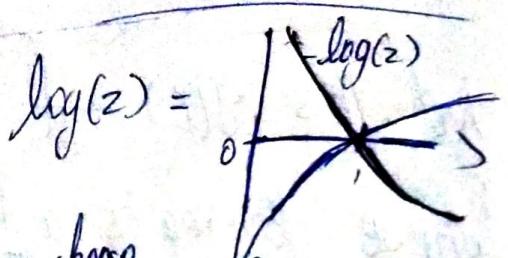
we are finding

$$z = h_0(x) = \mathbf{0}^T x$$

so if $h_0(x) = 0$ then cost ↑
 $h_0(x) < 0$ then cost ↓

why $-\log(z)$ because we need a convex function

we need value to $b/w [0-1]$ so we choose $-\log(z)$



so Finally coming to cost function for log reg & Gradient descent

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^i), y^i)$$

$$\text{Cost}(h_\theta(x^i), y^i) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases} \quad \text{①}$$

we can write it in a single line:

Note $y=0$ always.

$$\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x)) \quad \text{②}$$

① = ② \Rightarrow let's prove. (1-1)

if $y=1$ $\text{cost}(h_\theta(x), y) = -1 \log(h_\theta(x))$ - same as first one

if $y=0$ $\text{cost}(h_\theta(x), y) = -0 \log(1-h_\theta(x))$ - same as second

$$\textcircled{1} = \textcircled{2}$$

$\Rightarrow \textcircled{1}$

$$\text{so } \boxed{J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^i \log(h_\theta(x^i)) + (1-y^i) \log(1-h_\theta(x^i))}$$

why this cost fun?

* this has derived from statistics principle
Maximum likelihood estimation.

which is always used to find the best parameter.

* and it has the "convex" property.

To fit θ : feature scaling help to train faster for log reg

$\min_{\theta} J(\theta)$ and find θ

To make a prediction

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

using gradient decent:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$$J(\theta) = y^i \log(h_{\theta}(x)) + (1-y^i) \log(1-h_{\theta}(x))$$

$\frac{\partial}{\partial \theta_j} J(\theta) =$ let's forgot i and do for one data then add i

$$y^i \log(h_{\theta}(x)) = y^i \log\left(\frac{1}{1+e^{-\theta^T x}}\right) \Rightarrow$$

to derivate this use chain rule

$$\boxed{\frac{d}{dx} \log(x) = \frac{1}{x}} \quad \hookrightarrow f(g(x)) = f'(g(x)) \cdot g'(x)$$

$$\Rightarrow y \left(\frac{1}{1+e^{-\theta^T x}} \right) \frac{\partial}{\partial \theta_j} \left(\frac{1}{1+e^{-\theta^T x}} \right)$$

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

$$\Rightarrow y \cancel{\left(\frac{1}{1+e^{-\theta^T x}} \right)} \frac{\partial}{\partial \theta_j} \cancel{\left(\frac{1}{1+e^{-\theta^T x}} \right)}$$

$$\Rightarrow y \left(\cancel{\frac{1}{g(\theta^T x)}} \right) \frac{\partial}{\partial \theta_j} \left(\cancel{g(\theta^T x)} \right) + (1-y) \frac{1}{1-g(\theta^T x)} \frac{\partial}{\partial \theta_j} \left(1-g(\theta^T x) \right)$$

$$\Rightarrow \frac{\partial}{\partial \theta_j} \left[g(\theta^T x) \right] \left[y \left(\frac{1}{g(\theta^T x)} \right) - (1-y) \left(\frac{1}{1-g(\theta^T x)} \right) \right]$$

$$\Rightarrow x_j \cdot g(\theta^T x) (1-g(\theta^T x))$$

$$\Rightarrow x_j^o \left[y \left(\frac{1}{g(\theta^T x)} \right) - (1-y) \left(\frac{1}{1-g(\theta^T x)} \right) \right]$$

• $\frac{\partial}{\partial \theta_j} \left(y \left(\frac{1}{g(\theta^T x)} \right) \right)$ again chain rule
 derivation for $\frac{\partial}{\partial \theta_j} g(\theta^T x)$: $\frac{\partial}{\partial \theta_j} g(z) = \frac{\partial}{\partial z} g(z) \cdot \frac{\partial}{\partial \theta_j} \theta^T x$

$$\frac{\partial}{\partial \theta_j} g(\theta^T x) \Rightarrow \frac{\partial}{\partial z} g(z) = \frac{\partial}{\partial z} (1/e^{-z})$$

$$\Rightarrow \frac{1}{(1+e^{-z})^2} (e^{-z})$$

$$= \frac{1}{1+e^{-z}} \cdot \left(1 - \frac{1}{1+e^{-z}} \right)$$

$$= \underline{g(z) \cdot (1-g(z))} //$$

so $\frac{\partial}{\partial \theta_j} (g(\theta^T x)) = \cancel{\frac{\partial}{\partial \theta_j} (g(\theta^T x))} \cdot \frac{\partial}{\partial \theta_j} (\theta^T x)$

$$\Rightarrow g(\theta^T x) (1-g(\theta^T x)) \cdot x_j^o$$

so:

$$\Rightarrow g(\theta^T x) (1-g(\theta^T x)) x_j^o \left[y \left(\frac{1}{g(\theta^T x)} \right) - (1-y) \left(\frac{1}{1-g(\theta^T x)} \right) \right]$$

$\cancel{x_j^o} \left[y \cdot \cancel{g(\theta^T x)} (1-g(\theta^T x)) \cdot \frac{1}{g(\theta^T x)} - (1-y) \cdot g(\theta^T x) (1-g(\theta^T x)) \cdot \frac{1}{1-g(\theta^T x)} \right]$

$$\Rightarrow x_j^o \left[\cancel{g(\theta^T x)} - (1-g(\theta^T x)) \right]$$

$$\Rightarrow x_j^o \left[y - \cancel{y g(\theta^T x)} \right] - \left[\cancel{g(\theta^T x)} - \cancel{y g(\theta^T x)} \right]$$

$$\Rightarrow x_j^o \left[y - g(\theta^T x) \right]$$

$$\Rightarrow \frac{1}{m} \sum_{j=1}^m [x_j(y - g(\theta^T x))]$$

comes from $J(\theta) = -\frac{1}{m} \sum_{j=1}^m [y_j^2 - \dots]$

2 pages
ab
- ①

$$\Rightarrow \text{so } x_j(g(\theta^T x) - y)$$

add for multiple values (add i)

$$\Rightarrow \frac{1}{m} \sum_{j=1}^m x_j^i (g(\theta^T x^i) - y^i)$$

$$\Rightarrow \frac{1}{m} \sum_{j=1}^m x_j^i (h_\theta(x^i) - y^i).$$

So

$$\theta_j = \theta_j^0 - \alpha \sum_{j=1}^m (h_\theta(x^i) - y^i) x_j^i$$

for bias

$$\theta_{bias} = \text{bias} - \alpha \sum_{j=1}^m (h_\theta(x^i) - y^i)$$