

# Model and Cost function

Not's :

$m$  : no. of training ex.s

$x$ 's : input var.s / features

$y$ 's : output var.s / ~~features~~ target var. that we are trying to predict

$(x, y)$  : one training ex. / observ<sup>n</sup>

$(x^{(i)}, y^{(i)})$  :  $i^{\text{th}}$  training ex. / observ<sup>n</sup>

$(x^{(i)}, y^{(i)})$ ;  $i = 1, 2, \dots, m$  : training set

$i$  : index into the training set

$X$  : space of I/P val.s

$Y$  : space of O/P val.s

$n$  : no. of features

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

$\theta$ 's : para. of model

#  $h$  maps from  $x$ 's to  $y$ 's

## Classification.

The classification problem is just like the regression problem, except that the values we now want to predict take on a very small no. of discrete val-s.

i.e., we take some input & label it into 1 of ~~our~~ our known classifc's.

ex: Email: spam/not spam

Online Transaction: fraudulent (y/n)

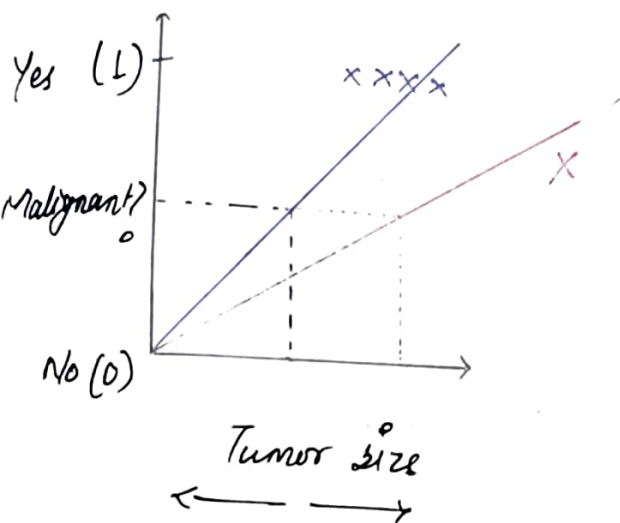
Tumor: Malignant/Benign

# for now, we will focus on binary classification problem.

## Hypothesis Representation

↳ fun<sup>o</sup> that we're going to use to represent our hypoth. when we have a classifc<sup>o</sup> probl.

# Why Linear Regr. won't work at a  
 Classific<sup>n</sup> problem:



i.e., on adding 1 irrelevant  
 data set, the  
 linear fit changes &  
 this doesn't work anymore

Threshold classifier o/p  $h(\theta)$  at 0.5:

If  $h(\theta) \geq 0.5$ , predict "y=1"

If  $h(\theta) < 0.5$ , predict "y=0"

$\therefore$  we k/ that  $y \in \{0, 1\}$  (discrete val.s),

we change our hypth to satisfy  $0 \leq h(\theta) \leq 1$

i.e.,

(classific<sup>n</sup> (binary)) :  $y = 0$  or  $1$

but in Lin. reg.,  $h(\theta)$  can be  $> 1$  or  $< 0$

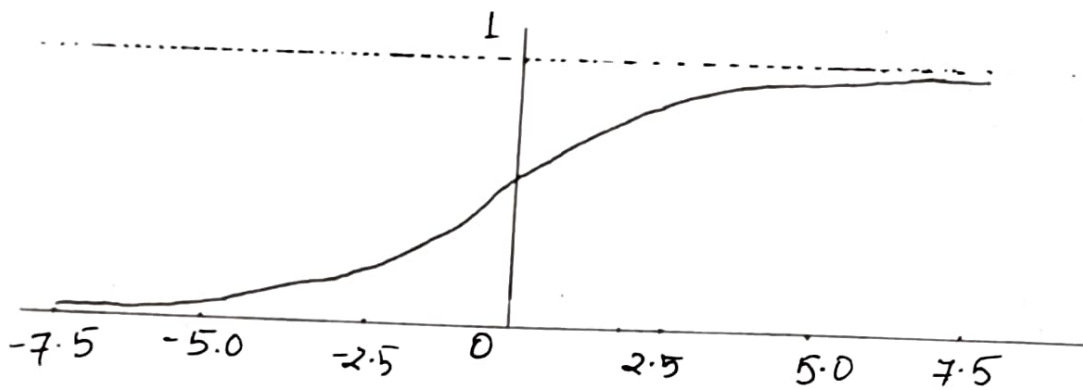
So, Logistic Regr. uses logistic or Sigmoid fun<sup>n</sup>  
 $0 \leq h(\theta) \leq 1$

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

$$z = \theta^T x$$

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

Sigmoid / Logistic function



The sigmoid fun<sup>n</sup> maps any real no. to the  $(0, 1)$  interval.

It asymptotes at 0/1 as it approaches  $-\infty/\infty$ .  
 continuously approaches but doesn't meet

Interpret<sup>n</sup> of Hypoth. o/p :

$h_{\theta}(x)$  = estimated probability that  $y=1$  on  $x$ .  
 on  $1/p$   $x$ .

ex: if  $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$   $x_0$  is always 1  
← feature

$$h_\theta(x) = 0.7$$

tell patient that for some (patient w some tumor size) input  $x$  70% chance of tumor being malignant.

Mathematically,

$$h_\theta(x) = P(y=1 \mid x; \theta)$$

"prob that  $y=1$  gives  $x$  parameterized by  $\theta$ "

# given  $x$ , we calc. the prob. of  $y=1$

# Also,

$$P(y=0 \mid x; \theta) + P(y=1 \mid x; \theta) = 1$$

$$P(y=0 \mid x; \theta) = 1 - P(y=1 \mid x; \theta) = 1 - 0.7$$

$$P(y=0 \mid x; \theta) = 0.3$$

prob that tumor is benign

$$z = 0, \quad e^0 = 1 \quad \Rightarrow \quad g(z) = \frac{1}{2}$$

$$z \rightarrow \infty, \quad e^{-\infty} \rightarrow 0 \quad \Rightarrow \quad g(z) = 1$$

$$z \rightarrow -\infty, \quad e^{\infty} \rightarrow \infty \quad \Rightarrow \quad g(z) = 0$$

## Decision Boundary

↳ is the line that separates the area where  $y=0$  and where  $y=1$ .

It is created by our hypoth fun<sup>n</sup>.

ex: ①

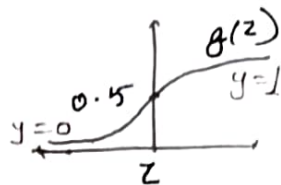
$$h_{\theta}(x) = g(\theta^T x) = P(y=1 | x; \theta)$$

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

$$\text{Now, } z=0, e^0=1$$

$$\Rightarrow g(z) = \frac{1}{2} = 0.5$$

so, predict "y=1" if  $h_{\theta}(x) \geq 0.5$   
( $\rightarrow \theta^T x \geq 0$ )



&

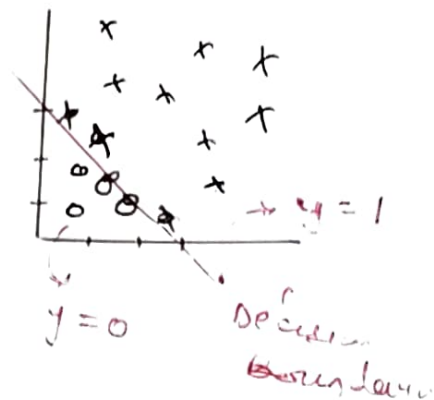
predict "y=0" if  $h_{\theta}(x) < 0.5$   
( $\rightarrow \theta^T x < 0$ )

② let, some training set:

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\text{let, } \theta_0 = -3, \theta_1 = 1, \theta_2 = 1$$

$$\text{so, } \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$



# we're trying to figure out  
where a hypoth. would end up  
predicting  $y=1/0$

Now, predict "y=1" if  $-3 + x_1 + x_2 \geq 0$   
 "y=1" if  $x_1 + x_2 \geq 3$   
 "y=0" if  $x_1 + x_2 < 3$

# Dec. Boundary is a prop. of hypoth  $h_\theta(x)$   
 & its params  $(\theta_0, \theta_1, \theta_2, \dots)$  & not a prop. of data set.

# so even if we remove the training set  $(x_3 \& \theta_3)$  for this hypoth., the dec. boundary would remain the same.

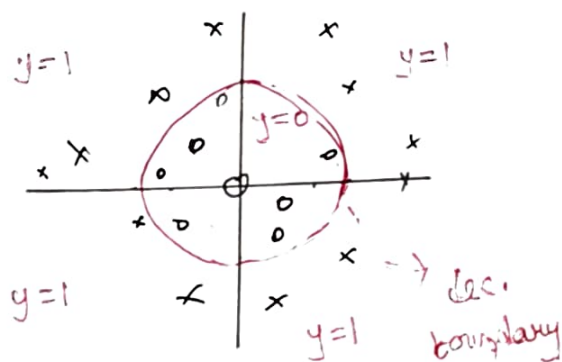
# the training set may be used to fit the params  $\theta$ .

## Non-Linear Decision Boundaries

Higher order polynomial features can have complex dec. boundaries.

i.e., for more complex problems, you can get dec. boundaries w much more complex shapes by adding higher order terms  $(x_1^2, x_1 x_2^3 x_3, \dots)$

ex:



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

predict "y=1" if  $x_1^2 + x_2^2 \geq 1$   
 # circle eq.



## Cost function.

# We cannot use the same cost function that we use for lin. regr. bc the logistic fun<sup>o</sup> will cause the o/p to be wavy, causing many local optima.  
i.e., it will not be a convex fun<sup>n</sup>.

~~Cost~~ Cost fun<sup>n</sup> for Logistic Regression:

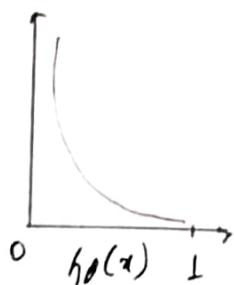
$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

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

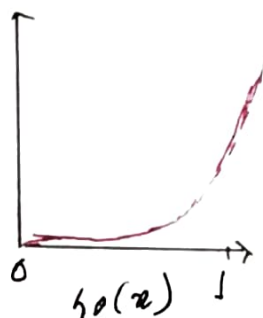
$$\text{Cost}(h_{\theta}(x), y) = -\log(1 - h_{\theta}(x)) \quad \text{if } y = 0.$$

$J(\theta)$  vs  $h_{\theta}(x)$  plots when

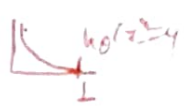
$y = 1$

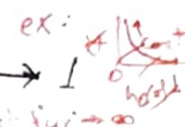


$y = 0$

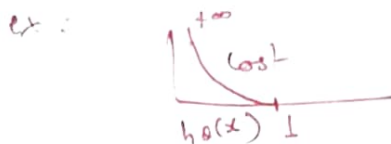




$\text{cost}(h_\theta(x), y) = 0$  if  $h_\theta(x) = y$  (=1 ex: 

$\text{cost}(h_\theta(x), y) \rightarrow \infty$  if  $y=0$  and  $h_\theta(x) \rightarrow 1$  ex: 

$\text{cost}(h_\theta(x), y) \rightarrow +\infty$  if  $y=1$  and  $h_\theta(x) \rightarrow 0$



If our correct answer  $y^*$  is 0, then the cost fun<sup>n</sup> will be 0 if our hypoth. fun also O/P is 0. If our hypoth. approaches 1, then the cost fun<sup>n</sup> will approach  $\infty$ .

If our correct answer  $y^*$  is 1, then the cost fun<sup>n</sup> will be 0 if our hypoth fun also O/P is 1. If the hypoth approaches 0, then the cost fun<sup>n</sup> will approach  $\infty$ .

# Writing cost fun<sup>n</sup> in this way guarantees that  $J(\theta)$  is convex for logistic regr.

(Simplified) Cost function and Grad. Desc.

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

$$\text{where, } \text{cost}(h(x), y) = \begin{cases} -\log(h_{\theta}(x)) & , \text{ if } y=1 \\ -\log(1-h_{\theta}(x)) & , \text{ if } y=0 \end{cases}$$



$$\boxed{\text{cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))}$$

So,

if  $y=1$

$$\begin{aligned} \text{cost}(h_{\theta}(x), y) &= -1 \cdot \log(h_{\theta}(x)) - \cancel{(1-1) \log(1-h_{\theta}(x))}^0 \\ &= -\log(h_{\theta}(x)) \end{aligned}$$

if  $y=0$

$$\begin{aligned} \text{cost}(h_{\theta}(x), y) &= -\cancel{0 \cdot \log(h_{\theta}(x))}^0 - (1-0) \log(1-h_{\theta}(x)) \\ &= -\log(1-h_{\theta}(x)) \end{aligned}$$

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right]$$

Vectorized implementation:

$$h = g(x\theta)$$

$$J(\theta) = \frac{1}{m} \cdot (-y^T \log(h) - (1-y)^T \log(1-h))$$

to fit params  $\theta$ :

$$\min_{\theta} J(\theta) \rightarrow \text{get } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \end{bmatrix}$$

to make a prediction given new  $x$ :

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

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

↑  
prob. that  $y=1$ ,  
given  $x$ ,  
parameterized by  $\theta$ .

Grad. Desc.

general form of grad. desc. is

Repeat

{

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

}

(simultaneously update all  $\theta_j$ )

# here,  $\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$ . So,

Repeat

{

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

}

Vectorized Implement<sup>n</sup>:

$$\theta_j := \theta_j -$$

$$\theta := \theta - \frac{\alpha}{m} X^T (y(X\theta) - y)$$

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

# we observe that the grad. desc. algo. for both lin. regr. & log. regr. is identical just that the def of their hypoth. func<sup>n</sup> has changed ( $\theta^T x$  /  $\frac{1}{1+e^{\theta^T x}}$ )

Multiclass Classification : one-vs-all.

↓

Classification of data for more than 2 categories.

i.e., instead of  $y = \{0, 1\}$ , our def. expands to

$$y = \{0, 1, \dots, n\}$$

Since  $y = \{0, 1, \dots, n\}$ , we divide our problem into  $n+1$  (+1 bc index starts at 0) binary classification problems; in each one, we predict the probab. that 'y' is a member of one of our classes.

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

$$h_0^{(0)}(x) = P(y=0 \mid x; \theta)$$

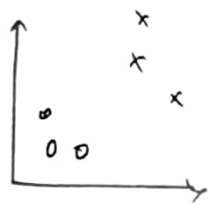
$$h_0^{(1)}(x) = P(y=1 \mid x; \theta)$$

$$h_0^{(n)}(x) = P(y=n \mid x; \theta)$$

$$\text{prediction} = \max_i (h_0^{(i)}(x))$$

We are basically choosing 1 class and lumping all the others into a secondary class. We do this repeatedly, applying binary log. regr. to each case, & then use the hypoth. that returns the highest val. as our prediction.

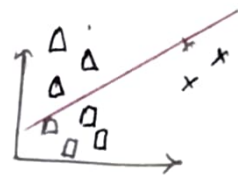
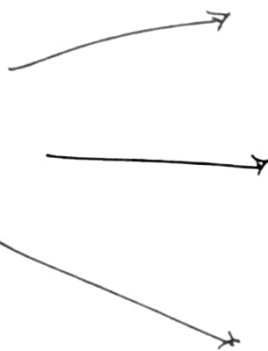
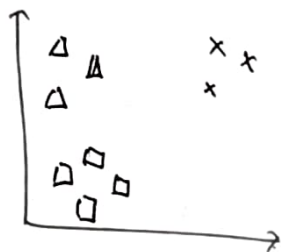
# Binary Classific<sup>n</sup>



# Multi-Class Classific<sup>n</sup>



One-vs-all



Class 1 :  $\Delta$

Class 2 :  $\square$

Class 3 :  $\times$

$$h_{\theta}^{(i)}(x) = P(y = i^o | x; \theta)$$

$$(i^o = 1, 2, 3)$$

One-vs-all

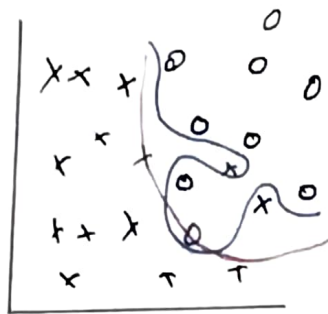
→ Train a log. reg. classifier  $h_{\theta}^{(i)}(x)$  for each class  $i^o$ .  
to predict the prob. that  $y = i^o$

→ On a new i/p  $x$ , to make a predic<sup>n</sup>, run  $h_{\theta}^{(i)}(x)$  for all  $i = 1, 2, 3, \dots$  & pick the class that maximizes  $h_{\theta}(x)$

→ Pick the class  $i$  that maximizes  $\max_i h_{\theta}^{(i)}(x)$

# Regularized Logistic Regression.

optimize objective:



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

optimize objective:

$$J(\theta) = - \left[ \frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

$$+ \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad \leftarrow \text{regularized log. reg.}$$

Grad. Descent.

Repeat

$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

$j = 1, 2, \dots, n$

}

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

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