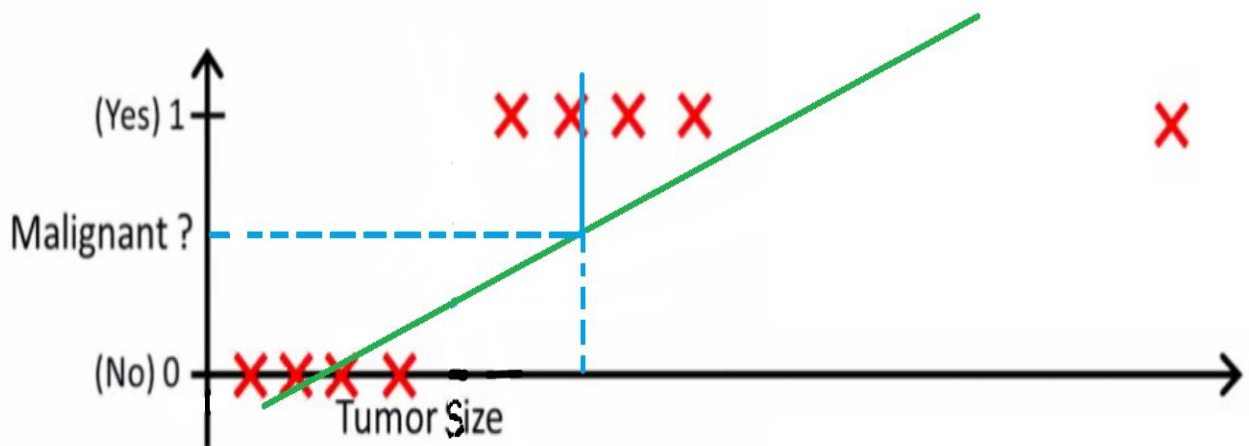
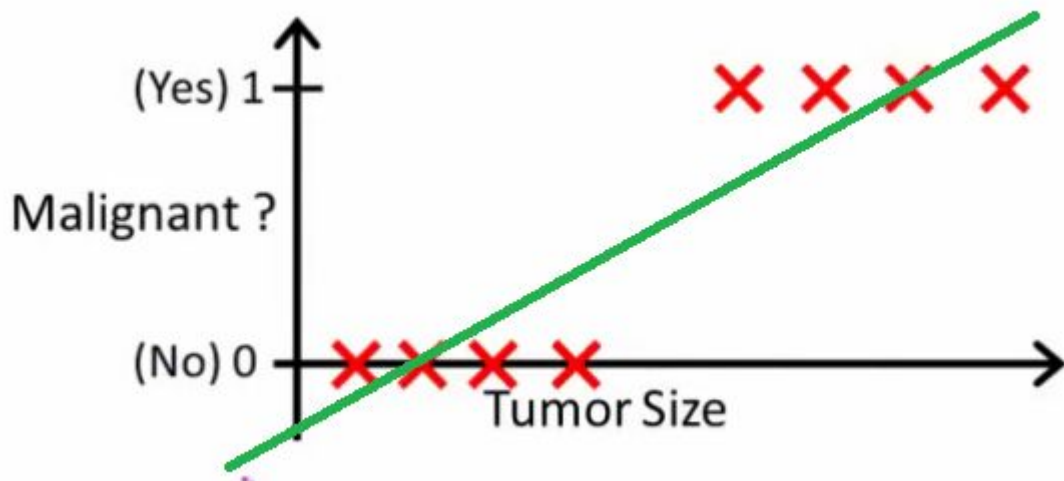
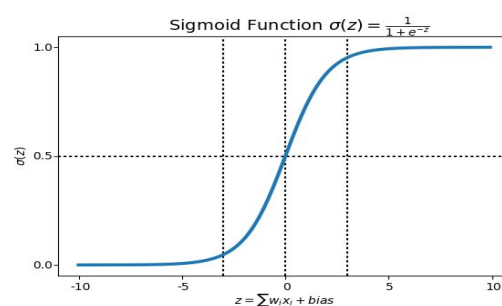


# Logistic Regression



## Sigmoid function



## logistic regression

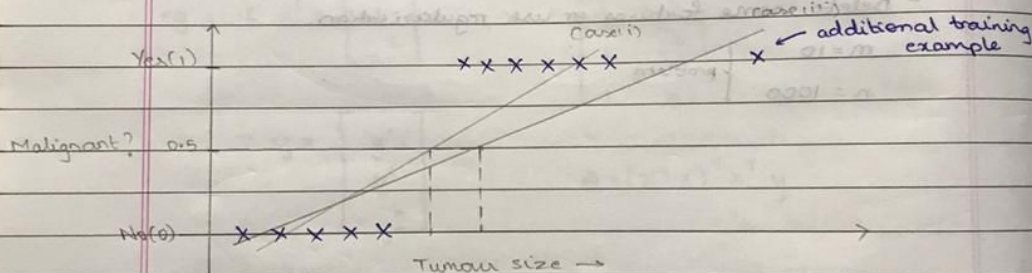
### Classification

- Email: Spam / Not spam?
- Online Transactions: Fraudulent (Yes/No)?
- Tumor: Malignant / Benign?

$y \in \{0, 1\}$       0: "Negative class"  
1: "Positive class"

This is binary classification with two values of  $y$ .

let's take an example



$$h_0(x) = \theta^T x$$

Case(i) Threshold classifies output  $h_0(x)$  at 0.5:

- $h_0(x) \geq 0.5$ , predict " $y = 1$ "
- $h_0(x) \leq 0.5$ , predict " $y = 0$ "

Case(ii) linear regression fails at threshold 0.5 and hence linear regression in classification problem is not suitable. Also  $h_0(x)$  can be  $> 1$  or  $< 0$  but here  $y = 0$  or  $1$ .

So we use logistic regression:  $0 \leq h_0(x) \leq 1$   
classification

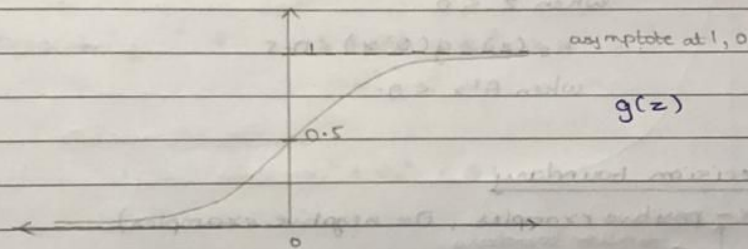
## logistic regression - hypothesis representation

Want  $0 \leq h_{\theta}(x) \leq 1$

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

$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{... sigmoid / logistic function}$$

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



Interpretation of hypothesis output

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

Example: if  $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumour size} \end{bmatrix}$

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

Tell the patient that 70% chance of tumour being malignant

$h_{\theta}(x) = p(y=1|x;\theta) \rightarrow$  "probability that  $y=1$ , given  $x$ , parametrized by  $\theta$ "

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

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



## logistic regression - decision boundary

- Suppose predict "y=1" if  $h_\theta(x) \geq 0.5$

$$\gg g(z) \geq 0.5$$

$$\text{when } z \geq 0$$

$$\therefore h_\theta(x) = g(\theta^T x) \geq 0.5$$

$$\text{when } \theta^T x \geq 0$$

- predict "y=0" if  $h_\theta(x) < 0.5$

$$\gg g(z) < 0.5$$

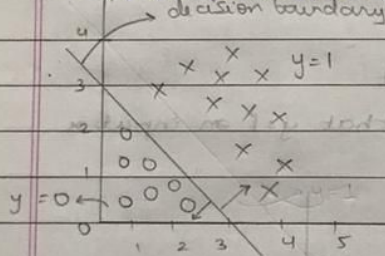
$$\text{when } z < 0$$

$$\therefore h_\theta(x) = g(\theta^T x) < 0.5$$

$$\text{when } \theta^T x < 0$$

### Decision boundary

(x - positive examples, 0 - negative examples)



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

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

predict "y=1" if  $\theta^T x \geq 0$

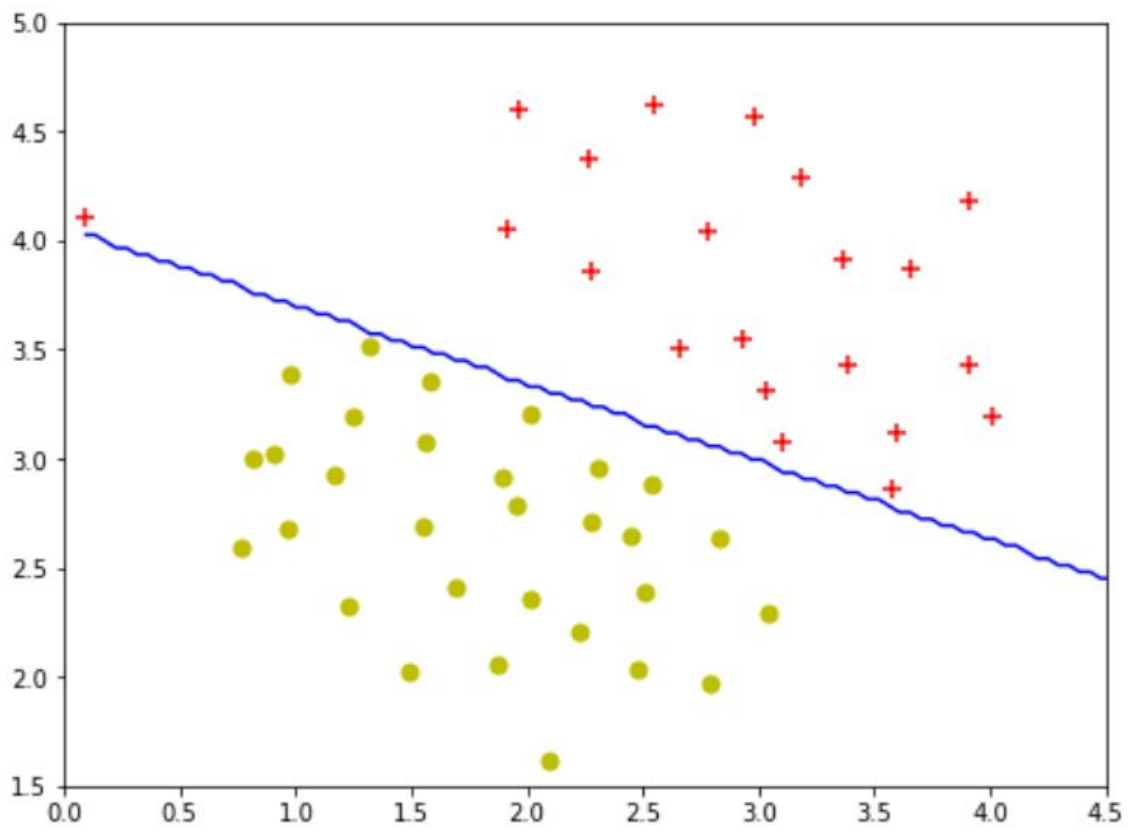
$$-3 + x_1 + x_2 \geq 0$$

$$x_1 + x_2 \geq 0$$

"y=0" if  $x_1 + x_2 < 0$

- \* Decision boundary is a property of the hypothesis,  $\theta$ ,  $\neq$  it is not a function/property of the data set.

# Linear decision boundary



## logistic regression - decision boundary

- Suppose predict "y=1" if  $h_\theta(x) \geq 0.5$

$$\gg g(z) \geq 0.5$$

$$\text{when } z \geq 0$$

$$\therefore h_\theta(x) = g(\theta^T x) \geq 0.5$$

$$\text{when } \theta^T x \geq 0$$

- predict "y=0" if  $h_\theta(x) < 0.5$

$$\gg g(z) < 0.5$$

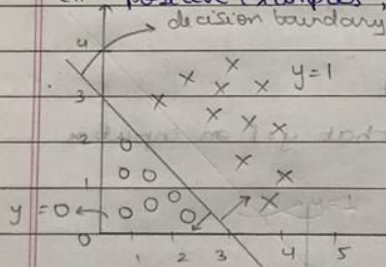
$$\text{when } z < 0$$

$$\therefore h_\theta(x) = g(\theta^T x) < 0.5$$

$$\text{when } \theta^T x < 0$$

### Decision boundary

(x - positive examples, 0 - negative examples)



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

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

predict "y=1" if  $\theta^T x \geq 0$

$$-3 + x_1 + x_2 \geq 0$$

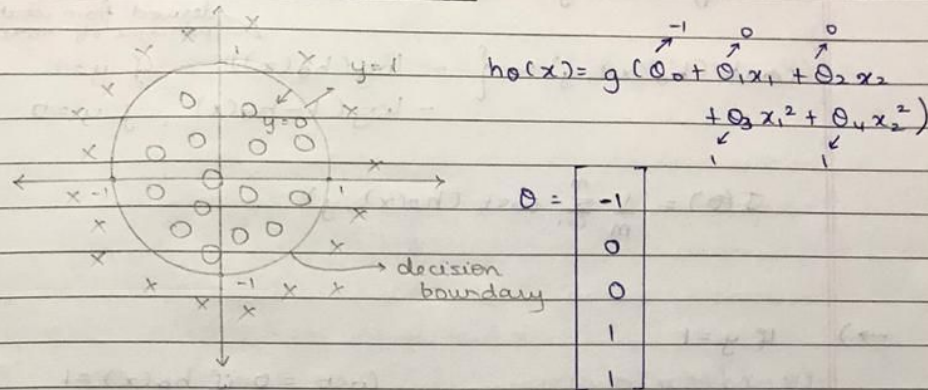
$$x_1 + x_2 \geq 0$$

"y=0" if  $x_1 + x_2 < 0$

- \* Decision boundary is a property of the hypothesis,  $\theta$ ,  $\neq$  it is not a function/property of the data set.



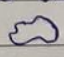
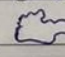
### Non-linear decision boundaries



Predict "y=1" if  $-1 + x_1^2 + x_2^2 \geq 0$   
 $x_1^2 + x_2^2 \geq 1$

"y=0" if  $x_1^2 + x_2^2 < 1$

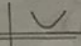
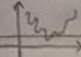
- \* Training set may be used for fitting the values of  $\theta$
- \* Decision boundary does not depend on the training set

With higher order polynomial we could have figures which are more complex (ellipse, , ) as our decision boundaries

### Cost function

We cannot use  $J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$

$\text{cost}(h_0(x), y) = \frac{1}{2} (h_0(x) - y)^2$  because for logistic

regression  $J(\theta)$  is non convex. i.e. the graph of  $J(\theta)$  vs  $\theta$  is not  but instead may be like   
 So we do not global minimum by gradient descent

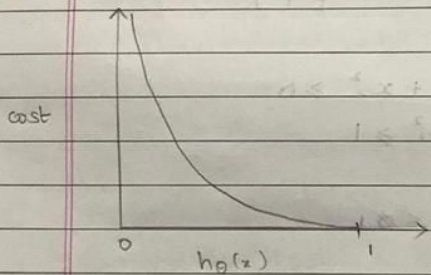
∴ logistic regression - cost function

derived from statistics  
principle of maximum likelihood  
estimation

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

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_0(x), y)$$

→) If  $y=1$



cost = 0 if  $h_0(x) = 1$

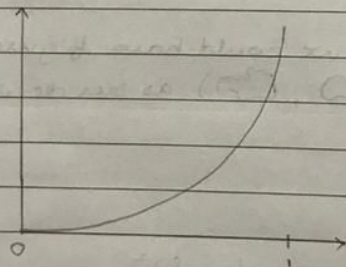
But as  $h_0(x) \rightarrow 0$

cost  $\rightarrow \infty$

$h_0(x) = 0$  but  $y = 1$

so the algorithm is penalized  
with a very high cost

→) If  $y=0$



cost = 0 if  $h_0(x) = 0$

But as  $h_0(x) \rightarrow 1$

cost  $\rightarrow \infty$

above explanation

Simplified cost function

$$\text{cost}(h_0(x), y) = -y(\log(h_0(x))) - (1-y)\log(1-h_0(x))$$

$$\text{if } y=1: \text{cost}(h_0(x), y) = -\log(h_0(x))$$

$$\text{if } y=0: \text{cost}(h_0(x), y) = -\log(1-h_0(x))$$



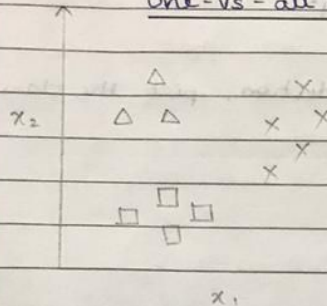
11010 001.4 11010

Repeat

- $y=1$        $y=2$        $y=3$        $y=4$

## Multi class classification

One-vs-all (one-vs-rest):

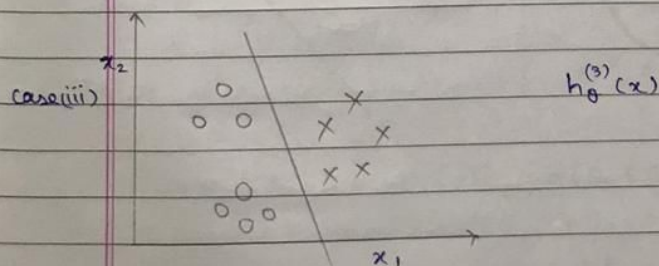
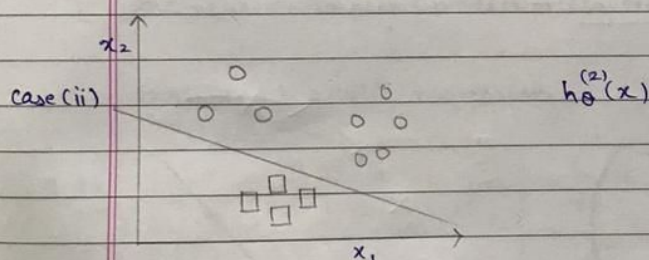
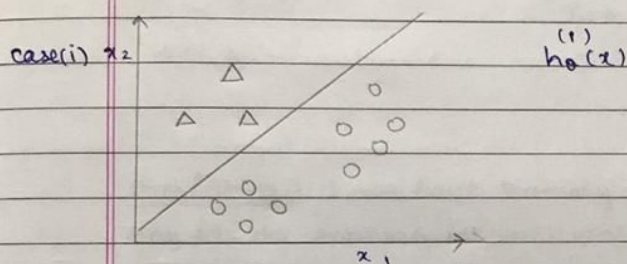


Class 1 =  $\Delta$ ,  $y=1$

Class 2 =  $\square$ ,  $y=2$

Class 3 =  $\times$ ,  $y=3$

we take this and turn it to  
separate binary classification



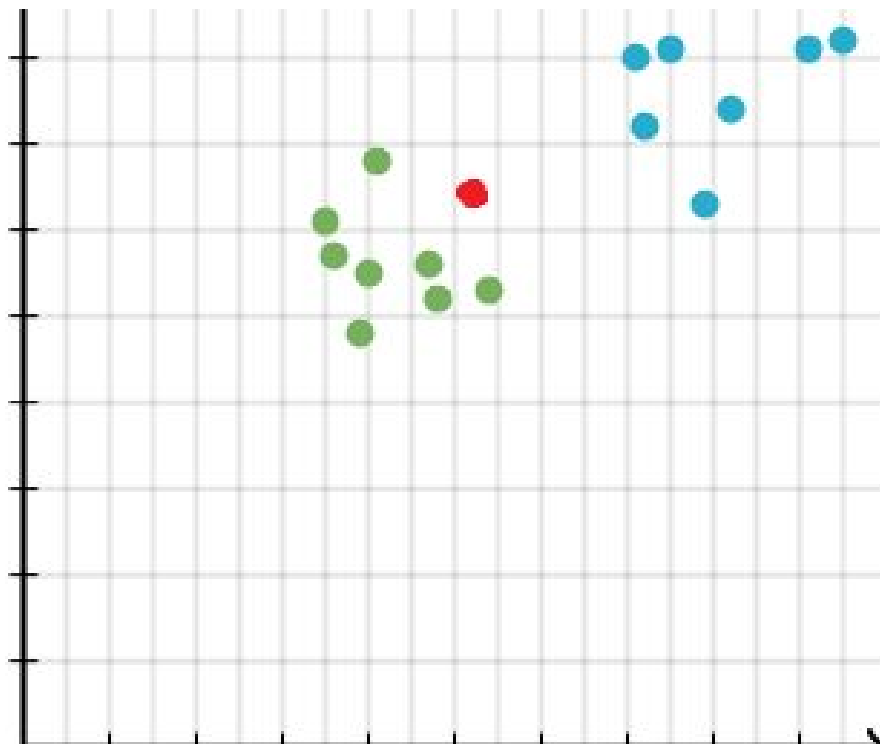
$$h_0^{(i)}(x) = P(y=i | x; \theta) \quad (i=1, 2, 3)$$

# K-nearest Neighbours

The k-nearest neighbors (KNN) algorithm is a simple, easy-to-implement supervised machine learning algorithm that can be used to solve both classification and regression problems.

KNN is a non-parametric, lazy learning algorithm. Its purpose is to use a database in which the data points are separated into several classes to predict the classification of a new sample point.

KNN Algorithm is based on **feature similarity**.





When we say a technique is non-parametric , it means that it does not make any assumptions on the underlying data distribution.

KNN is also a lazy algorithm. What this means is that it does not use the training data points to do any generalization.

## **The KNN Algorithm**

1. Load the data
2. Initialize K to your chosen number of neighbors
3. For each example in the data
  - 3.1 Calculate the distance between the query example and the current example from the data.
  - 3.2 Add the distance and the index of the example to an ordered collection
4. Sort the ordered collection of distances and indices from smallest to largest (in ascending order) by the distances
5. Pick the first K entries from the sorted collection
6. Get the labels of the selected K entries
7. If regression, return the mean of the K labels
8. If classification, return the mode of the K labels

# **A few Applications and Examples of KNN**

1. Should the bank give a loan to an individual? Would an individual default on his or her loan? Is that person closer in characteristics to people who defaulted or did not default on their loans?

2. In political science — classing a potential voter to a “will vote” or “will not vote”, or to “vote Democrat” or “vote Republican”.