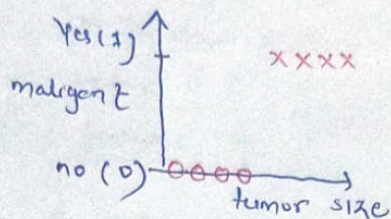


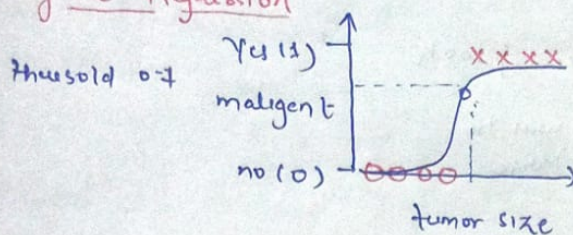


# **Week\_1 Machine learning**

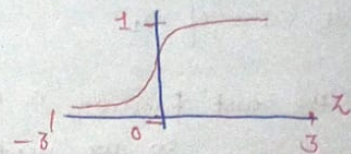
Classification	Ques	Answer
	Is this email spam	"yes" "no"
	Is this transaction fraud	"y" "n"
"Binary classification"		1 0
Two classes or categories		positive class negative class



### Logistic Regression



### Sigmoid Function or Logistic Function



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

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

$$f_{\vec{w}, b}(\vec{x}) = g(\underbrace{\vec{w} \cdot \vec{x} + b}_z) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

"logistic regression"

### Interpretation of logistic regression output

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

"probability that class is 1"

$$f_{\vec{w}, b}(\vec{x}) = 0.7 \quad \Rightarrow 70\% \text{ chances that } y \text{ is } 1$$



when is  $f(\vec{w}, b(\vec{x})) \geq 0.5$

$$g(z) \geq 0.5$$

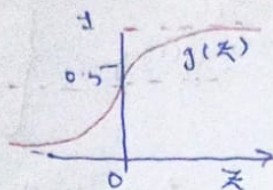
$$z \geq 0$$

$$\vec{w} \cdot \vec{x} + b \geq 0$$

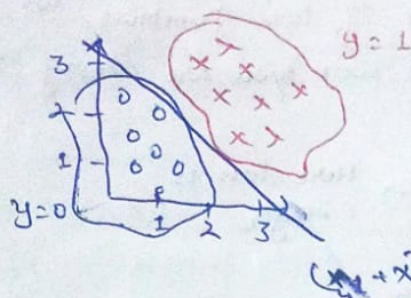
$$\hat{y} = 1$$

$$\vec{w} \cdot \vec{x} + b < 0$$

$$\hat{y} = 0$$



Decision Boundary



$$f(\vec{w}, b(\vec{x})) = g(z) =$$

$$g(w_1 x_1 + w_2 x_2 + b)$$

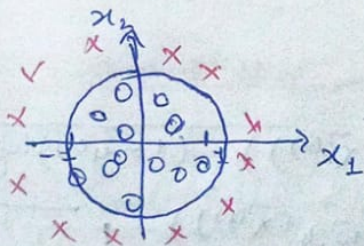
$$\begin{matrix} 1 & 1 & -3 \end{matrix}$$

Decision boundary  $| z = \vec{w} \cdot \vec{x} + b = 0 |$

$$z = x_1 + x_2 - 3 = 0$$

$$x_1 + x_2 = 3$$

Non-linear decision Boundary



$$f(\vec{w}, b(\vec{x})) = g(z) = g(w_1 x_1^2 + w_2 x_2^2 + b)$$

decision boundary  $z = x_1^2 + x_2^2 - 1 = 0$

$$x_1^2 + x_2^2 = 1$$

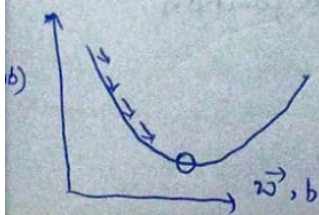
Cost function for Logistic Regression

Squared error cost

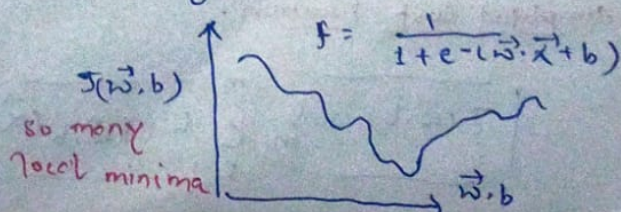
$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f(\vec{w}, b(\vec{x}^{(i)})) - y^{(i)})^2$$

Not good choice

Linear Regression



Logistic Regression



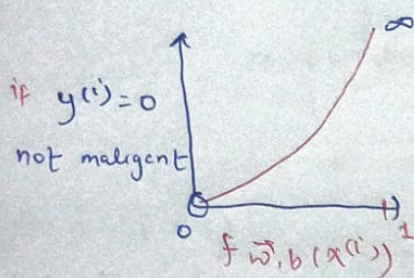
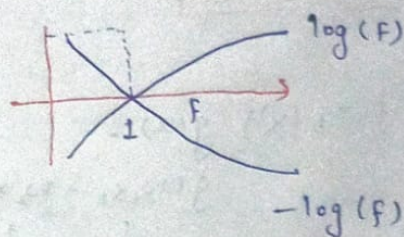
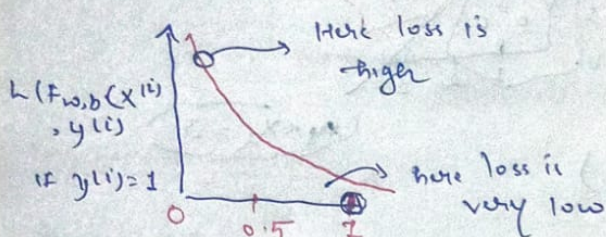


### Logistic loss function

$$L(\vec{w}, b)(\vec{x}^{(i)}, y^{(i)}) = \begin{cases} -\log(F_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - F_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

Loss function tells us how well you are doing in one training sample

After summation of loss function we get cost function which tells us how well you are doing in whole training set.



The further prediction  $F_{\vec{w}, b}(\vec{x}^{(i)})$  is from target  $y^{(i)}$  the higher the loss.

### Simplified loss function.

$$L(F_{\vec{w}, b}(\vec{x}^{(i)}, y^{(i)}) = \frac{-y^{(i)} \log(F_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - F_{\vec{w}, b}(\vec{x}^{(i)}))}{1}$$

if  $y^{(i)} = 1$ :

$$L(F_{\vec{w}, b}(\vec{x}^{(i)}, y^{(i)}) = -1 \log(F_{\vec{w}, b}(\vec{x}^{(i)}))$$

if  $y^{(i)} = 0$ :

$$L(F_{\vec{w}, b}(\vec{x}^{(i)}, y^{(i)}) = -\log(1 - F_{\vec{w}, b}(\vec{x}^{(i)}))$$

### Simplified cost function.

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(F_{\vec{w}, b}(\vec{x}^{(i)}, y^{(i)}))$$

upper function.



## Training logistic regression

find  $\vec{w}, b$  given new  $\vec{x}$ , output  $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

$$p(y=1 | \vec{x}; \vec{w}, b)$$

minimizing the cost function { Gradient descent algorithm }

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))]$$

repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b) \quad \frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

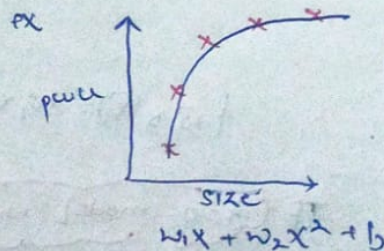
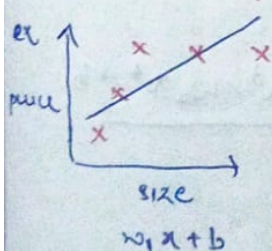
$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b) \quad \frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

this all looks like same as linear regression but here definition of  $f(x)$  is different

Linear regression =  $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

Logistic regression =  $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

## The Problem of Overfitting



Just right

Underfit-

Does not fit the training set well

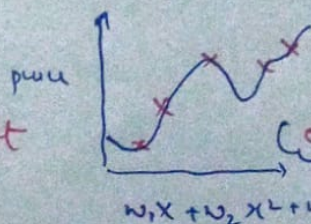
(high Bias)

• fits training set pretty well

generalization: predict well

in new test set

Fits the training set extremely well

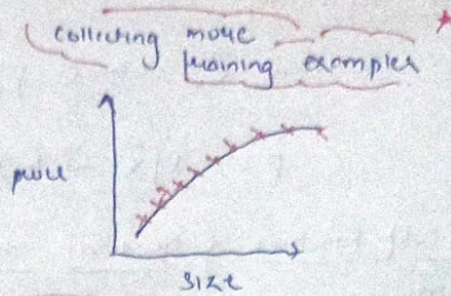
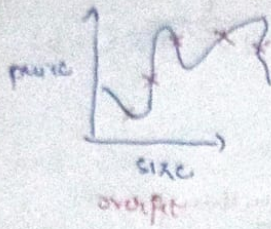


(overfitting)

(high Variance)



## Addressing Overfitting



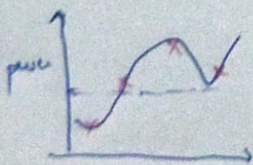
→ Select features to include/exclude → feature selection → Reduce features if you have less data

\* Regularization: Reduce the size of parameters  $w_j$  eliminating features by setting value 0 {or small}  
for ex  $w_3 \approx 0$   
 $w_4 \approx 0$

$$\min_{\vec{w}, b} J(\vec{w}, b) = \min_{\vec{w}, b} \left[ \underbrace{\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2}_{\text{mean squared error}} + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{regularization term}} \right]$$

fit data →      ← keep  $w_j$  small

$\lambda$  balance both goals.



$f_{\vec{w}, b}(\vec{x}) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$

If  $\lambda = 0$  model will overfit

If  $\lambda = 10^{10}$  model will underfit  
 $w_1 \approx w_2 \approx w_3 \approx w_4 = 0$