# WEEK - 3

⇒ We cant use linear regression for classification models.

⇒ LOGISTIC REGRESSION
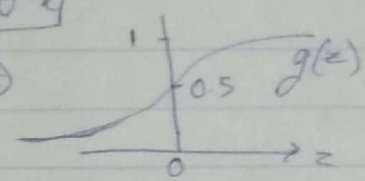$$\left( \text{CLASSIFICATION PROBLEM} \right)$$

⇒ we want $0 \leq h_\theta(x) \leq 1$ (binary classification)

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

$$g(z) = \frac{1}{(1+e^{-z})} \quad \text{(Logistic/Sigmoid Function)}$$

$$\therefore \boxed{z = \theta^T x}$$

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



⇒ For +ve output, $P(y=1 \mid x; \theta)$ 
$\quad$ -ve output, $P(y=0 \mid x; \theta)$ $\Big] += 1$

$\qquad\qquad\qquad$ ↑
$\qquad\qquad$ Probability

⇒ In sigmoid $f^n$ → for $h_\theta(x)$ to be $\geq 0.5$,
$\quad$ (+ve outcome) $\qquad\qquad z \geq 0$
$$\therefore \theta^T x \geq 0$$

⇒ Decision Boundary → The line which separates the +ve area from the -ve area.

$$5 - x_1 > 0$$
$$5 > x_1$$

⇒ <u>Decision Boundary</u> ⇒ line that separates the region where the hypothesis predicts Y equals 1 from the region where the hypothesis predicts that y is equal to 0.

⇒ If we use the cost $f^n$ of linear regression as the cost $f^n$ for logistic regression → we get a non - convex $f^n$.
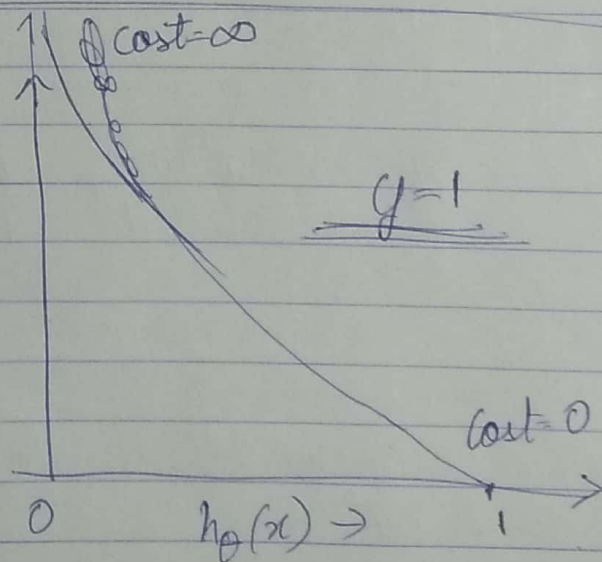
It is very difficult to reach the local minima.

⇒ Cost $f^n$ of Logistic Regression $\boxed{J(\theta) = \frac{1}{m} \sum_{i=1}^{m} Cost}$

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

<u>for $y=1$</u>

⇒ Cost $= 0$, if $y=1$ & if $h_\theta(x) = 1$
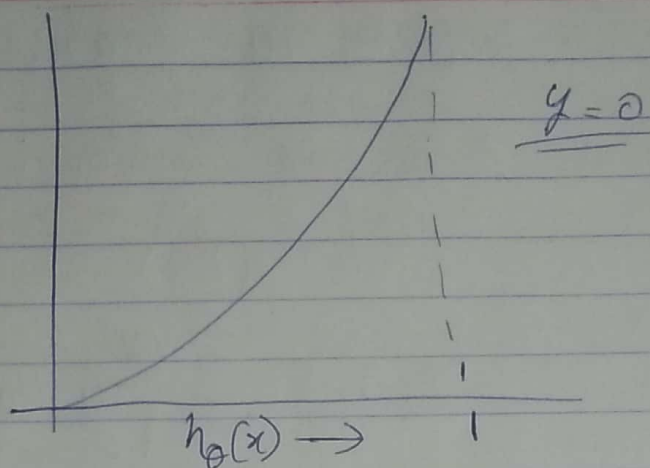
⇒ for $h_\theta(x) = 0$ & $y=1$
→ Cost → ∞

Cost=∞

$y=1$

Cost=0

$0$          $h_\theta(x) →$          $1$

$h_\theta(x) = 0$ → -ve outcome

$\Rightarrow$ for $y = 0$

$\Rightarrow$ if $y = 0$ & $h(x) \to 1$
then cost $\to \infty$

$\Rightarrow$ If $y = 0$ & $h(x) = 0$
Cost $= 0$



$y = 0$

$h_\theta(x) \to \quad 1$

$\Rightarrow$ Logistic Regression Cost $f^n$

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

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

Compact version $\text{cost} = (1/m) * \text{sum}(-y.*\log(\text{hypo})$
$\qquad\qquad\qquad\qquad - (1-y).*\log(1-\text{hypo}))$;

$\Rightarrow$ This cost $f^n$ can be derived from statistics using the principle of maximum likelihood estimation.

$\Rightarrow$ To minimise $J(\theta)$, we use gradient descent as it has the same eg as in linear regression.

Repeat {

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

for all $\theta_j$

}

$\Rightarrow (1/m) * ((\text{hypo} - y)' * X)$;

$\Rightarrow$ we can use feature scaling in logistic Reg.

$\Rightarrow$ Vectorized Implementation of $J(\theta)$ is

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

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

$\Rightarrow$ Vectorized Implementation of Gradient descent

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

$\Rightarrow$ Other ways to minimize $J(\theta)$ $\rightarrow$
1) = Conjugate Gradient
2) = BFGS
3) = L-BFGS

we do not need to pick $\alpha$, to do it themselves

& they are faster than gradient descent
But they are more complex.

$\Rightarrow$ For Multiclass Classification Problems,
we use one-vs-all classifiers.

We use binary logistic Reg. on class to find
its probability.

We are basically choosing one class and then
lumping all the others others into a single
second class.

Cost $f^n$ (Regularized)

$\Rightarrow$ cost = $(1/m)$ * sum($-y .* log(hypo) - (1-y) .* log(1-hypo)$ + term
$(1/2*m)$ * sum(theta(2:end).*2);

$\Rightarrow$ <u>Overfitting</u> $\rightarrow$ If we have too many features, the learned hypothesis may fit the training set very well, but fails to generalize to new examples (predict) aka <u>High Variance</u>

$\Rightarrow$ <u>Underfitting</u> or High Bias $\rightarrow$ not a very good fit to the data. (Using too few features)

$\Rightarrow$ To Address overfitting —
1) = Reduce no. of features
2) = Regularization (keep all features, but reduce the magnitude/value of $\theta_j$.

$\Rightarrow$ <u>Regularized Cost $f^n$</u> $\rightarrow$ keeping all parameters small

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^i) - y^i)^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

regularization parameter (helps balance the two sides)

extra term for regularization

$\Rightarrow$ If the value of $\lambda$ is very high, then algo results in underfitting

$\ast$ (grad) $\theta$ Gradient Descent with Regularization :-

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

for theta $\neq 0$

this no. is slightly less than 1

grad = $\frac{1}{m}$ * ($x'$ * (hypo-y)) + $(1/m)$ * temp

⇒ **Normal eq with Regularization**

Normal eq → $\theta = (X^T X)^{-1} X^T y$.

**if with Regularization**

if $\lambda > 0$,

this makes it invertible

$$\theta = \left( X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix} \right)^{-1} X^T y.$$

$(n+1) \times (n+1)$

$\left(\begin{array}{c}\text{only Identity}\\\text{matrix with}\\\text{0 at first place}\end{array}\right)$

⇒ **Regularized Logistic Regression**

$$J(\theta) = \left[ -\frac{1}{m} \sum_{i=1}^{m} y^i \log(h_\theta(x^i) + (1-y^i) \log 1 - h_\theta(x^i) \right] + \lambda \frac{\sum_{j=1}^{n} \theta_j^2}{2m}$$

⇒