

# Notes

- ① Linear Regression in Matrix Form.
- ② Normal Equation.
- ③ Gradient Descent
- ④ Gradient Descent for Linear Regression

- ① Linear Regression in Matrix Form.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p$$

For  $n$  training Data:

$$\hat{y}^{(1)} \xrightarrow{\text{Sample}} = \hat{\beta}_0^{(1)} + \hat{\beta}_1 x_1^{(1)} + \hat{\beta}_2 x_2^{(1)} + \dots + \hat{\beta}_p x_p^{(1)}$$

$$\hat{y}^{(2)} = \hat{\beta}_0^{(2)} + \hat{\beta}_1 x_1^{(2)} + \hat{\beta}_2 x_2^{(2)} + \dots + \hat{\beta}_p x_p^{(2)}$$

$$\vdots$$

$$\hat{y}^{(n)} = \hat{\beta}_0^{(n)} + \hat{\beta}_1 x_1^{(n)} + \hat{\beta}_2 x_2^{(n)} + \dots + \hat{\beta}_p x_p^{(n)}$$

$$\Rightarrow \begin{bmatrix} \hat{y} \end{bmatrix}_{n \times 1} = \begin{bmatrix} X \end{bmatrix}_{n \times (p+1)} \begin{bmatrix} \hat{\beta} \end{bmatrix}_{(p+1) \times 1}$$

We want to minimize.

$$\|y - \hat{y}\|_2^2 \rightarrow 0$$

$\downarrow$

we estimated  $\hat{y}$  with  $\beta$  and  $X$

i.e.:  $y - X\hat{\beta} = 0$ .

$$\Rightarrow \hat{\beta} = (X^T X)^{-1} X^T y$$

Pros:

Accurate Analytical Solution

Cons:

→ Matrix  $X$  has to be invertible.

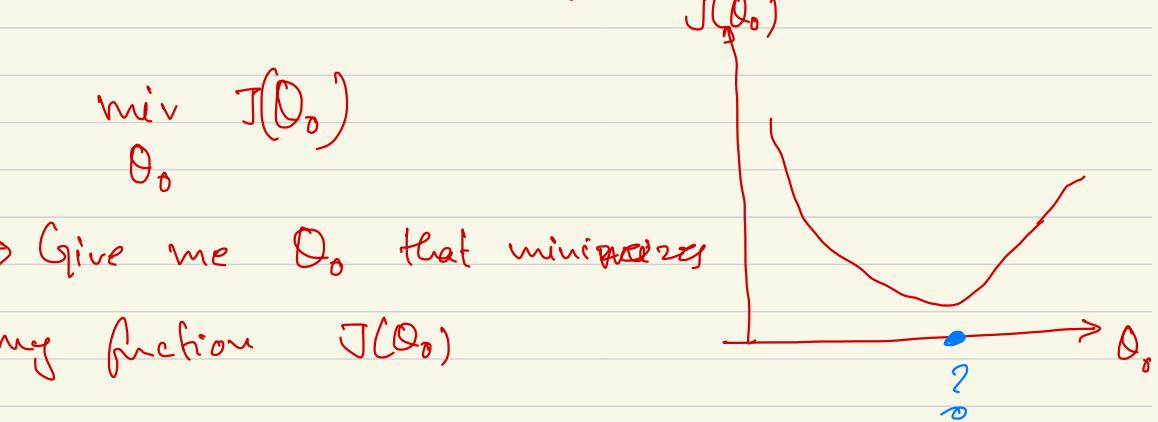
→ As  $n$  and  $p$  increases, computationally expensive

## Gradient Descent.

- Iterative Optimization Algorithm for finding a local minimum of a differentiable function.
- Generally involves two steps:
  - ① Finding the direction of gradient descent & gradient
  - ② Updating arguments based on ①

To understand intuition behind GD we will use a 1-D function, i.e. function which depends on 1 variable.

Let  $J(\theta_0)$  is some function we want to minimize.



- Give me  $\theta_0$  that minimizes my function  $J(\theta_0)$
- We start with some random guess initial value for the parameter we want to find.

→ Usually

- ① Then we find Gradient of function at that value with respect to each
- $$\frac{d J(\theta_0)}{d \theta_0}$$
- parameters.

If there are more than 1 variable  
we take partial derivatives w.r.t each.

Example:  $J(\theta_0, \theta_1)$

$\Rightarrow$  Compute:

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} \text{ and}$$

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$$

② In  $k+1$  th iteration

Update the parameter  $\theta$ .

$$\theta_0^{k+1} := \theta_0^k - \alpha \frac{\partial J(\theta_0)}{\partial \theta_0} \quad \text{learning rate}$$

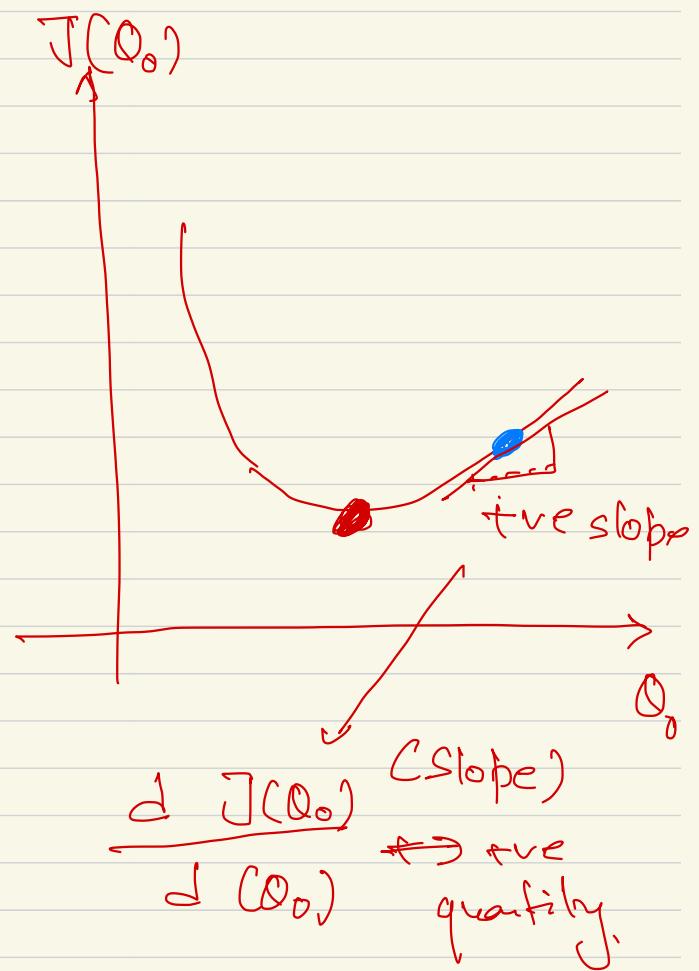
If there are more than 1, we update all acc to:

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1, \dots, \theta_p)}{\partial \theta_j}$$

$$\partial \theta_j$$

$$j = 0, 1, 2, \dots, p$$

What's happening in Step ② Update:



$$\hat{Q}_0^{k+1} := Q_0^k - \alpha \underbrace{\frac{d J(Q_0)}{d (Q_0)}}_{\downarrow Q_0}$$

*Small  
tive*

*Smaller value of  $Q_0$  than in  $k^{\text{th}}$*

If our initial guess were on other side of minimum.

Slope  $\rightarrow$  ve.

$\rightarrow$  We will move to tve direction

When do we stop?

→ One general criteria is:

① Stop until our slope is almost 0  
At minimum slope = 0.

or

② Fixed # of iteration.

or

③ Until the  $J(\theta) - J(\theta^{k+1})$  is a  
some very small value.

Some Interesting points about Gradient Descent and learning parameter

① Gradient Descent can stuck at local minimum

② Often times on first few

iteration the cost decreases

rapidly then starts decreasing

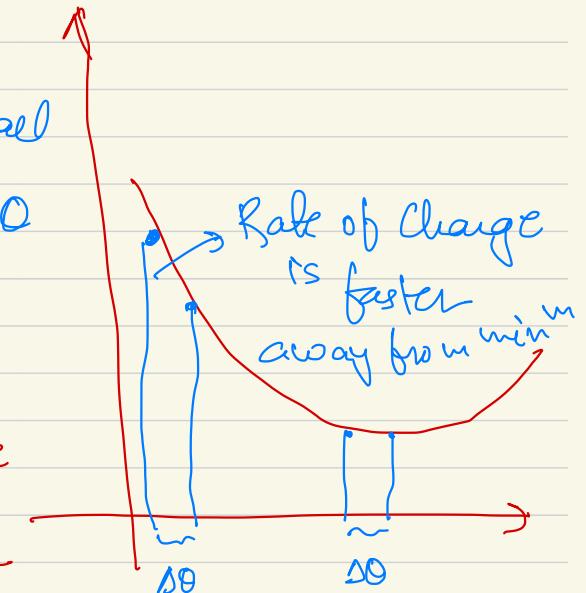
slowly. This is because

the further the away  $\theta$  is from optimal, gradient term is higher



③ For small learning rate

we are traversing through small  
rate of change (Because we change 0  
by a small value).



④ For larger learning rate

the rate of change is faster

but might also bounce in opposite sides of  
optimal value.

### Gradient Descent in Linear Regression

Step 1.

$$J(\beta_0, \beta_1) = \frac{1}{2n} \sum_{i=1}^n (h_{\beta}(x^{(i)}) - y^{(i)})^2 \quad (1)$$

where:  $h_{\beta}(x^{(i)})$  = Approx  $\hat{y}$  of  $y$  using  $i$ th  $X$   
and parameters  $\beta$

$$\hat{y} = \hat{\beta}_0 + \beta_1 x_i$$

Taking Partial Derv. of ① wrt to  $\beta_j$   
 $j = 0, 1, \dots$

$$\Rightarrow \frac{\partial J(\beta_0, \beta_1)}{\partial \beta_j} = \frac{1}{2n} \sum_{i=1}^n (h_\beta(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2n} \frac{\partial}{\partial \beta_j} \sum_{i=1}^n (\beta_0 + \beta_1 x^{(i)} - y^{(i)})^2$$

For  $j = 0$ .

$$\Rightarrow \frac{\partial J(\beta_0, \beta_1)}{\partial \beta_0} = \frac{1}{n} \sum_{i=1}^n (h_\beta(x^{(i)}) - y^{(i)}) \underbrace{\text{Error}}$$

$$\frac{\partial J(\beta_0, \beta_1)}{\partial \beta_1} = \frac{1}{n} \sum_{i=1}^n (h_\beta(x^{(i)}) - y^{(i)}) \underbrace{x^{(i)}}_{\text{Error variable.}}$$

Step 2

Update each  $\beta$  using Update Rule

$$\beta_j^{k+1} = \beta_j^k - \alpha \frac{\partial J(\beta)}{\partial \beta_j}$$

Step 3 Repeat ① and ② until convergence.

GRADIENT DESCENT

FOR

LOGISTIC REGRESSION

# Gradient Descent for logistic Regression:

Recall: for 1 i/p <sup>feature</sup> in logistic Regression

$$P(x) = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}}$$

$\rightarrow P(x) > 0.5 \Rightarrow x \in Y=1 \text{ else, } x \in Y=0$

For P inputs

$$P(x) = \frac{e^{B_0 + B_1 x_1 + B_2 x_2 + \dots + B_p x_p}}{1 + e^{B_0 + B_1 x_1 + \dots + B_p x_p}}$$

$\underbrace{\qquad\qquad\qquad}_{B^T X}$

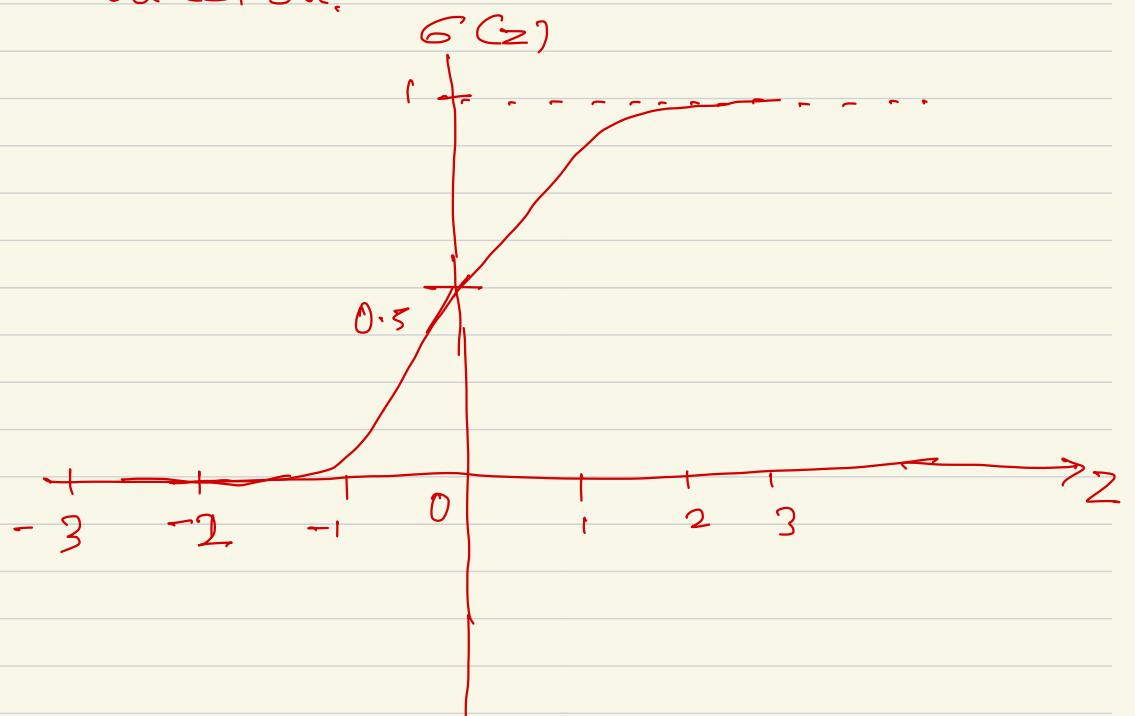
$$\beta \rightarrow R^P$$

$$\text{Then: } P(x) = \frac{e^{B_0 + \beta^T x}}{1 + e^{B_0 + \beta^T x}}$$

$\underbrace{\qquad\qquad\qquad}_{z}$

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

## Sigmoid Function:



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

When we are calculating Prob. in Log Reg:

$$P(Y) = \sigma(z) \text{ where } z = \beta_0 + \beta^T X$$

$$P(Y=1) = \sigma(z) = \sigma(\beta_0 + \beta^T X)$$

$$P(Y=0) = 1 - \sigma(z)$$

For Sigmoid

$$1 - \sigma(z) = \sigma(-z)$$

$$\therefore P(Y=1) = \sigma(z) = \frac{1}{1+e^{-z}}$$

$$P(Y=0) = 1 - \sigma(z) = \sigma(-z) = \frac{1}{1+e^{+z}}$$

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

Decision Boundary:

$$\text{Decision}(x) = \begin{cases} 1 & \text{if } P(Y=1|x) > 0.5 \\ 0 & \text{Otherwise} \end{cases}$$

We have established:

$$\hat{y} = \sigma(z) = \sigma(\beta_0 + \beta^T x)$$

i.e.  $\sigma(z) \geq 0.5$

↳ Note: If  $z > 0$  we decide  $\hat{y} = 1$

Recall: Probabilities are calculated using sigmoid in logistic Regression.

$$L(y, \hat{y})$$

True      Predicted

## Loss Function

Recall Bernoulli

$$f(k, p) = p^k (1-p)^{1-k}$$

We need to find  $\beta_0$  and  $\beta$  such that our loss is minimized

In other words,

We need to find  $\beta_0$  and  $\beta$  such that

$$\hat{y} = \sigma(\beta_0 + \beta^T x) = \sigma(z)$$
 to be closer

to  $y$  which is either 1 or 0

Bernoulli's Distribution

lets Define a rule that maps  $\hat{y}$  to 1 or 0

$$P(y|x) = \hat{y}^y (1-\hat{y})^{1-y} \Rightarrow \text{Likelihood}$$

let's take a look at nature of this expression.

When  $y=1$

$$P(y|x) = \hat{y}^1 (1-\hat{y})^{1-1} \quad \checkmark$$
$$= \hat{y} \quad \checkmark$$

Recall: sigmoid will map  $\hat{y} \rightarrow 1$  when  $\sigma(z) = P(y) > 0.5$

When  $y = 0$

$$P(y|n) = \hat{y}^0 (1-\hat{y})^{1-0}$$

$$P(y|n) = \hat{y}^0 (1-\hat{y})^{1-0} = \frac{\hat{y}}{1-\hat{y}} \approx \hat{y} \quad \text{So writing}$$

Note:

Our sigmoid function maps  $\hat{y}$  to 1 or to 0  
 ~~$\hat{y} \neq 0$~~  based on the Decision Rule.

Recall: Sigmoid will map  $\hat{y}$  to 0,  
when  $\sigma(z) = P(\hat{y}) \leq 0.5$ .

Taking log on both:

$$\log P(y|n) = y \log \hat{y} + (1-y) \log(1-\hat{y})$$

→ This is called log-likelihood that we want to minimize such that our loss  $L(y, \hat{y})$  is minimized.

Now; In terms of loss i.e something we should minimize, if we flip sign of log-likelihood, that's our loss.

$$\text{loss}(y, \hat{y}) = -\log P(y|n)$$

→ We could do this because log is a monotonic monotonic function.

$$L_{CE}(y, \hat{y}) = -\log P(y|n)$$

Cross-Entropy loss

$$L_{CE}(y, \hat{y}) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

Remember:  $\hat{y} = \sigma(z)$

$$= \sigma(\beta_0 + \beta^T x)$$

Example:

Let's see cross-entropy makes sense or not.

Case:  $\rightarrow$  When we made right estimate:

Actual ( $y$ ) = 1      Predicted ( $\hat{y}$ ) = 0.7

$$L_{CE}(y, \hat{y}) = -[1 \cdot \log(0.7) + (1-1) \cancel{\log(1-0.7)}]$$

$$= -\log(0.7)$$

$$= 0.36$$

Case: When we made wrong estimate:

Actual ( $y$ ) = 0      Predicted  $\hat{y} = 0.7$

$$LCE(y, \hat{y}) = -[0 \cdot \log(0.7) + (1-0) \log(1-0.7)]$$

$$= -[\log(0.3)]$$

$$= 1.02$$

Note: loss is smaller when we make right prediction and Vice-versa

//

## GRADIENT DESCENT.

We want to minimize the loss over all training samples.  $x$   
lets say we have  $n$  samples

We can rewrite  $\beta_0 + \beta^T x$  as:

$m$  training  
samples

$$f \left[ \begin{array}{cccccc} 1 & X_{1,1} & X_{1,2} & X_{1,3} & \dots & X_{1,p} \\ 1 & X_{2,1} & X_{2,2} & X_{2,3} & \dots & X_{2,p} \\ \vdots & \vdots & & & & \\ 1 & X_{n,1} & X_{n,2} & X_{n,3} & \dots & X_{n,p} \end{array} \right]$$

~ ~ ~ ~ ~

Dummy.  $X_0^D \Rightarrow$  I/P with  
Dummy 1 column

$$\beta_0 + \beta^T X$$

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}$$

$$\rightarrow \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} \rightarrow \beta^D$$

↓  
Coeff vector with coeff of

Dummy I/P included.

$$\Rightarrow \beta_0 + \beta^T X \Rightarrow (\beta^D)^T X$$

$$\hat{\beta}^D = \underset{\beta^D}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \text{CE}\left(f(x^{(i)}; \beta^D), \hat{y}^{(i)}\right)$$

$\hat{y}$   
is  $i^{\text{th}}$  Training sample:

$\Rightarrow$  Minimize <sup>average</sup> loss over all samples.  
 Cost  $\downarrow$  Function.

All we need is this  $\beta^D \in \mathbb{R}^{P+1}$   
 $\sim$   
 Coefficients.

## Gradient Descent:

- ① Find Gradient
- ② Update parameters

} Recall our approach  
in linear Regn.

$$L(\hat{y}_i, y_i) = L(f_{B^D}(x), \hat{y}_i)$$

General Form

$$\frac{\partial L(\hat{y}_i, y)}{\partial \beta^D} = \left[ \begin{array}{c} \frac{\partial f_{B^D}(x)}{\partial \beta_0} \\ \frac{\partial f_{B^D}(x)}{\partial \beta_1} \\ \frac{\partial f_{B^D}(x)}{\partial \beta_2} \\ \vdots \\ \frac{\partial f_{B^D}(x)}{\partial \beta_p} \end{array} \right]$$

Find Gradient

$$\nabla L(f_{B^D}(x), y) = \left[ \begin{array}{c} \vdots \end{array} \right]$$

Update:

$$\underbrace{\beta^D_{t+1}}_{\downarrow \text{iteration}} = \beta^D_t - \alpha \nabla L(f_{B^D}(x), y)$$

For logistic Regression:

We know

$$L(y, \hat{y}) \rightarrow L_{CE}(\hat{y}, y)$$

$$\therefore L_{CE}(\hat{y}, y) = - \left[ y \log \tilde{\epsilon}(\beta^T x^0) + (1-y) \log (1 - \tilde{\epsilon}(\beta^T x^0)) \right]$$

$$\Rightarrow \frac{\partial L_{CE}(\hat{y}, y)}{\partial \beta_j} = \left[ \underbrace{\tilde{\epsilon}(\beta^T x^0) - y}_{\hat{y}} \right] x_j \\ = [\hat{y} - y] x_j$$

This means:  $j = 0, 1, 2, \dots, p$

Derivative of cross entropy loss w.r.t  
j<sup>th</sup> coeff.  $\beta_j$  is equal to product of  
j<sup>th</sup> input  $x_j$  and difference between predicted  
( $\hat{y}$ ) and actual  $y$

Once, we know this Gradient

⇒ We can calculate this for all training samples and calculate average

$$\frac{\partial}{\partial \beta_j} \underbrace{\frac{1}{n} \sum_{i=1}^n L_{CE}(f_{\theta}(x^{(i)}), y^{(i)})}_{J(\theta)}$$

$\frac{\partial}{\partial \beta_j}$

$$= \frac{1}{n} \sum \left\{ \frac{\partial L_{CE}(f_{\theta}(x^{(i)}), y^{(i)})}{\partial \beta_j} \right\}$$

→ We know how to calc. this.

② We can update parameters using Gradient.

$$\underbrace{\frac{\partial L_{CE}(y, \hat{y})}{\partial \beta_j}}_{=} = \frac{1}{n} \sum (y^{(i)} - \hat{y}^{(i)}) x_j^{(i)}$$

$\frac{\partial}{\partial \beta_j}$

$i = 1, 2, \dots, n \rightarrow$  # of Training example,

$j = 0, 1, 2, \dots, p \rightarrow$  I/P variable

Note:

When  $\hat{j} = 0$ ,  $X_0 = 1$

Once we have Gradient:

$$\beta_j^{t+1} = \beta_j^t - \alpha \frac{\partial L_C(y, \hat{y})}{\partial \beta_j}$$

$\approx$

$$= \beta_j^t - \alpha \frac{1}{n} \sum (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$$

$\leftarrow$  Random Init.       $\leftarrow$  # of Training.       $\leftarrow$

Note we do

$$\hat{y} - y \text{ not}$$

$$y - \hat{y}$$

We are more interested on the sign of  $(\hat{y}^{(i)} - y^{(i)})$  than the actual value. Sign give us the direction of the gradient. That is what we are interested in so that we can update our parameters  $\beta$  to get to an optimal.

MSE estimation is a special case of MLE.

→ Recall in Linear Regression, we used Residual Sum of Squares (RSS) mean, i.e Mean Squared Error to estimate parameter  $\beta$  of linear Regression.

→ But we used Max. Likelihood Estimator (MLE) to estimate parameter  $\beta$  in Logistic, Poisson's and other Regressions.

→ Our goal in this section is to show MSE is actually a special case of MLE.

To show that ~~is~~ we will first briefly discuss about MLE.

② Take an Example of Gaussian Dist<sup>n</sup> to understand MLE.

## Maximum likelihood Estimator.

→ Given some Data samples  $(X_i)$ ,  $i = 1, \dots, n$ , MLE  
allows us to Estimate the parameters of distribution  
that maximizes the chances of getting those Data  
samples.

Given:  $x_1, x_2, x_3, \dots, x_n$ , The likelihood  
function is defined as:

$$L(\theta | x_1, x_2, x_3, \dots, x_n)$$

Joint probability of observed Data as a  
function of the parameters  $\theta$  of the chosen statistical model.

In MLE, we want to estimate parameters  $\theta$  such that  
the likelihood function is maximized.

$$\text{i.e } \hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta | X)$$

To understand Concept of LE and MLE, lets  
take an Example of Statistical Model i.e  
Gaussian.

## Example:

Consider we are measuring size of Tumor ( $X$ ) in patients.



Now, if we want to model this data is coming from a Gaussian, we need to estimate two parameters

- Mean
- $\Sigma^2$  (Variance)

$$\mathcal{N}$$

$$\theta : \{\mu, \sigma^2\}$$

→ Well this is easy right?

Our intuition says, just calculate mean and variance of the sample (Training) Data.

$$\text{i.e. } \hat{\mu} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$$

$$\sigma \approx x - \hat{\mu}$$

But ① is our intuition correct?

② Do maximum likelihood estimators give us some estimation of  $\mu$  and  $\sigma$ ?

Gaussian Dist<sup>n</sup>:

$$f(n) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{n-\mu}{\sigma}\right)^2}$$

Prob Density function.

Let's see how MLE estimates  $\mu$ .

Sample

Let's say only one ~~observed~~ is available.  $x_1 = 2$

Assume we have estimated  $\sigma^2 = 2$ .

$$\sigma = \sqrt{2}$$

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$



$P(X)$

Fig: Here

lets say 2 data points are available:  
 $x = 2, 6$

$$\text{Then: } L(\mu | x=2, x=6) = P(x=2), P(x=6)$$



Because we say samples are independently drawn.

$$= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \left(\frac{x-u}{\sigma}\right)^2} \cdot \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \left(\frac{y-u}{\sigma}\right)^2}$$

Fig 2 here.

Similarly: If  $n$  samples are given:

$$\mathcal{L}(u, \sigma | x = x_1, x_2, \dots, x_n) \\ = \prod_{i=1}^n p(x_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \left(\frac{x_i-u}{\sigma}\right)^2}$$

$$= \prod_{i=1}^n \text{Prob. density function of the model.}$$

If you want to prove that for Gaussian, optimal  $u$  estimator from sample Data is mean (average) of the sample data.

We could take partial derivative of  $\mathcal{L}(u, \sigma)$

w.r.t  $u$  and  $\sigma$  and setting them to zero.

$$\Rightarrow u = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\sigma^2 = \frac{(x_1 - u)^2 + (x_2 - u)^2 + \dots + (x_n - u)^2}{n}$$

# MLE as a special case of MLE

Now, we know about MLE, let's show that for

MSE is a special case of MLE

When we assume the statistical model to be Gaussian.

In Linear Regression:

$$y = \beta_0 + \beta_1 n_1 + \dots + \beta_p n_p + \epsilon$$

$\epsilon$  P inputs not samples

Now; In LR we assume our model is Gaussian,

So:  $y \sim N(u, \sigma^2)$

→ Note this  $\sigma^2$  is not the  $\sigma^2$  of noise.  $\epsilon$ .

So: Likelihood of  $y$ :

$$L(u, \sigma^2 | X = n_1, n_2, \dots, n_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \prod_{i=1}^n e^{-\frac{(y_i - u)^2}{2\sigma^2}}$$

We are estimating  $u$  ~~and~~ by our parameter  $\beta$

$$\Rightarrow L(\beta, \sigma^2) = \prod_{i=1}^n \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) \prod_{i=1}^n e^{-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}}$$

∴ log-likelihood

$$\log L(\beta, \sigma^2) = \log \left[ \prod_{i=1}^n \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) \prod_{i=1}^n e^{-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}} \right]$$

$$= n \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \sum_{i=1}^n -\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}$$

Now If we say negative of likelihood as our loss:

Then:  $\text{Loss} = -\log \mathcal{L}(\beta, \sigma^2)$

$$\text{Loss} = -n \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^T x_i)^2 \quad (-B)$$

Now Recall that: We are interested in optimal solution not optimal value.

So, we are interested in those values of  $\beta$  (parameters of L.R) which minimize loss  $(-B)$  or equiv. max. likelihood in  $(A)$ .

Which we can find by setting partial derivative of Loss  $\Rightarrow$  to 0,

That is parameters  $\beta$  that minimizes  $(B)$

also minimize any scalar  $\Rightarrow$  multiplication or addition to function  $(B)$ .

i.e

$$\begin{aligned}\hat{\beta} &= \underset{\beta}{\operatorname{argmin}} \mathcal{L}(\beta, \sigma^2) \\ &= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \beta^T x_i)^2\end{aligned}$$

$$= \underset{\beta}{\text{argmin}} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

↓  
Sum of square Error.

In  
 Hence: Linear Regression, ~~MLE~~ MSE  
 is a special case of MLE when we  
 assume the model is Gaussian.

In Conclusion:

In Linear Regression, Logistic Regression,  
 Poisson's Regression,

we are always estimating the  
 parameters of the model which maxi-  
 mizes the likelihood of the training  
 data.

So, all the methods are based on  
 Maximum likelihood Estimation.