



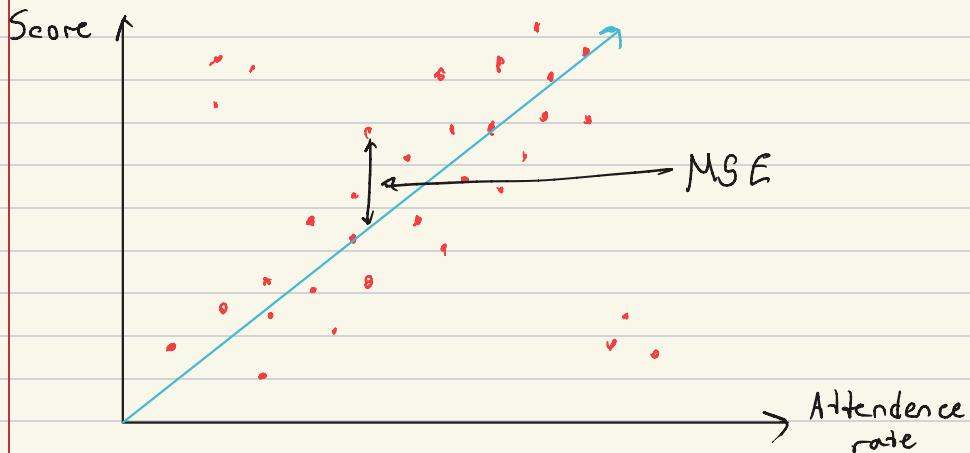
Mini project 00

Linear regression

## Part 1 : Linear Regression

Input (feature)  $\rightarrow$  attendance\_rate

Output (Target)  $\rightarrow$  capstone-score



line  $\rightarrow y = m \cdot x + b$   $\hat{y}_i$  (Prediction)

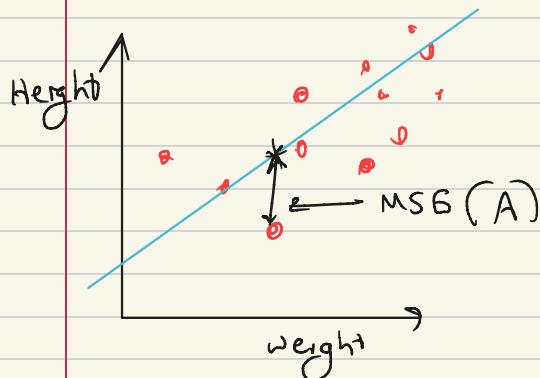
Hypothesis  $\rightarrow$  Capstone-score =  $y = B_0 + B_1 \cdot \text{Attendance}$   
features

MSE  $\rightarrow \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$

\* We have to find the correct  $m$  &  $b$  to minimise MSE

## Gradient Descent

- We optimize the best line to fit using gradient descent



Predicted height = intercept + slope  $\times$  weight

$$Y = B_0 + B_1 \cdot X$$

Finding the intercept using gradient descent

$$Y = B_0 + B_1 \cdot X$$

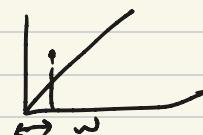
Slope

$$B_1 \rightarrow 0.64 \text{ (From Least squares)}$$

$$B_0 \rightarrow 0 \text{ (Any number)}$$

Predicted height (Y)

$$\begin{aligned} Y &= 0 + 0.64 \times \text{weight (0.5)} \\ &= 0 + 0.64 \times 0.5 \\ &= 0.32 \end{aligned}$$



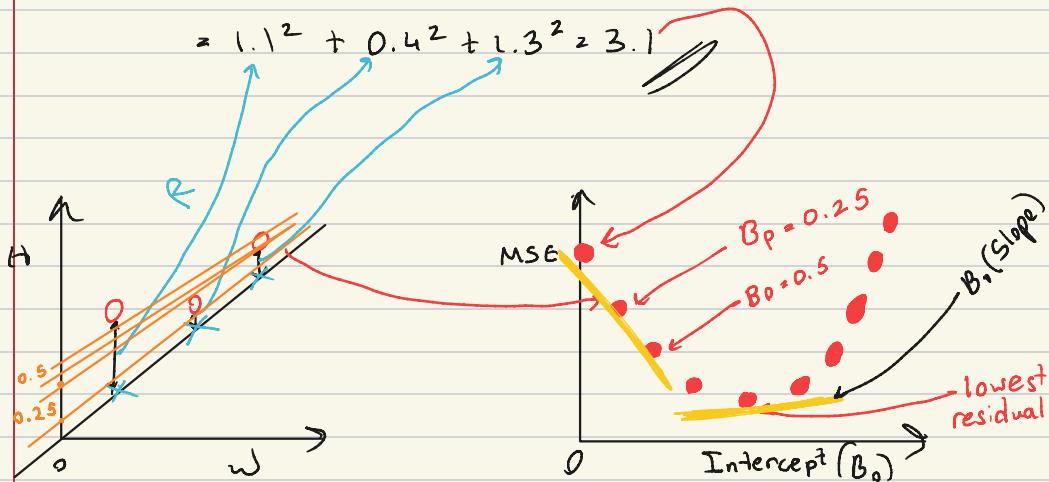
Then residual  $\rightarrow y - \hat{y}$

$$= 1.4 - 0.32 = 1.1$$

↑  
actual H      ↑  
predicted H

\* Now we need sum of the squared residuals  $\rightarrow$  MSE Function

$$= 1.1^2 + 0.4^2 + 1.3^2 = 3.1$$



\* Gradient descent only does a few calculations far from the solution

And increases the number of solution closer to optimal value

1<sup>st</sup> data point (•)

$$\text{Sum of squared residuals} = (\underbrace{\text{actual height}}_{\text{1st data point (•)}} - \underbrace{\text{predicted height}}_{(\text{intercept} + 0.64 \times 0.5)})^2$$

This is MSE {

$$+ 2^{\text{nd}} \text{ data point (•)} + 3^{\text{rd}} + \dots$$

\* We can give any value for intercept ( $B_0$ ) and get residual sum

\* Now can take the derivative of the function and determine the slope at any value for the intercept

$$\text{Sum of squared residuals} = (1.4 - (\text{intercept} + 0.64 \times 0.5))^2$$

Taking derivatives

1st data point

$$\frac{\partial}{\partial B_0} \text{Sum of residuals}^2 = \frac{\partial}{\partial B_0} (1.4 - (B_0 + 0.64 \times 0.5))^2$$

+ 2<sup>nd</sup> dev + 3<sup>rd</sup>

$$= 2(1.4 - (\text{intercept} + 0.64 \times 0.5)) \times 1$$

$$= \frac{\partial}{\partial B_0} \cdot (1.4 - (\text{intercept} + B_0 + 0.64 \times 0.5))$$

+ 2<sup>nd</sup> dev + 3<sup>rd</sup>

~~2<sup>nd</sup> dev~~

Chain rule for solving derivatives

\* We use chain rule method for solving derivatives easily

Slope  $\rightarrow$  derivative

\* multiplying rates of both outer and inner parts of the equation

$$1) (1.4 - (B_0 + 0.64 \times 0.5))^2$$

$$= (1.4 - (B + 0.32))^2$$

$$= (1.08 - B)^2$$

Outer function  $(u)^2 \rightarrow 2u$

Inner function:  $1.08 - B \rightarrow -1$

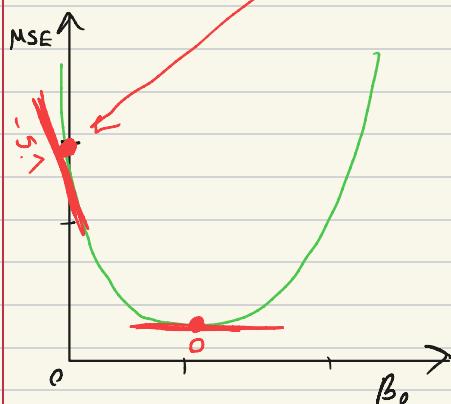
$$= 2(1.4 - (B + 0.64 \times 0.5)) \times -1$$

$$= -2(1.4 - (B + 0.64 \times 0.5))$$

Total sum of the derivative is  $B_0 = 0$

$$\text{Sum} = 1^{st} + 2^{nd} + 3^{rd}$$

Slope  $\rightarrow -5.7$  far from 0  $\rightarrow$  maximal value



- \* When the slope of curve is 0 we find the optimal value for  $B_0$ , we
- \* we should keep getting closer to 0

\* Determine a step size to keep getting close to 0

$$\text{Step size} = -5.7 \times n$$

$n \rightarrow \text{learning rate}$

$$\text{Step Size} = -5.7 \times 0.1 = -0.57$$

$$\text{New } B_0 = \text{Old } B_0(0) - \text{Step size} (-0.57)$$

-  $0 - (-0.57)$

=  $0.57$

$\rightarrow$  This is passed as a new value for  $B_0$  and the derivative is calculated again

\* This is done until slope is 0

\* Gradient descent will stop after it reaches below 0.001

\* Also have a limit of number of steps

- Max No. of steps  $\rightarrow 1000$  or greater

## Equations

### 1) Cost function (MSE)

$$J(B) = \frac{1}{2m} \sum_{i=1}^m (h_B(x^{(i)}) - y^{(i)})^2$$

### 3) Gradient with respect $B_0$ (intercept)

$$\frac{\partial J(B)}{\partial B_0} = \frac{1}{m} \sum_{i=1}^m (h_B(x^{(i)}) - y^{(i)})$$

3) Gradient with respect to  $\beta_j$  (Slope / weights)

$$\frac{\partial J(\beta)}{\partial \beta_j} = \frac{1}{m} \sum_{i=1}^m (\hat{y}_\beta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

## Part 2: Polynomial Regression

\* Instead of fitting a straight line, we fit a polynomial curve by adding powers of  $x$

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots + \beta_d x^d$$

$d \rightarrow$  degree

$d = 1 \rightarrow$  straight line

$d = 2 \rightarrow$  parabola

$d = 3 \rightarrow$  cubic curve

## Linear regression with multiple features

- We have multiple features

$$h_{\beta}(x^{(i)}) = \beta_0 + \beta_1 x_1^{(i)} + \beta_2 x_2^{(i)} + \dots$$

- We make it in Matrix form

$$h_{\beta}(x^{(i)}) \rightarrow \text{Feature Matrix}$$

shape  $(n, m+1)$

$n \rightarrow$  training examples

$m \rightarrow$  no. of features

$$\beta \rightarrow (m+1, 1)$$

$$y \rightarrow (n, 1)$$

$$\hat{y} = X \cdot \beta$$

Ex:-

$m = 3 + 1 \rightarrow +1$  for bias column, because to match inner dimensions to multiply

$$\tilde{X} = \begin{bmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ 1 & x_{31} & x_{32} \end{bmatrix} \quad \underbrace{\beta}_{(3, 1)} = \begin{bmatrix} B_0 \\ B_1 \\ B_2 \end{bmatrix}$$

$$\hat{y} = \tilde{X} \cdot \beta$$

### Cost function

$$\text{MSE} \rightarrow J(\beta) = \frac{1}{2n} (\tilde{X}\beta - y)^T (\tilde{X}\beta - y)$$

### Gradient rule

$$\nabla J(\beta) = \frac{1}{n} \tilde{X}^T (\tilde{X}\beta - y)$$

Update

$$\beta = \beta - \eta \cdot \nabla J(\beta)$$

## Feature Scaling

- \* We have to use feature scaling to standardize value when there are multiple values

Ex :- quiz score 0 - 100 vs study hours are different

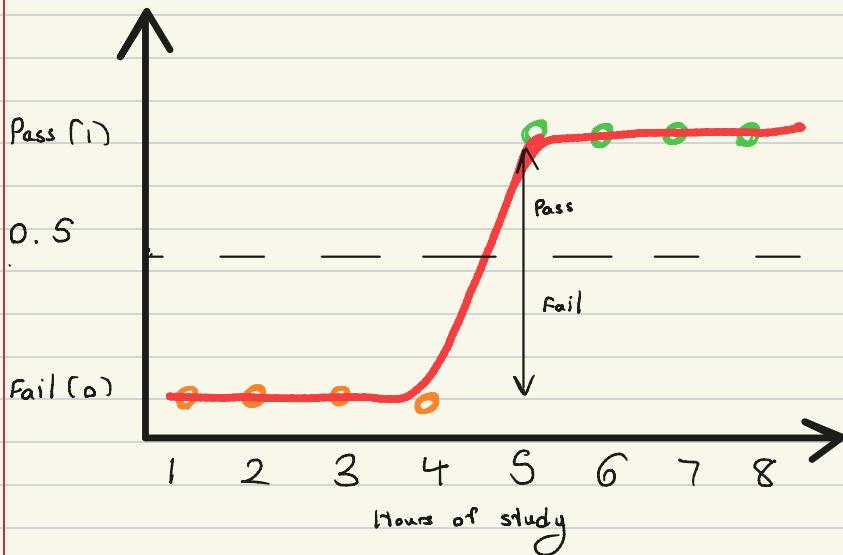
$$x' = \frac{x - \mu}{\sigma}$$

- \* makes all features have mean 0 and standard deviation 1 → faster convergence

## Part 6: Algorithms in Classifications

### Logistic regression

- This predicts if something is true or False



- Instead of fitting a straight line into the data, Logistic Regression fits an S-Shaped curve to the data

S shaped curve → Sigmoid function

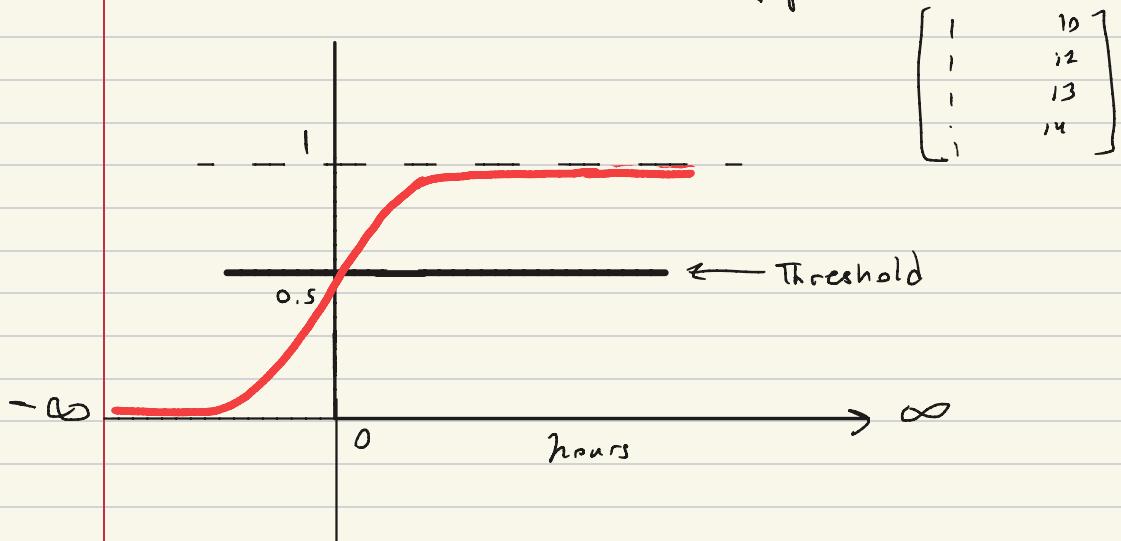
\* Sigmoid function takes any real number and makes it into a range between 0 and 1

Essential for interpreting output as a probability

## Mathematical Implementations

$\theta^T x^{(i)} = z^{(i)}$  → We use sigmoid function  
n. parameters      n features  
to convert ~~z~~ into a number between 0 & 1

$\Theta \rightarrow$  theta will include both bias and slopes =  $x_{bias} + w$ )



→ Sigmoid Function

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

## Model hypothesis

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

estimated probability for  $x^{(i)}$       sigmoid function      gives the probability

\* To optimize parameters sigmoid function is not enough

$$J(\beta) = \sum_{i=1}^m \left[ h_{\beta}(x^{(i)}) \right]^{y^{(i)}} \left[ 1 - h_{\beta}(x^{(i)}) \right]^{1-y^{(i)}} = 1/0$$

actual value      predicted prob

\* This function returns a value close to 0 and 1

- This function is not very optimal for calculating gradient descent

- 1) Negligible to minimize gradient descent  $\rightarrow -(-)$
- 2) Numerically unstable  $\rightarrow \log(1/n)$
- 3) Scaling the function  $\rightarrow \perp$

## Binary Cross-Entropy Loss

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

## Gradient of Cost function

$$\text{Derivative of } S(\theta) \rightarrow \nabla S(\theta) = \frac{1}{m} X^T (h_\theta(X) - y)$$

Update gradient

$$\theta := \theta - n \cdot \nabla S(\theta)$$

\* This tells us how much the gradient should change to reduce error

