

# Data Science & AI



**Machine Learning**



**Supervised Learning**

**Lecture No.- 03**



**By- Krish Naik Sir**

# Recap of Previous Lecture



Topic

Topic

Topic

Topic

Topic

## Regression

# Topics to be Covered



Topic

Polynomial

Topic

Logistic

Topic

Topic

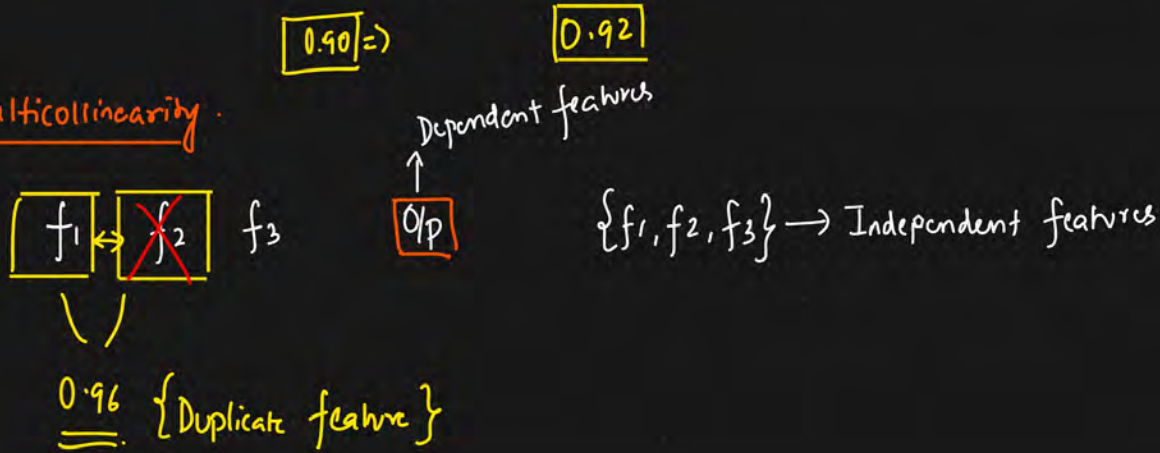
Topic

# Machine Learning

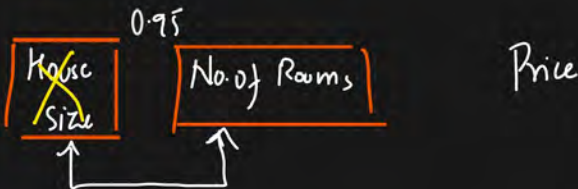
## Agenda

- i) Multicollinearity ✓ [Bias Variance Tradeoff] ✓
- ii) Polynomial Regression ✓
- iii) Logistic Regression ✓
- iv) Performance metrics [Accuracy, True Positive, False Positive, Precision, Recall] ✓.

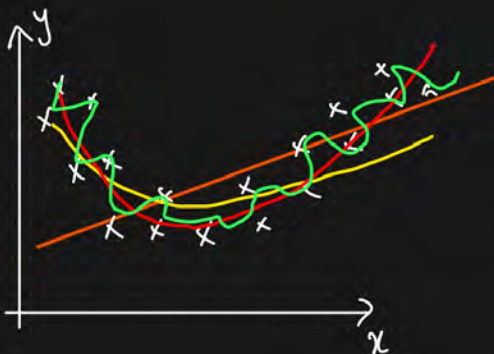
### ① Multicollinearity



### House Price Prediction



### ② Polynomial Regression



#### Underfitting

TRAINING Acc  $\downarrow \downarrow = 60\%$   
Test Acc  $\downarrow \downarrow = 65\%$  55%

#### Overfitting

low Bias  
Training Acc  $\uparrow \uparrow \Rightarrow 95\%$   
Test Acc  $\downarrow \downarrow \Rightarrow 60\%$   
High Variance

$$h_0(x) = \theta_0 + \theta_1 x_1$$

Underfitting

Normalized model

$\Rightarrow$

High Bias  
High Variance

High Bias  
Low Variance  $\Rightarrow$  Underfitting



Overfit-

Generalized or Perfect Model

TRAINING Acc  $\uparrow\uparrow$  } Low Bias  
TEST Acc  $\uparrow\uparrow$  } Low Variance

## Polynomial Regression

Polynomial degree = 0

$$h_{\theta}(x) = \theta_0 \times x_1^0 = \theta_0 \times 1 = \theta_0 \rightarrow \text{Intercept}$$

Polynomial degree = 1

$$h_{\theta}(x) = \theta_0 \times x_1^0 + \theta_1 x_1^1 = \theta_0 + \theta_1 x_1 \rightarrow \text{Simple Linear Regression}$$

Polynomial degree = 2

$$h_{\theta}(x) = \theta_0 \times x_1^0 + \theta_1 x_1^1 + \theta_2 x_1^2 \leftarrow$$

Polynomial degree = 3

$$h_{\theta}(x) = \theta_0 x_1^0 + \theta_1 x_1^1 + \theta_2 x_1^2 + \theta_3 x_1^3$$

⋮

$x_1 \quad x_2 \quad x_3$

Polynomial degree = 2

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_1^2 + \theta_5 x_2^2 + \theta_6 x_3^2$$

⑧ Logistics Regression  $\rightarrow$  Classification  $\rightarrow$  Binary Classification

Dataset

UPSC

No. of study hours

Pass/Fail

1

Fail

2

Fail

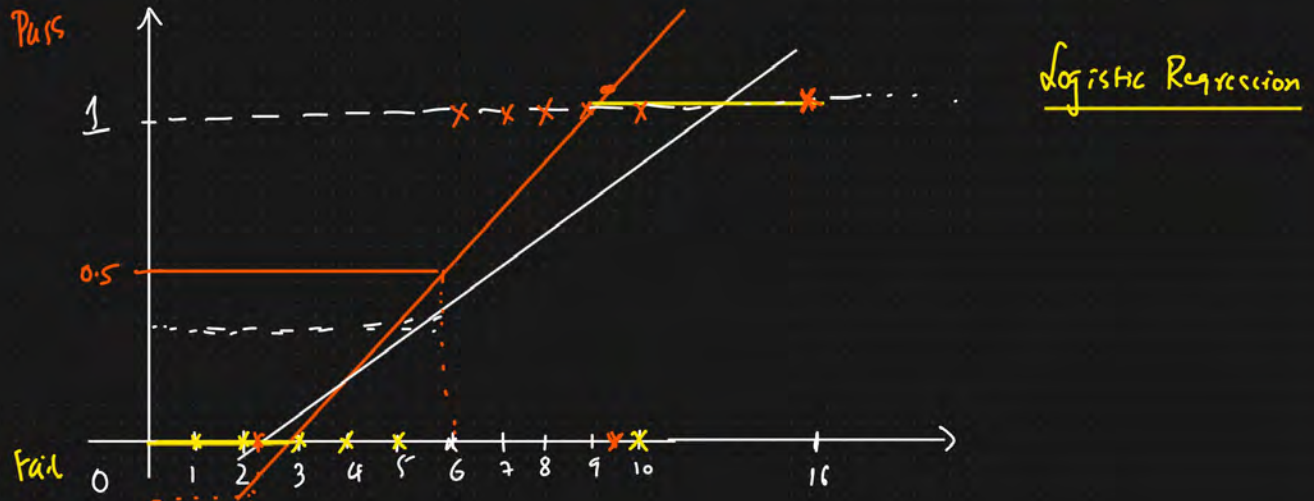
3

Fail

4

Fail

5	Fail
6	Pass
7	Pass
8	Pass
9	Pass
10	Fail



Why we cannot use Linear Regression for classification

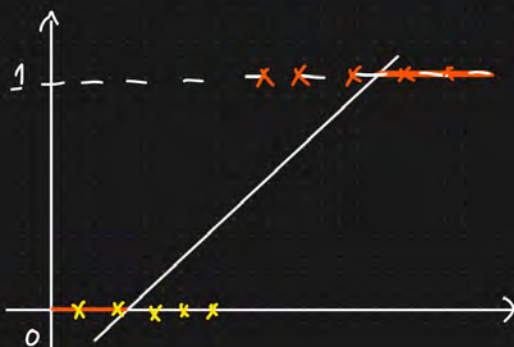
- ① Best fit line changes because of outlier  $\rightarrow$  prediction goes wrong
- ② The outcome comes  $> 1$  and  $< 0$

To solve this problem we use Logistic Regression

$$\downarrow$$

$$\boxed{0 \text{ to } 1} \Rightarrow \text{Squashing Technique}$$

How Logistic Regression solves classification Problem



$$L = \boxed{h_0(x) = \theta_0 + \theta_1 x_1} \Rightarrow \text{Best fit line}$$

$$\downarrow$$

$$[\text{Sigmoid Activation function}]$$

$$\downarrow$$

$$\boxed{0 \text{ to } 1}$$



## Logistic Regression

$$h_{\theta}(x) = \sigma(\theta_0 + \theta_1 x_1)$$

$$\sigma = \frac{1}{1 + e^{-z}} \Rightarrow 0 \text{ to } 1.$$

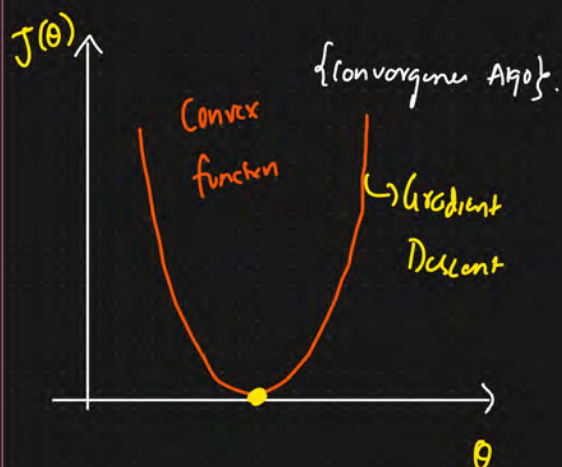
$$h_{\theta}(x) = \frac{1}{1 + e^{-z}}$$

$$\Rightarrow h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}}$$

## Linear Regression Cost fn

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

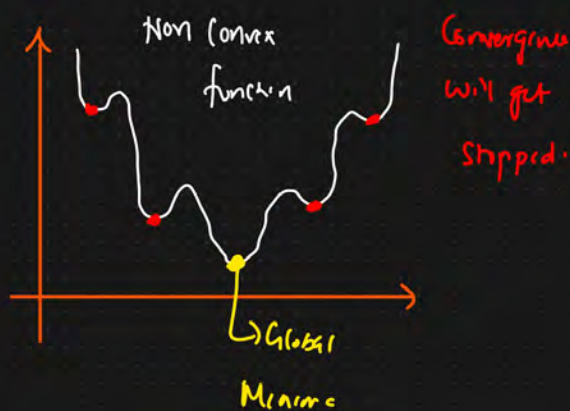
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$



## Logistic Regression Cost fn

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}}$$



log loss  $\rightarrow$  Cost function  $\rightarrow$  Logistic Regression

$h_{\theta}(x) \Rightarrow$  predicted  $o/p$

$$J(\theta_0, \theta_1) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

$$J(\theta_0, \theta_1) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x)) \Rightarrow \text{Convex function}$$



Minimize  $J(\theta_0, \theta_1)$  changing  $\theta_0$  and  $\theta_1$

↓  
1 global Minima.



## \* Performance Metrics, Accuracy, Precision, Recall and F-Beta

### Topics to be covered

① Confusion Matrix ✓

② Accuracy ✓

③ Precision

④ Recall

⑤ F-Beta Score

### DATASET

		Actual	n Predicted	
$x_1$	$x_2$	y	y	
-	-	0	1	→ Wrong Prediction
-	-	1	1	→ Correct Prediction
-	-	0	0	
-	-	1	1	
-	-	1	1	
-	-	0	1	
-	-	1	0	

### ① Confusion Matrix

	1	0	→ Actual Value
1	3	2	
0	1	1	

↓  
Predicted Value

	1	0
1	TP	FP
0	FN	TN

$$\text{Accuracy Score} = \frac{TP + TN}{TP + FN + FP + TN}$$

$$\text{Accuracy Score} = \frac{3 + 1}{3 + 2 + 1 + 1} = \frac{4}{7} =$$



## DATASET →

$$1000 \text{ Data points } \left\{ \begin{array}{l} \rightarrow 900 - 1 \\ \rightarrow 100 - 0 \end{array} \right\} \text{ Imbalanced dataset}$$

Dumb Model-1  $\Rightarrow$  Accuracy = 90%.

③ Precision =  $\frac{TP}{TP+FP}$  Out of all actual values how many are correctly predicted.

	1	0
1	TP	FP
0	FN	TN

False Positive  $\rightarrow$  Important ↓↓↓↓.

④ Recall =  $\frac{TP}{TP+FN}$   $\Rightarrow$  FN ↓↓↓.

Use case 1 : Spam classification  $\left\{ \begin{array}{l} \rightarrow \text{Spam} \\ \rightarrow \text{Not Spam.} \end{array} \right.$

	1	0
1	TP	FP
0	FN	TN

FP ↓↓

FN ↓↓

Mail  $\rightarrow$  Spam {Truth}  $\Rightarrow$  Good Scenario.  
Mail  $\rightarrow$  Spam

0  $\leftarrow$  Mail  $\rightarrow$  Not a Spam  $\Rightarrow$  False positive.  
Model  $\rightarrow$  Spam

[Precision]

Mail  $\rightarrow$  Spam  $\Rightarrow$  FN

Model  $\rightarrow$  Not a Spam

Use case 2: Predict whether a person has diabetes or not  $\Rightarrow$  FNN

$\downarrow$   
Recall

X Bank  
 $\downarrow\downarrow\downarrow$

Assignment: Tomorrow the stock market is going to crash.

$\downarrow\downarrow$

F - Beta Score

$\nearrow$  Harmonic Mean.

$$F1 \text{ Score} = 2 \times \left[ \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \right]$$

#### Question1:

In a simple linear regression model, the R-squared ( $R^2$ ) value is 0.70, and the Adjusted R-squared ( $R^2$ ) is 0.68. What does this difference between  $R^2$  and Adjusted  $R^2$  suggest?

- a) The model is underfitting the data.
- b) The model is perfectly accurate in making predictions.
- c) The independent variable is not relevant to the dependent variable.
- d) The model includes unnecessary polynomial terms(features).



**Question2:**

**What is the primary purpose of the Adjusted R-squared ( $R^2$ ) in regression analysis?**

- a) To quantify the proportion of the variance in the dependent variable explained by the independent variables.
- b) To provide a measure of the model's accuracy in predicting the dependent variable.
- c) To account for the number of predictors in the model, penalising excessive complexity.
- d) To calculate the residual sum of squares (RSS) of the regression model.



# THANK - YOU