



Joint Tech Internship Community Program

Assignment 1

SUBMITTED BY :

JAHAGANAPATHI SUGUMAR

CANDIDATE ID : 2024060193

THE GIVEN TABLE :

Make Year	Brand	Variant	Mileage	Fuel	Transmission	Resale Price (INR)
2015	BMW	520D	80000	Diesel	Automatic	2500000
2016	Audi	A6	92000	Petrol	Automatic	1900000
2018	Mercedes Benz	E200	61000	Petrol	Automatic	3400000
2014	Skoda	Superb	95000	Petrol	Automatic	600000
2020	Benz	E200	35000	Petrol	Automatic	12000000

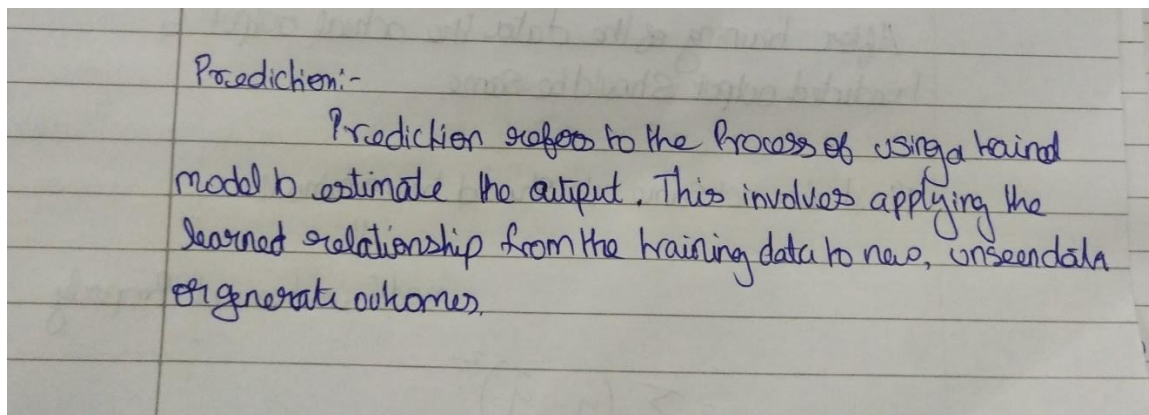
Feature :

From the given table the inputs are Make Year, Brand, Variant, Mileage, Fuel, Transmission. these input determines the price of the car. And price depends on the inputs.

LABEL :

From the given table the Label is the Resale price which is predicted with the help of its feature.

Prediction :



Outlier :

From the given table since there exist of one value which is different from the majority values it is univariate outlier

Make Year	Brand	Variant	Mileage	Fuel	Transmission	Resale Price (INR)
2015	BMW	520D	80000	Diesel	Automatic	2500000
2016	Audi	A6	92000	Petrol	Automatic	1900000
2018	Mercedes Benz	E200	61000	Petrol	Automatic	3400000
2014	Skoda	Superb	95000	Petrol	Automatic	600000
2020	Benz	E200	35000	Petrol	Automatic	12000000

Dictation can be done by the Z-Score if it is ± 3.29 or beyond

$$\text{Z-Score} = \frac{x_i - \bar{x}}{s}$$

$\bar{x} \rightarrow \text{mean} \rightarrow 4,080,000$

$x_i \rightarrow \text{outlier value} \rightarrow 12000000$

$s \rightarrow \text{standard deviation} \rightarrow 4,063,200.71$

$$\text{Z-Score} = \frac{12000000 - 4,080,000}{4,063,200.71} = 1.95$$

it is not considered an outlier based on the Z-score analysis.

1) Domain knowledge 3) Math/Statistics

2) Visualization

Two standard deviation

Outliers:-

→ Outliers are cases that have data values that are very different from the data values for the majority of cases in the data sets.

→ An outlier is an observation that is substantially different from the other observations.

→ Outliers are important because they can change the results of our data analysis.

Types of outliers

univariate

multivariate

univariate:-

univariate outliers are cases that have an unusual value for a single variable.

univariate can be predicted with it
Z-Score is ± 3.29 or beyond

Z-score = $(i - \text{mean}) / \text{standard deviation}$

$$Z_i = \frac{x_i - \bar{x}}{s}$$

multivariate

→ multivariate outliers are cases that have an unusual combination of values for a number of variables

→ It can be detected by the Mahalanobis Distance (D^2) Probability is less than 0.001

Mahalanobis distance:-

$$D_m = \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$$

where μ is mean

Σ is covariance matrix

Generalization of Z-scores to multi-dimensional space.

→ Replace univariate mean with multivariate mean

→ Replace standard deviation with covariance

Z-score

Mahal fun

$$Z_i = \frac{x_i - \bar{x}}{s} \rightarrow \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$$

Training Data :

From the given table assuming that 80 percent of the data is training data.

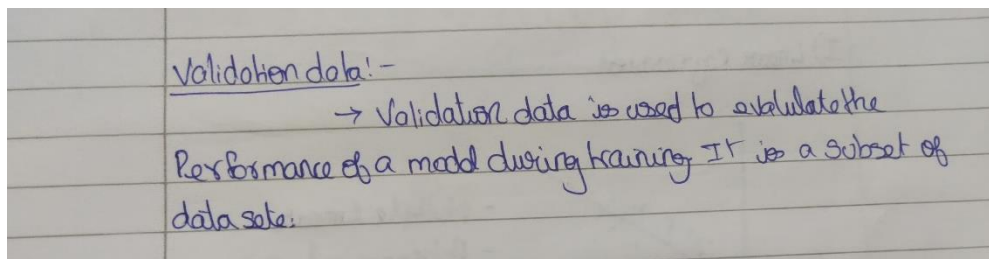
Make Year	Brand	Variant	Mileage	Fuel	Transmission	Resale Price (INR)
2015	BMW	520D	80000	Diesel	Automatic	2500000
2016	Audi	A6	92000	Petrol	Automatic	1900000
2018	Mercedes Benz	E200	61000	Petrol	Automatic	3400000
2014	Skoda	Superb	95000	Petrol	Automatic	600000

Test Data:

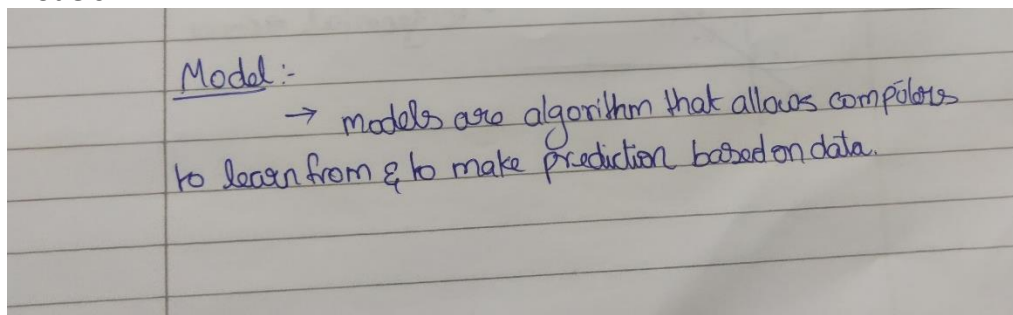
From the given table assuming that 80 percent of the data is training data and remaining 20 percent is Test data.

Make Year	Brand	Variant	Mileage	Fuel	Transmission	Resale Price (INR)
2020	Benz	E200	35000	Petrol	Automatic	12000000

Validation Data :



Models:



Types of models:-

- 1) Supervised Learning
- 2) unsupervised.

Supervised Learning:-

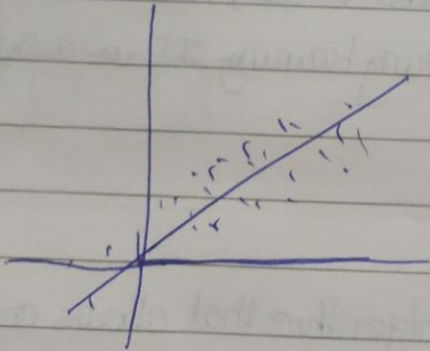
- └─> Regression
- └─> Classification

Regression:-

→ we try to make a relation between Dependent variable & Independent variable.

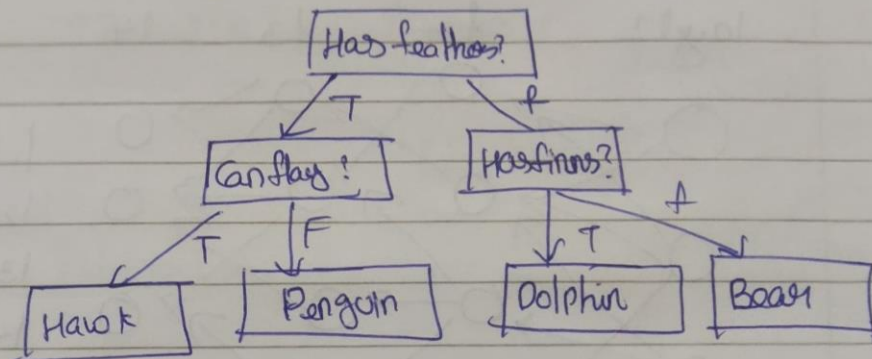
→ continuous output

I) Linear Regression



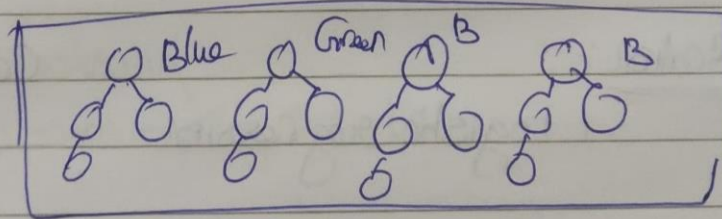
- Multiple Linear
- Polynomial \Rightarrow curve

II) Decision Tree



more nodes \rightarrow more accurate

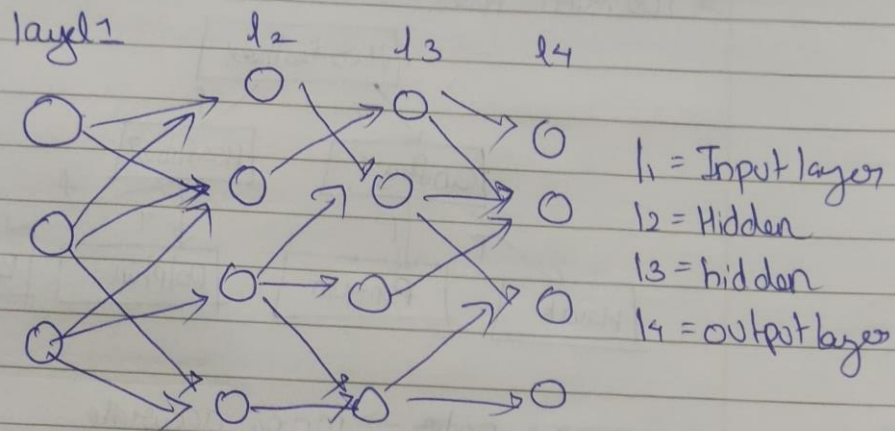
III) Random forest



Blue

majority wins Model

iv) Neural Network



hidden layer represent the function where the input goes through and leading the output.

Classification:-

(discrete)

I. Logistic regression

- To find the Probability of finite number of outcomes
- varies between 0 & 1

II. Support Vector machine

- N-dimensional space that distinctly classifies data points.

III. Naïve Bayes

→ Based on Bayes Theorem

$$P(B|A) = \frac{P(A|B) P(B)}{P(A)}$$

← Posterior Probability P(A) → Predictor

↑ Likelihood ↑ Prior Probability

on Supervised Learning:-

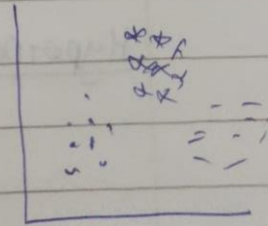
→ find Patterns from input data without reference to labeled outcomes

- 1) Clustering
- 2) Dimensionality Reduction

Clustering:-

Group data Points into clusters

- > K-means
- > Hierarchical
- > Mean Shift
- > Density-based



Hyperparameter :

Hyperparameter:-

→ Hyperparameter are Parameters whose values are set before the learning process begins

→ used to control the learning process

common Hyperparameter:-

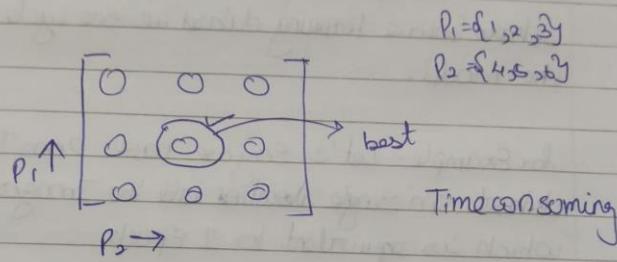
- 1) Learning Rate
- 2) Epochs
- 3) Batch size

Hyperparameter Tuning:-

- 1) Grid Search
- 2) Random Search

HT refers to the process of choosing the optimum set of hyperparameters for a machine learning model

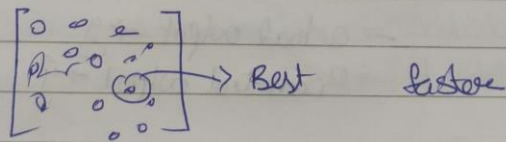
Grid Search cv:-



Random Search cv:-

→ randomly selects a specified number of combinations & evaluates

→ randomness makes the search process faster



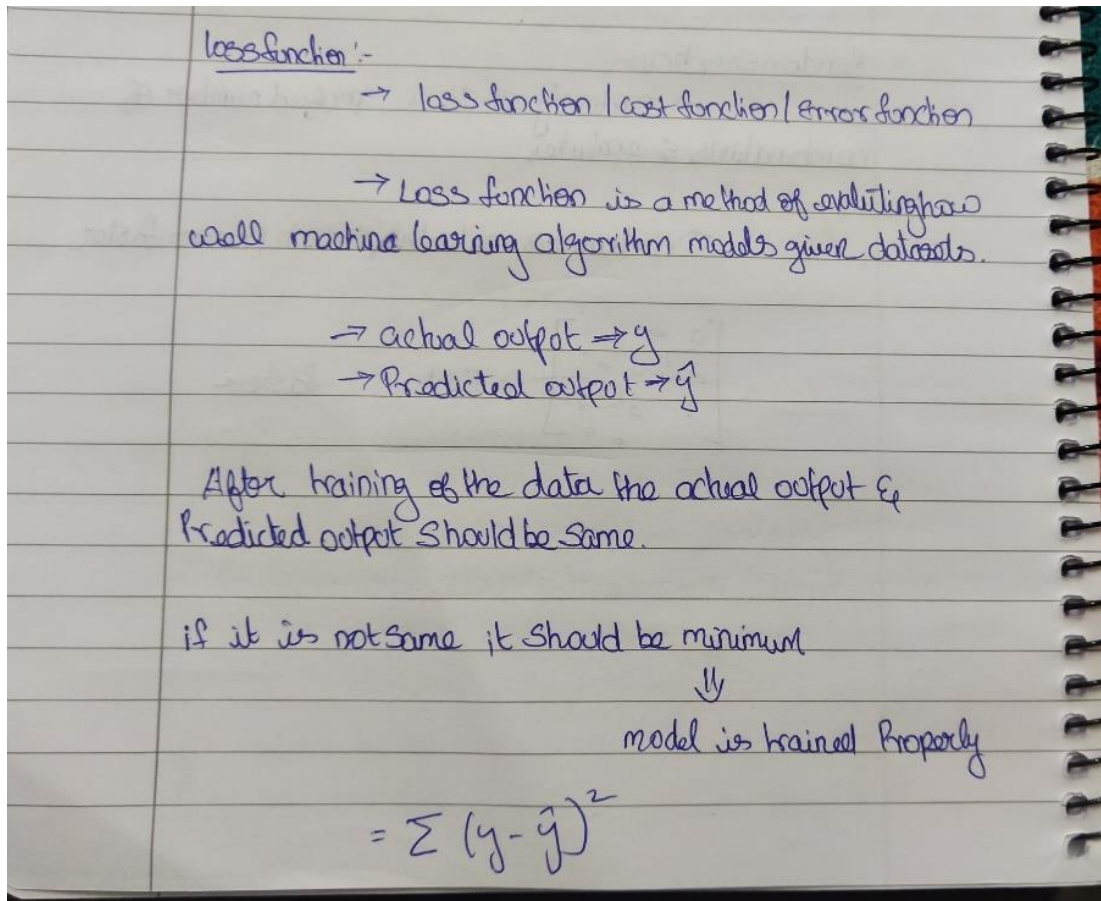
Epoch :

Epoch:-

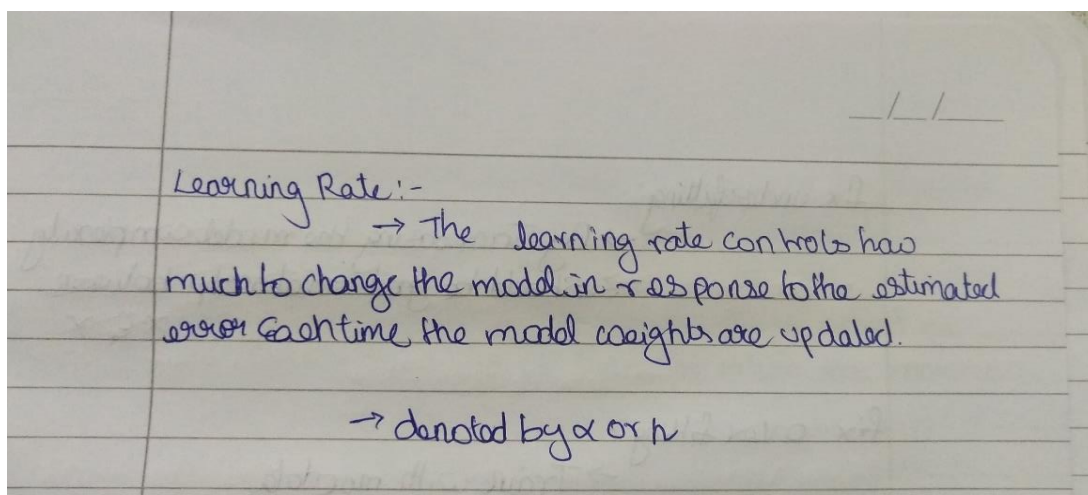
→ Epoch is a complete iteration through the ~~whole~~ online training dataset in one cycle for the training the model.

for Example:- Let's say we have 2000 Images as a dataset. In single iteration all the Images will be trained which is equivalent to 1 epoch.

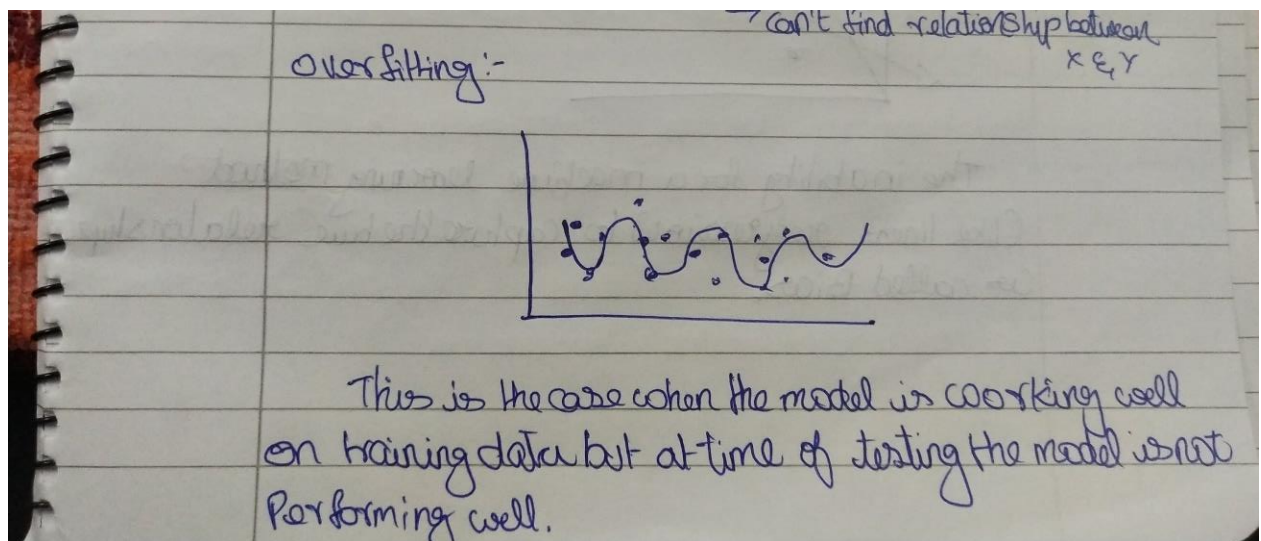
Loss Function:



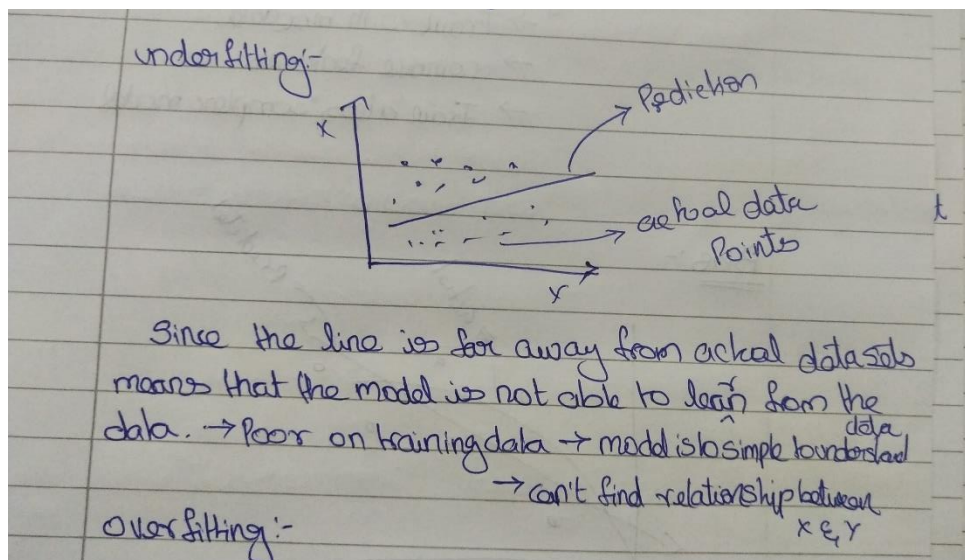
Learning Rate :



Overfitting :



Underfitting :



fix underfitting:-

- By Increasing the model complexity
- Should be good relationship between ~~R^2~~ R^2

fix overfitting:-

- Train with more data
- remove features
- Train a less complex model

Regularization :

Regularization:-

To ↓ Overfitting

Regularization avoids overfitting by adding a penalty to the model loss function:-

$$\text{Regularization} = \text{Loss function} + \text{Penalty}$$

3 techniques

→ L₂ regularization

→ L₁

→ Elastic Net

L₂ regularization:-

$$\text{Ridge Regression Cost function} = \text{Loss function} + \frac{1}{2} \lambda \sum_{j=1}^m w_j^2$$

where w is the slope of the line

λ control the strength of regularization

if λ controls the strength

$$\lambda = 0,$$

regularization will be eliminated

L₁ Regularization:-

$$\text{Lasso Regression Cost function} = \text{Loss function} +$$

$$\lambda \sum_{j=1}^m |w_j|$$

Elastic net Regularization

→ uses both L₁ & L₂

$$J(\beta_1, \beta_2; \beta_n) = \sum_{i=1}^n (y_i - \sum_{j=1}^m x_{ij} \beta_j)^2 + \lambda \left(\alpha \sum_{j=1}^m |\beta_j| + \frac{1-\alpha}{2} \sum_{j=1}^m \beta_j^2 \right)$$

⇓
L₂

$$\alpha \Rightarrow 0 \Rightarrow L_2$$

$$\alpha \Rightarrow 1 \Rightarrow L_1$$

Cross-Validation :

Cross validation :-

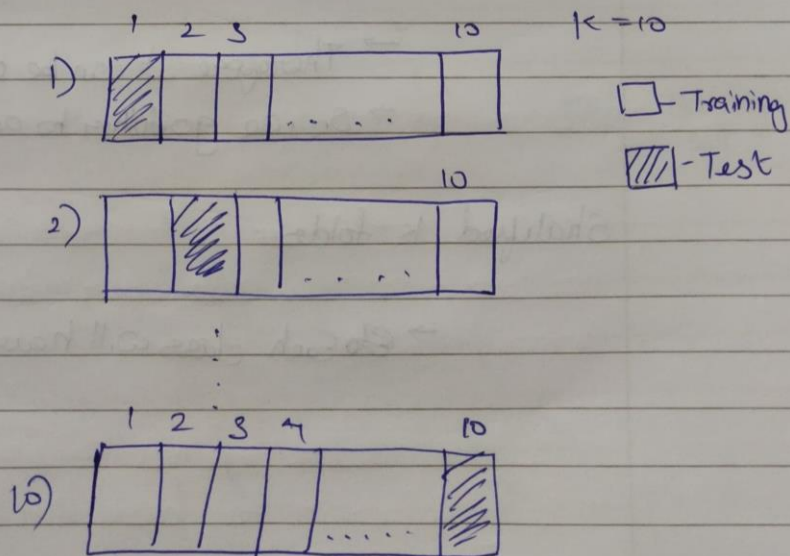
→ It involves dividing the available data into multiple fold or subsets, using one of these folds as a validation set & training the model on the remaining folds.
→ ~~Ensure the model selection~~

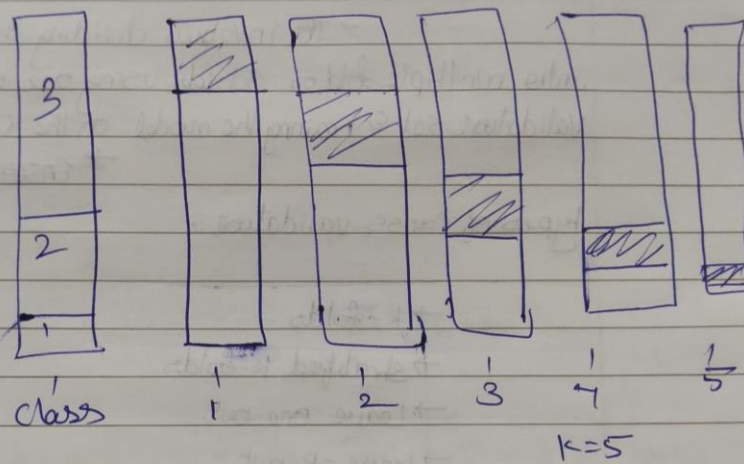
Types of Cross validation :-

- k-folds
- stratified k-folds
- Leave-one-out
- Leave-P-out

k-fold :-

→ we divide the entire dataset into k-folds





disadvantage:-

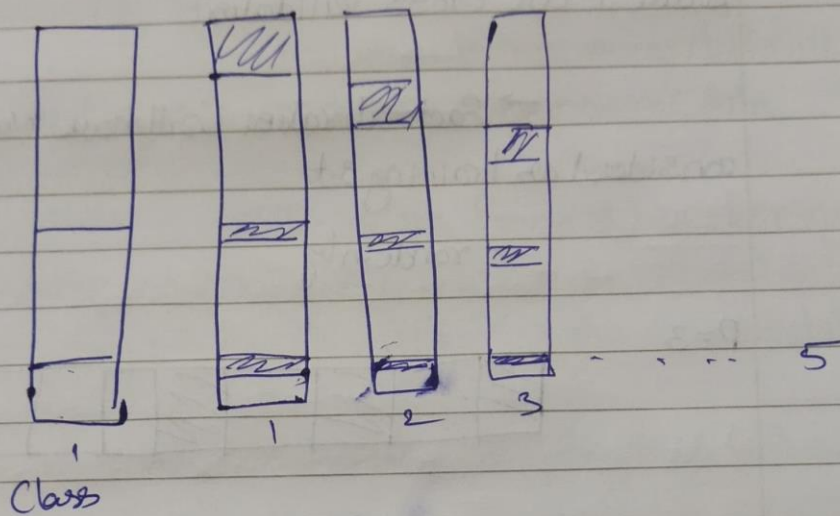
→ Each iteration is not taking Equal representation of the class \therefore

→ Therefore it can be over fitting again

→ So we go further to another Type

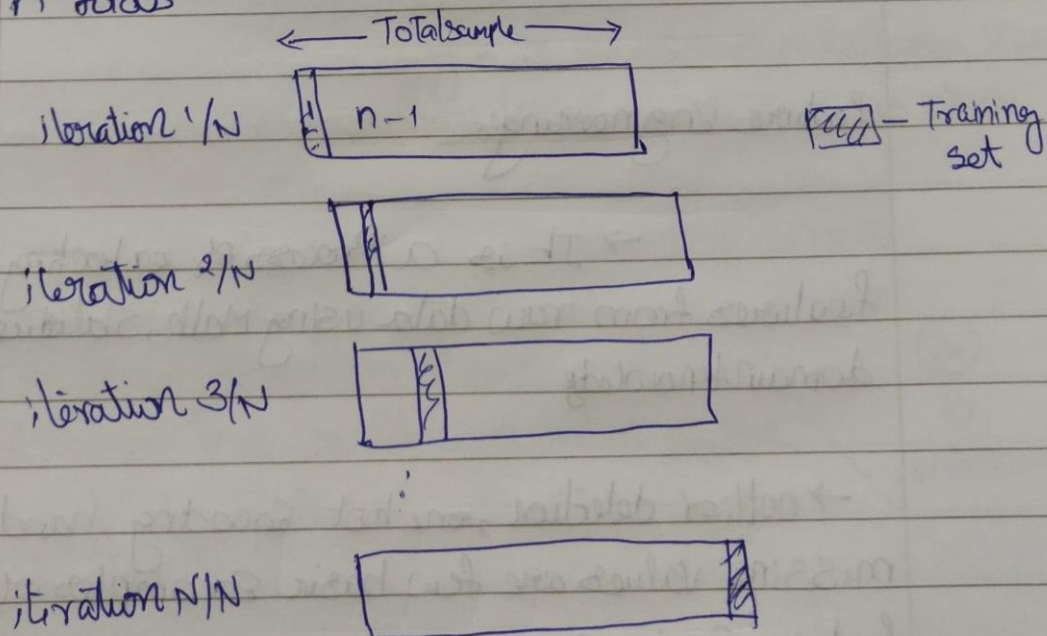
Stratified k-folds:-

→ Each class will have Equal representation



Leave one out cross validation :-

→ Total cross validation is divided into n folds

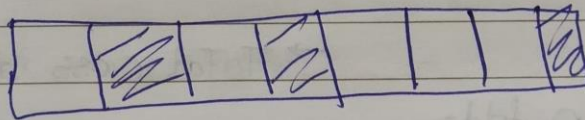
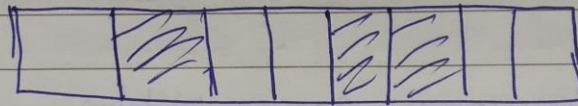
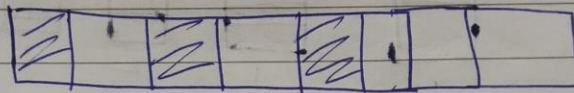


Leave P out cross validation

→ Each iteration will have No. P's will be considered as training set.

→ randomly

$P=3$



Feature Engineering :

Feature Engineering:-

→ It is a Process of extracting useful features from raw data using Math, Statistics and domain knowledge

→ outlier detection, one hot Encoding, handling missing values are few basic examples of Feature Engineering

Dimensionality Reduction :

Dimensionality Reduction:-

→ ML methods have some difficulty when dealing with such high-dimensional data

→ It is the process of transforming high-dimensional data into a lower dimensional data that still preserves the essence of the original data

→ used to reduce the no. of features

Types:-

- 1) Feature Selection
- 2) Feature Extraction

1) Feature Selection

→ Finding k of the total of n features that give us the most information & we discard the other $(n-k)$ dimensions.

2) Feature Extraction

→ Finding a new set of k features that are the combination of the original n features

$n = 10$

$10 \Rightarrow 7$

③

↓

discarded

Feature extraction

- Principal component Analysis (PCA)
- Linear Discriminant Analysis (LDA)

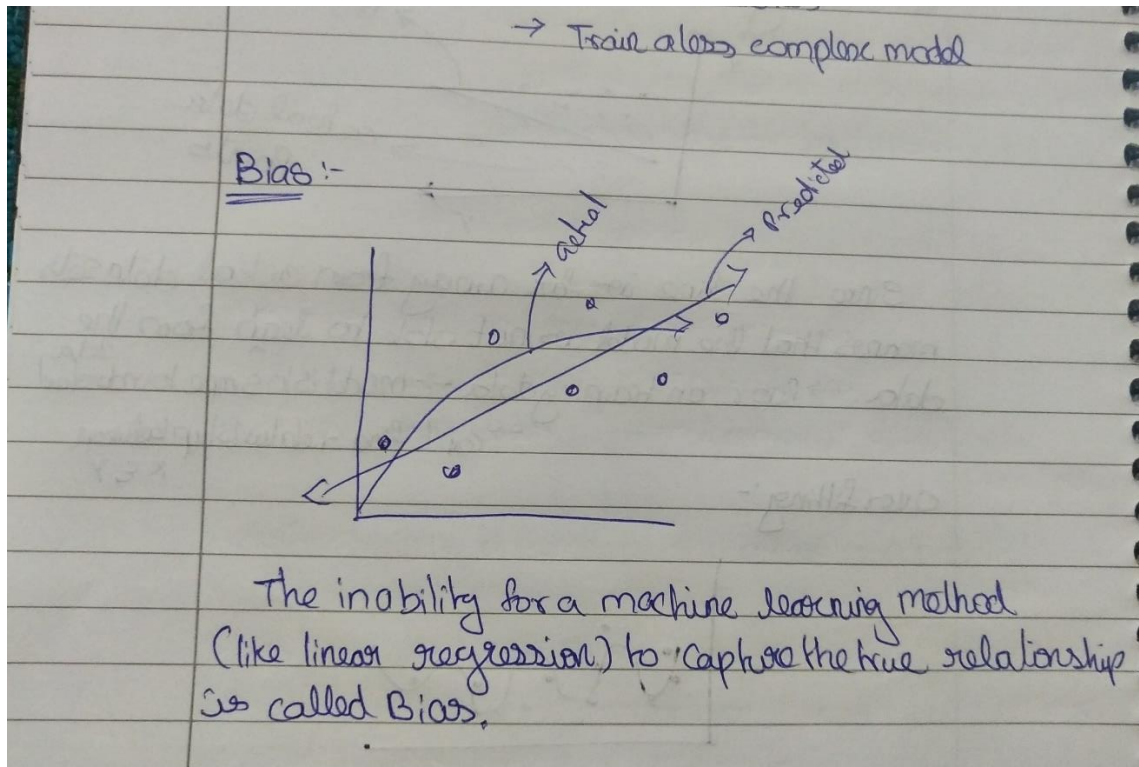
PCA:-

- Introduced by Karl Pearson
- It works on the condition that while the data in a higher dimensional space is mapped to data in a lower dimensional space, the variance of the data in the lower dimensional space should be minimum.

Steps :-

- Constructs covariance matrix of the data
- Find the Eigen vectors
- Eigen vectors corresponding to the largest eigen values are used to reconstruct a large fraction of variance of the original data.

Bias :



→ It is the difference between the expected Prediction of our model & the true value.

→ These differences between actual or expected values and the Predicted values are known as Bias.

Low Bias:- The model will closely match the training dataset

High Bias:- The model will not match the training dataset

high bias \Rightarrow underfitting

Variance :

Variance:-

→ variance refers to the changes in the model when using different portions of the training or test data

→ Variance is the variability of the model that how much it is sensitive to another subset of the training dataset.

Low variance:-

→ Low variance means that the model is less sensitive to changes in the training data

→ case of underfitting

High variance:- High variance means that the model is very sensitive to changes in the training data & can result in significant change in the estimate of the target function when trained on different subsets

→ case of overfitting.

high Bias, Low Variance \Rightarrow under fitting

high Var, Low Bias \Rightarrow over fitting

high var, high Bias \Rightarrow inconsistent & inaccurate Predictions

Low Bias, Low Variance \Rightarrow consistent & accurate Predictions.

Bias variance trade off:-

→ An model or Algorithm can't be more complex & less complex at the same time