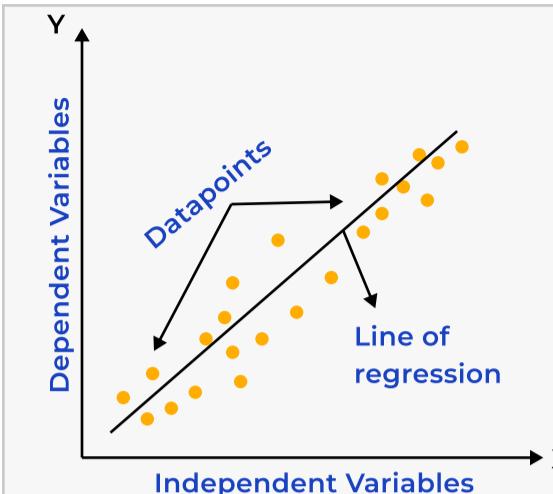


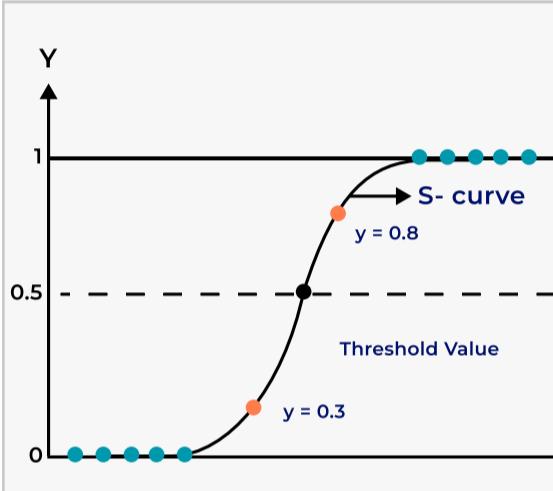
## 01. LINEAR REGRESSION



DESCRIPTION	
<ul style="list-style-type: none"> <li>Simplest algorithm which maps linear relation between input and continuous output by finding the best fine line</li> </ul>	
$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_n * x_n$	
<ul style="list-style-type: none"> <li>Assumptions: Linearity, No multicollinearity, Homoscedasticity, Normality of errors, no autocorrelation</li> </ul>	

TASK	DATA TYPE		OBJECTIVE FUNCTION	HYPERPARAMETER	ALGORITHM CHARACTERISTICS			WHEN TO USE		
	Regression	Numerical, need to change categorical variables to dummy variables			MSE	Learning rate( $\eta$ )	TYPE			
ALGORITHM CHARACTERISTICS	OUTLIER	MULTICOLLINEARITY	DATA TREATMENT	MODEL INTERPRETABILITY	REGULARIZATION TECHNIQUE					
Sensitive, pulls best fit line toward itself	Sensitive	Required	OUTLIER REMOVAL	MISSING VALUE TREATMENT:	FEATURE SCALING	model coefficients indicate magnitude and direction of each feature effect	Parametric	High	Low	<ul style="list-style-type: none"> <li>Model interpretation is required</li> <li>less volume of data (with less noise)</li> <li>EDA(correlation and scatter plots) indicate linear relationship between input and output</li> </ul>

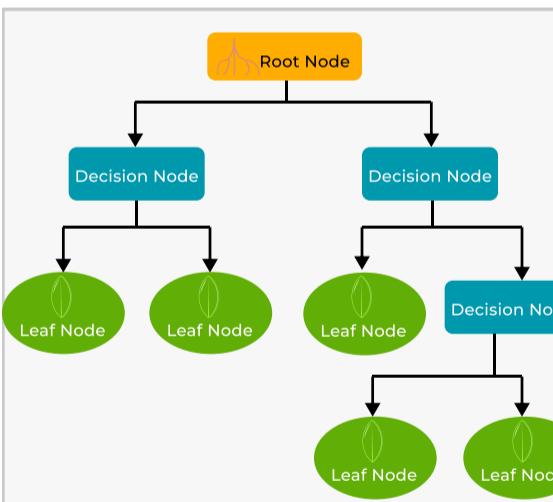
## 02. LOGISTIC REGRESSION



DESCRIPTION	
<ul style="list-style-type: none"> <li>Models linear relationship between input and a categorical output (0 or 1) using sigmoid function</li> </ul>	
$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_n * x_n$	

TASK	DATA TYPE		OBJECTIVE FUNCTION	HYPERPARAMETER	ALGORITHM CHARACTERISTICS			WHEN TO USE		
	Classification	Numerical, need to change categorical variables to dummy variables			Log-Loss or Binary Cross-Entropy (derived from MLE)	Learning rate( $\eta$ )	TYPE			
ALGORITHM CHARACTERISTICS	OUTLIER	MULTICOLLINEARITY	DATA TREATMENT	MODEL INTERPRETABILITY	REGULARIZATION TECHNIQUE					
Sensitive, pulls best fit line toward itself	Sensitive	Required	OUTLIER REMOVAL	MISSING VALUE TREATMENT:	FEATURE SCALING	model coefficients indicate magnitude and direction of each feature effect	Parametric	High	Low	<ul style="list-style-type: none"> <li>less data</li> <li>model interpretation is required</li> <li>classes look geometrically distinguishable</li> </ul>

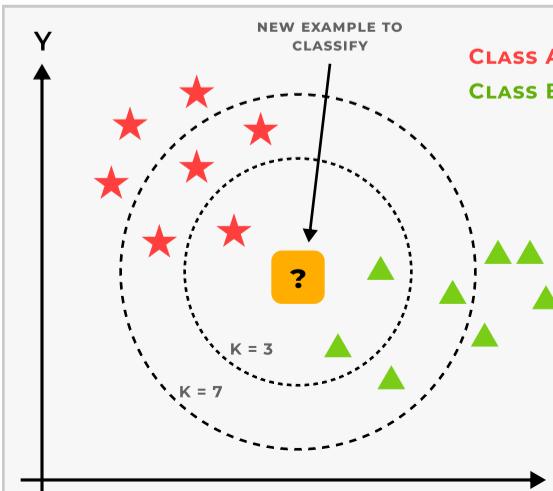
## 03. DECISION TREE



DESCRIPTION	
<ul style="list-style-type: none"> <li>Uses a set of rules on features to produce predictions.</li> <li>Each rule is considered as a node with split being a binary decision which terminates at leaf.</li> <li>Criterion like information gain, gini index are used to decide branches while splitting</li> </ul>	

TASK	DATA TYPE		OBJECTIVE FUNCTION	HYPERPARAMETER	ALGORITHM CHARACTERISTICS			WHEN TO USE		
	Regression Classification	Mixed			Gini Impurity, Entropy for classification MSE for Regression	Max Depth, Min Samples Leaf, Min Samples Split	TYPE			
ALGORITHM CHARACTERISTICS	OUTLIER	MULTICOLLINEARITY	DATA TREATMENT	MODEL INTERPRETABILITY	REGULARIZATION TECHNIQUE					
No effect due to splitting of data	Robust	Not Required	OUTLIER REMOVAL	MISSING VALUE TREATMENT:	FEATURE SCALING	decision tree plots can be used to interpret model	Non-Parametric	Low	High	<ul style="list-style-type: none"> <li>if data has lots of categorical variables</li> <li>Large amount of training data points</li> </ul>

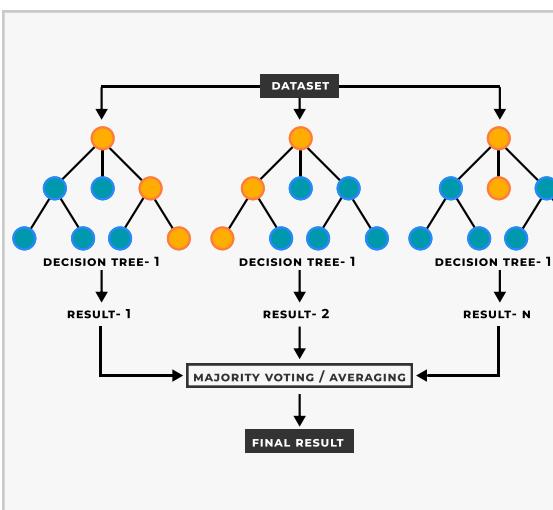
## 04. KNN



DESCRIPTION	
<ul style="list-style-type: none"> <li>Works by finding the K nearest neighbors to the query data point and using their value(class label / regression value) to make a prediction for the target variable.</li> <li>A larger K value will result in a smoother decision boundary, while a smaller K value will result in a more complex boundary.</li> </ul>	

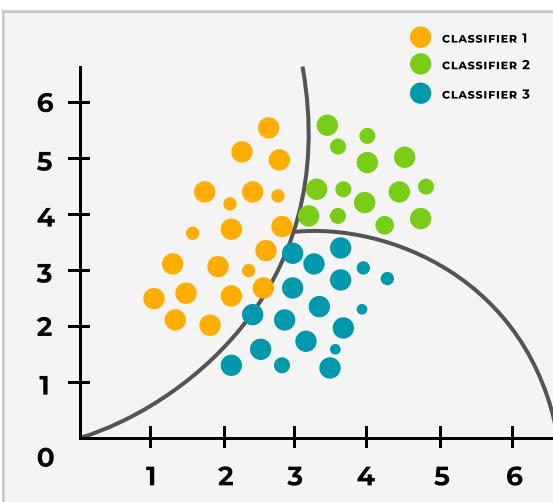
TASK	DATA TYPE		OBJECTIVE FUNCTION	HYPERPARAMETER	ALGORITHM CHARACTERISTICS			WHEN TO USE		
	Regression Classification	Numerical, need to change categorical variables to dummy variables			N/A	distance metric: K (# neighbours), - Euclidian, - Manhattan, or Minkowski	TYPE			
ALGORITHM CHARACTERISTICS	OUTLIER	MULTICOLLINEARITY	DATA TREATMENT	MODEL INTERPRETABILITY	REGULARIZATION TECHNIQUE					
Less impact if k is big enough	Sensitive	Required	OUTLIER REMOVAL	MISSING VALUE TREATMENT:	FEATURE SCALING	feature importance can not be interpreted	Non-Parametric	low (if k is small) high (if k is large)	High(k is small) Low (k is large)	<ul style="list-style-type: none"> <li>small dataset with low dimensions</li> <li>where latency is not a concern</li> </ul>

## 05. RANDOM FOREST



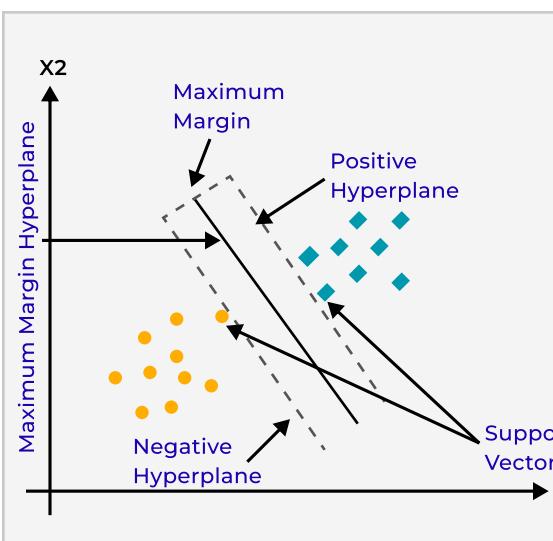
DESCRIPTION	TASK	DATA TYPE	OBJECTIVE FUNCTION	HYPERPARAMETER	ALGORITHM CHARACTERISTICS			WHEN TO USE		
<ul style="list-style-type: none"> <li>Ensemble method that combines the output of multiple decision trees each of which is trained on randomly sampled data with repetition.</li> </ul>	Regression Classification	Mixed	<ul style="list-style-type: none"> <li>Gini Impurity, Entropy for classification</li> <li>MSE for Regression</li> </ul>	number of trees, column Sample, row sample size, Depth of base learners	Parametric	BIAS High <small>(bootstrapping and restricted variable participation during split increase bias)</small>	VARIANCE Low	<ul style="list-style-type: none"> <li>large dataset</li> <li>When simple models does not produce desired results</li> </ul>		
<b>ALGORITHM CHARACTERISTICS</b>					<b>DATA TREATMENT</b>	<b>MODEL INTERPRETABILITY</b>				
OUTLIER No Impact as it get averaged out due to aggregation	MULTICOLLINEARITY Robust due to column sampling	OUTLIER REMOVAL Not Required	MISSING VALUE TREATMENT: Not Required	FEATURE SCALING Not Required	Interpretable Feature importance = weighted information gain of feature across base learners	<b>REGULARIZATION TECHNIQUE</b>		# base trees row and column sampling rate		
<b>DATA TREATMENT</b>					<b>MODEL INTERPRETABILITY</b>					
<b>REGULARIZATION TECHNIQUE</b>					# base trees row and column sampling rate					

## 06. NAÏVE BAYES



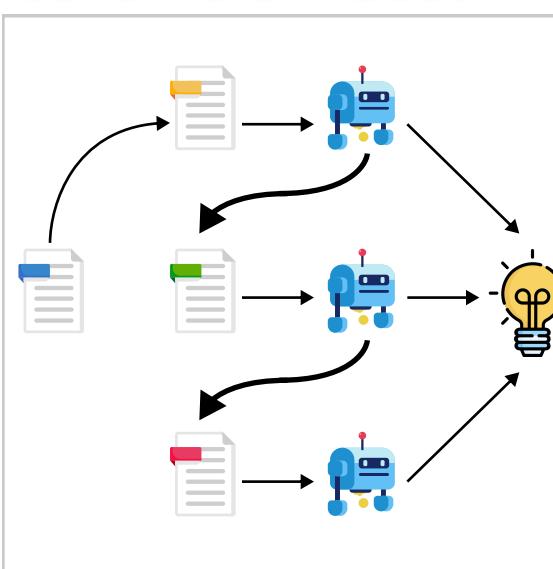
DESCRIPTION	TASK	DATA TYPE	OBJECTIVE FUNCTION	HYPERPARAMETER	ALGORITHM CHARACTERISTICS			WHEN TO USE
<ul style="list-style-type: none"> <li>Probabilistic machine learning algorithm based on Bayes' theorem</li> <li>Assumes that all the features are independent of each other.</li> </ul>	Classification	Numerical, need to change categorical variables to dummy variables	N/A	Laplace smoothing ( $\alpha$ )	Parametric	BIAS High	VARIANCE Low	<ul style="list-style-type: none"> <li>suitable for solving text classification problems</li> </ul>
<b>ALGORITHM CHARACTERISTICS</b>					<b>DATA TREATMENT</b>			
OUTLIER High impact in Gaussian NB Low impact in Multinomial/Bernoulli NB	MULTICOLLINEARITY Robust	OUTLIER REMOVAL Required	MISSING VALUE TREATMENT: Not Required	FEATURE SCALING Not Required	<b>MODEL INTERPRETABILITY</b>			
<b>DATA TREATMENT</b>					<b>MODEL INTERPRETABILITY</b>			
<b>REGULARIZATION TECHNIQUE</b>					# Laplace Smoothing			

## 07. SVM



DESCRIPTION	TASK	DATA TYPE	OBJECTIVE FUNCTION	HYPERPARAMETER	ALGORITHM CHARACTERISTICS			WHEN TO USE
<ul style="list-style-type: none"> <li>Works by finding the hyperplane that best separates the data into classes, maximizing the margin between the classes.</li> <li>It can even handle non-linear data, by using a technique called kernel trick to map the data into a higher-dimensional space where a linear boundary can be found.</li> </ul>	Classification Regression	Numerical, need to change categorical variables to dummy variables	Hinge Loss	C gamma kernel	Parametric	BIAS Low	VARIANCE High	<ul style="list-style-type: none"> <li>less data, and less features - kernel SVM</li> <li>less data, more features - linear SVM</li> </ul>
<b>ALGORITHM CHARACTERISTICS</b>					<b>DATA TREATMENT</b>			
OUTLIER Less sensitive, as separating hyperplane is decided by support vector	MULTICOLLINEARITY Sensitive	OUTLIER REMOVAL Not Required	MISSING VALUE TREATMENT: Required	FEATURE SCALING Required	<b>MODEL INTERPRETABILITY</b>			
<b>DATA TREATMENT</b>					<b>MODEL INTERPRETABILITY</b>			
<b>REGULARIZATION TECHNIQUE</b>					# L2 regularization with C parameter			

## 08. GRADIENT BOOSTING



DESCRIPTION	TASK	DATA TYPE	OBJECTIVE FUNCTION	HYPERPARAMETER	ALGORITHM CHARACTERISTICS			WHEN TO USE
<ul style="list-style-type: none"> <li>It uses multiple weak models (DT's) in sequential manner where each model tries to predict the error left over by the previous model. GBDT uses gradient descent to optimize the loss function and find the best combination of trees</li> </ul>	Classification Regression	Mixed	<ul style="list-style-type: none"> <li>Gini Impurity, Entropy for classification</li> <li>MSE for Regression</li> </ul>	<ul style="list-style-type: none"> <li>number of trees</li> <li>Learning Rate</li> <li>Regularization Parameters</li> </ul>	Non-Parametric	BIAS Low	VARIANCE High	<ul style="list-style-type: none"> <li>High dimensional and large datasets</li> </ul>
<b>ALGORITHM CHARACTERISTICS</b>					<b>DATA TREATMENT</b>			
OUTLIER Sensitive, because it builds each tree on previous trees residuals. Outliers will have much larger residuals than non-outliers	MULTICOLLINEARITY Sensitive	OUTLIER REMOVAL Not Required	MISSING VALUE TREATMENT: Required	FEATURE SCALING Required	<b>MODEL INTERPRETABILITY</b>			
<b>DATA TREATMENT</b>					<b>MODEL INTERPRETABILITY</b>			
<b>REGULARIZATION TECHNIQUE</b>					# base trees • row and column sampling rate			

## BAGGING

ensemble methods use predictive power of multiple trees(weak learners) instead of 1.  
But....

Follows parallel learning  
i.e base learners are formed independently

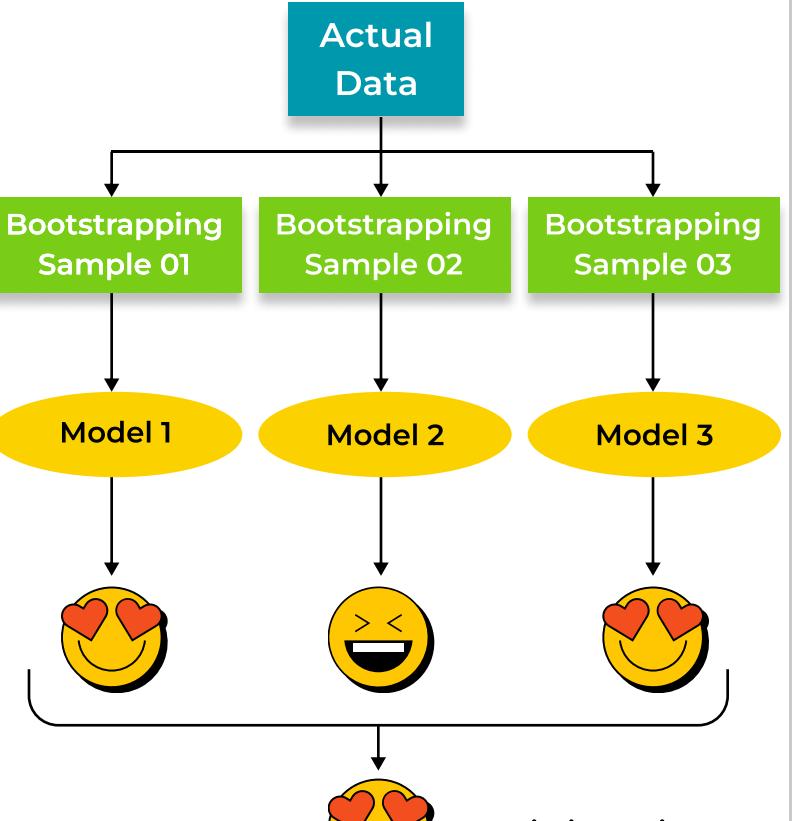
Random sampling with replacement

Both give final prediction by averaging N learners But...

..equal weight is given to all learners

..it reduced variance, and helps with overfitting

## BAGGING ENSEMBLE METHOD



## BOOSTING

## BOOSTING

Follows sequential learning  
i.e base learners are dependent on previous weak base learner

Random sampling with replacement over weighted data

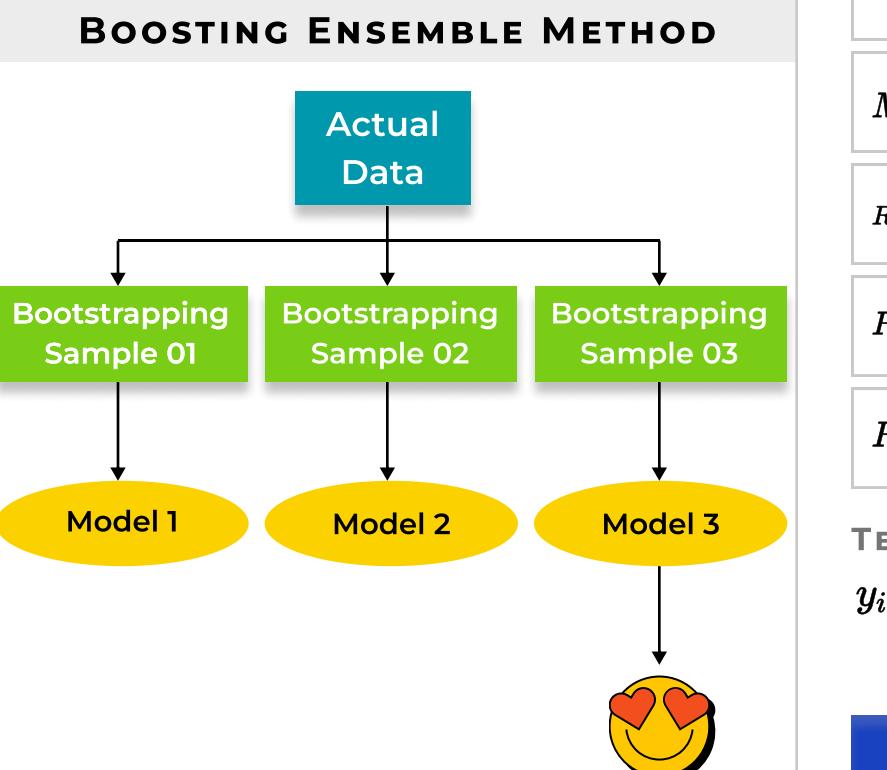
Both give final prediction by averaging N learners But...

..more weight is given to learners with better performance ( weighted average )

Both provide good scalability But ...

..it reduces bias but more prone to overfitting which can be avoided by tuning parameters

## BOOSTING ENSEMBLE METHOD



BUILD PARALLEL

BUILD SEQUENTIALLY

## REGULARIZATION

## L1 REGULARIZATION

How does it penalizes cost function ?

adds sum of absolute values of weights

adds sum of squared values of weights

## What's the formulation ?

Cost function = Loss +  $\lambda \sum |w|$

Cost function = Loss +  $\lambda \sum w^2$

## How does it impact weight coef?

Produces sparse solution i.e. non important feature weights become 0. (can be used for feature selection)

Produces non sparse solution i.e. reduces the value of non important feature weight but doesn't make them 0

## Is it impacted by outliers ?

Robust to outliers

Impacted by outliers as squared term is involved

## REGRESSION METRIC

$$MAE = \frac{1}{n} \sum_{i=1}^n |x_i - \hat{x}_i|$$

$$MSE = \frac{\sum (y_i - \hat{y}_i)^2}{n}$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n \left( \frac{d_i - f_i}{\sigma_i} \right)^2}$$

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

$$R^2_{adj} = 1 - (1 - R^2) \frac{n - 1}{n - p - 1}$$

## CLASSIFICATION METRIC

$$Accuracy = \frac{\text{Correct Predictions}}{\text{Total number of predictions}}$$

$$Precision = \frac{TP}{TP + FP} = \frac{\text{True Positive}}{\text{Predicted Positive}}$$

$$Recall = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} = \frac{\text{True Positive}}{\text{Actual Positives}}$$

$$Specificity = \frac{\text{True Negative}}{\text{True negative} + \text{False Positive}}$$

$$F1 - score = \frac{2 \cdot Precision \cdot Recall}{\frac{1}{Precision} + \frac{1}{Recall}} = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$

## TERMINOLOGIES:

$y_i \rightarrow$  True value ,  $\hat{y}_i \rightarrow$  Predicted value ,  $n \rightarrow$  number of data points

$\bar{y}_i \rightarrow$  mean value of all the data points ,  $k \rightarrow$  number of independent variables

## PARAMETRIC MODELS

Fixed number of parameters to build the model.

Considers strong assumptions about underlying distribution of data

Requires lesser data

## NON - PARAMETRIC MODELS

Not Fixed i.e flexible

No or fewer assumptions

requires more data

## BIAS VARIANCE TRADEOFF

$$\text{Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

## ANALOGY - SHOOTING EXAMPLE

High bias - aiming at wrong place

High variance - unsteady aim

Low bias, low variance: Aiming at the target and hitting it with good precision.



Low bias, high variance: Aiming at the target, but not hitting it consistently.



High bias, low variance: Aiming off the target, but being consistent.



High bias, high variance: Aiming off the target and being inconsistent.



## High Variance Model ?

Overfitting i.e. High error on test dataset (performing well on train dataset)

## High bias model means?

Underfitting i.e. High error on train as well as test data

