# Data Science and Artificial Intelligence

## Machine Learning

**Bias and Variance**

**Lecture No. 3**

By- SIDDHARTH SABHARWAL SIR

# Recap of Previous Lecture

Topic — Feature selection Technique

Topic — Ensemble method → Bagging
→ Boosting

Topic

Topic

Topic

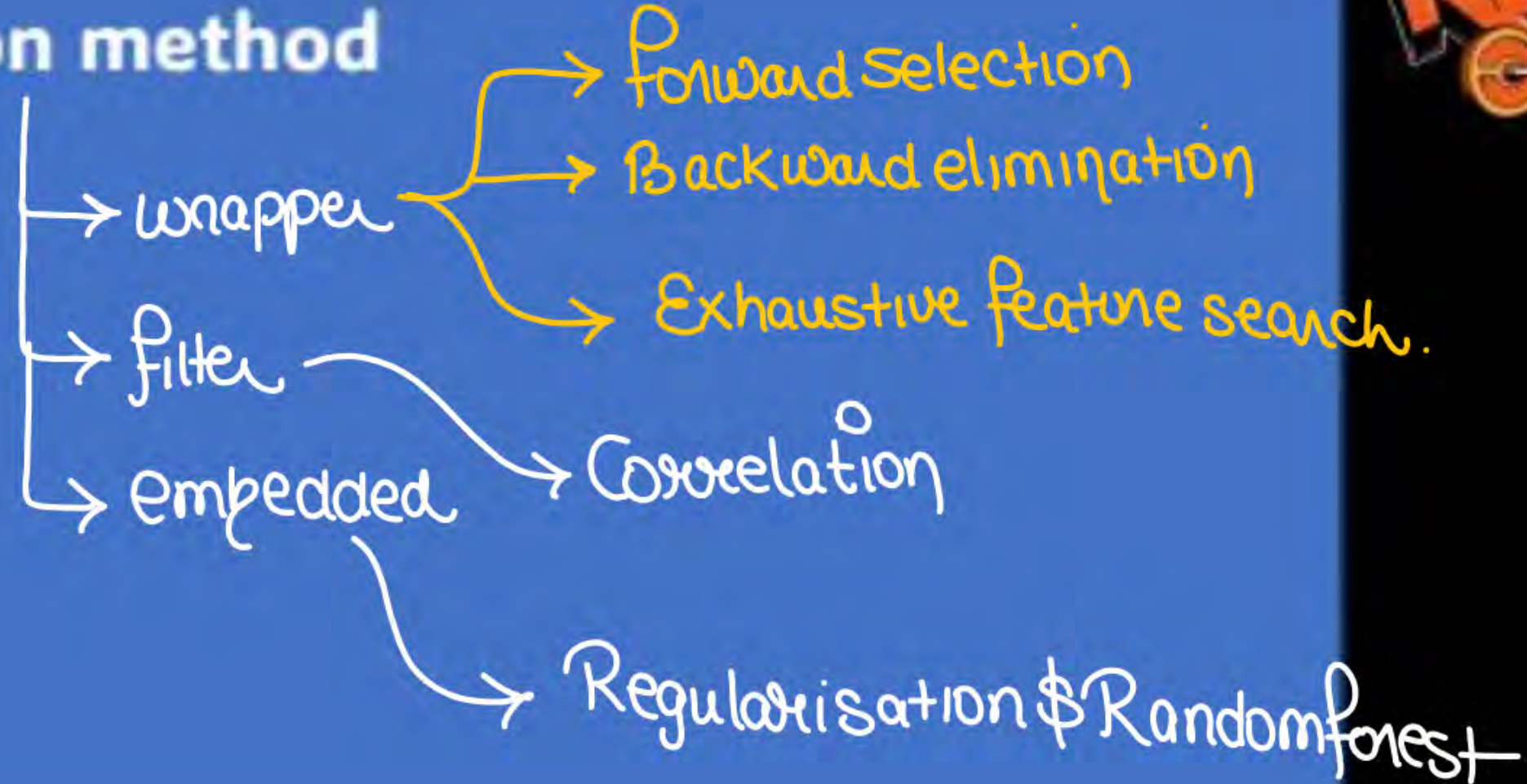# Topics to be Covered

Topic — Random Forest

Topic — MLE/MAP

Topic

Topic

Topic

"Strength and growth come only through continuous effort and struggle."

- NAPOLEON HILL -

**Feature selection method**

Wrapper →
- Forward selection
- Backward elimination
- Exhaustive feature search.

Filter → Correlation

Embedded → Regularisation $ Random forest
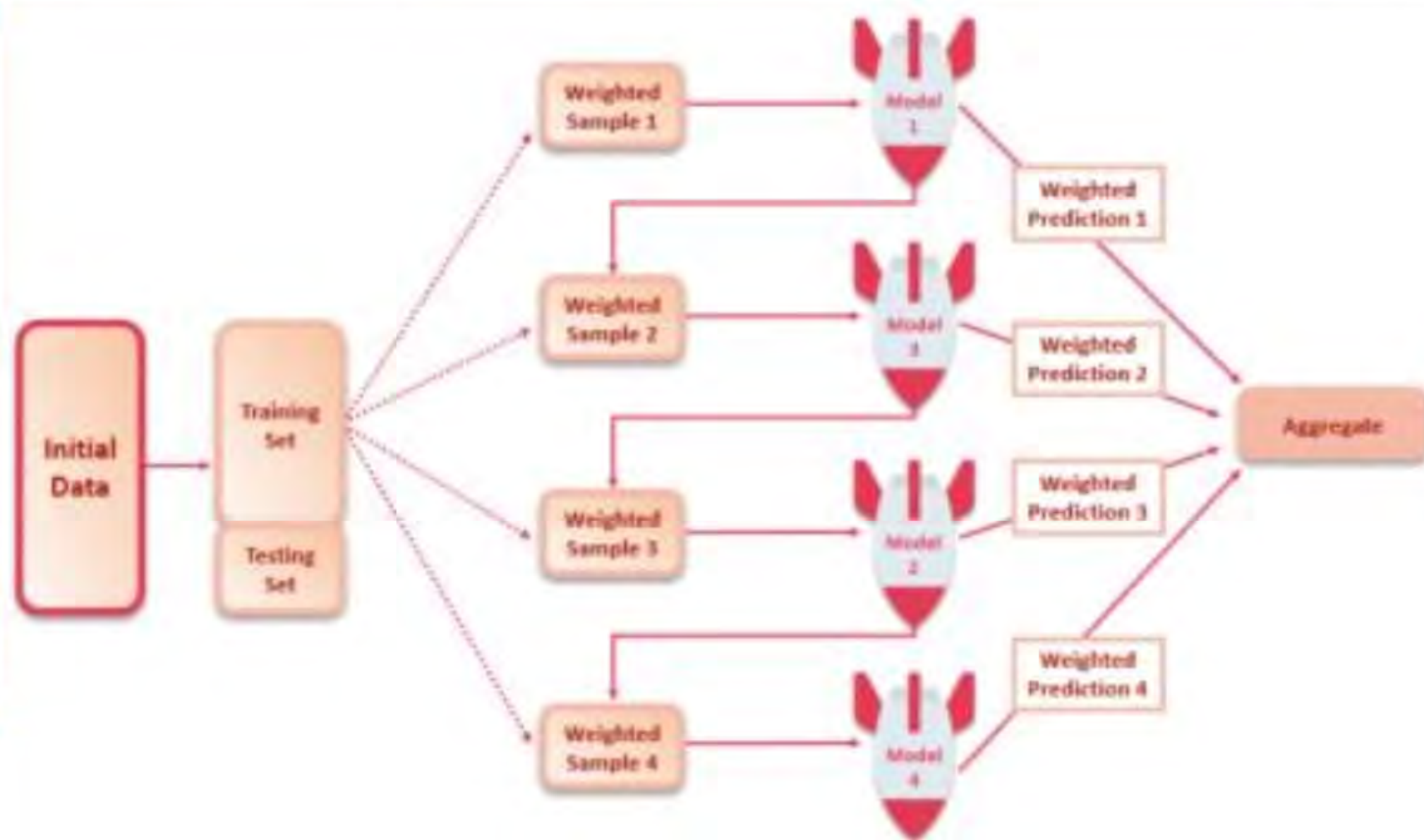
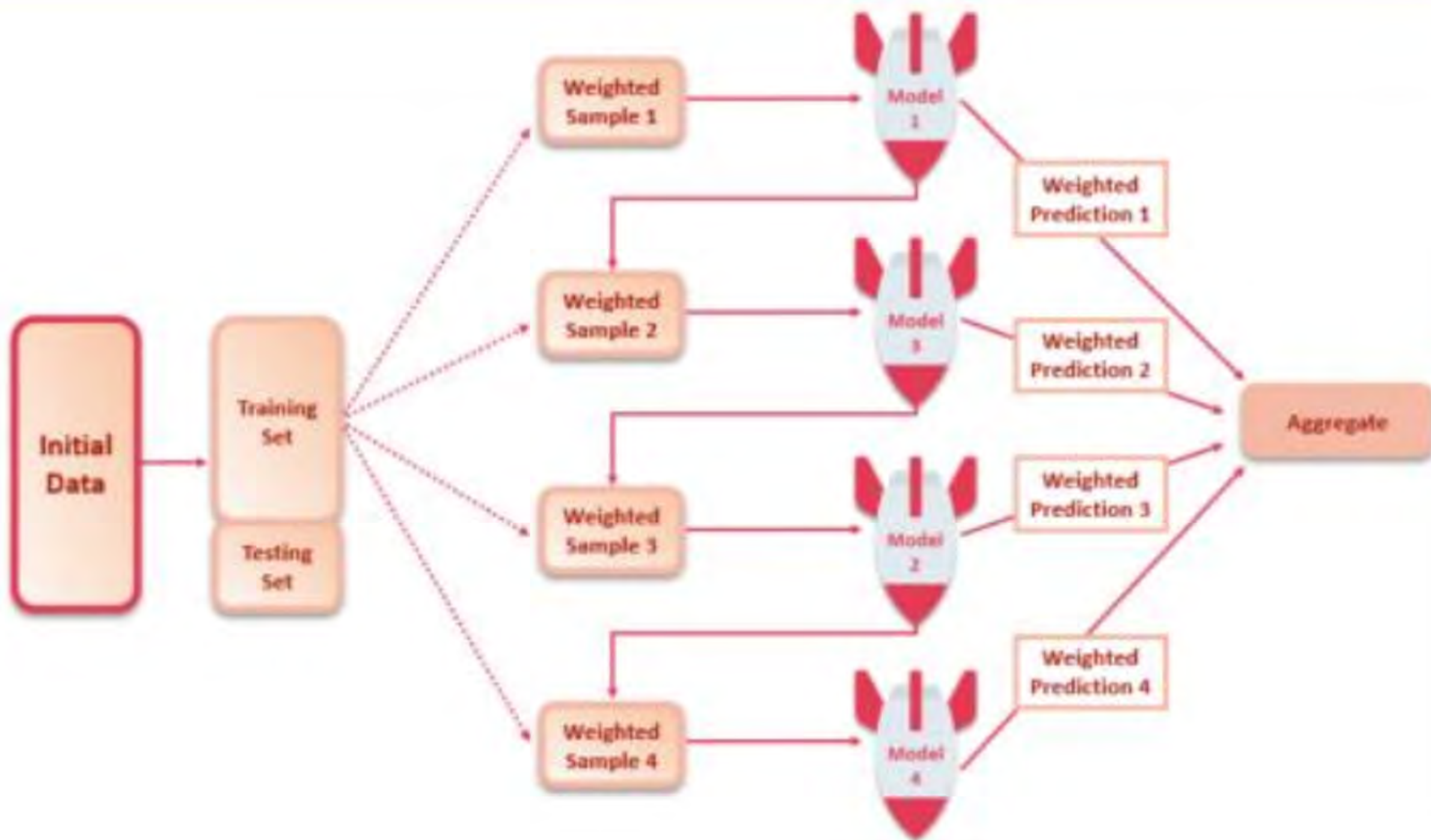Bagging (Parallel Run)

Boosting → Sequential

# Ensemble learning

❖ **Types of Ensemble Classifier – Boosting:**

❖ This is like Bagging.

❖ But this is not a parallel process rather a sequential process...

❖ Here we first learn a model and find the error on the data and then train next model where we have more error...

# Ensemble learning - Boosting

1. Samples generated from the training set are assigned the **same weight** to start with. These samples are used to train a homogeneous weak learner or base model.
2. The prediction error for a sample is calculated – **the greater the error, the weight of the sample increases**. Hence, the sample becomes more important for training the next base model.
3. The individual learner is weighted too – **does well on its predictions, gets a higher weight assigned to it**. So, a model that outputs good predictions will have a higher say in the final decision.
4. The weighted data is then passed on to the following base model, and steps 2) and 3) are repeated until the **data is fitted well enough to reduce the error below a certain threshold**.
5. When new data is fed into the boosting model, it is passed through all individual base models, and **each model makes its own weighted prediction**.
6. Weight of these models is used to generate the final prediction. The predictions are scaled and **aggregated to produce a final prediction**.

$$\cdot \Rightarrow \text{So} \left( \hat{y} = \hat{y_1} w_1 + \hat{y_2} w_2 + \hat{y_3} w_3 + - - - \right)$$

Both class/Reg $\Rightarrow$

Algo 1   $w_1 = \cdot 1 \rightsquigarrow$ ✓ Class 1

2   $w_2 = \cdot 4 \rightsquigarrow$ ✓ Class 0

3   $w_3 = \cdot 2 \rightsquigarrow$ ✓ Class 0

4   $w_4 = \cdot 7 \rightsquigarrow$ ✓ Class 0

5   $w_5 = \cdot 3 \rightsquigarrow$ Class 1 ✓

6   $w_6 = \cdot 4 \rightsquigarrow$ Class 1 ✓

$$\hat{y} = \left( \hat{y_1} w_1 + w_2 \hat{y_2} + - - - \right)$$

$$\hat{y} = \text{Class 1}(\cdot 8) + \text{Class 0}(1 \cdot 3)$$

$\Rightarrow$ assign class with larger coeff

$$\bullet \Rightarrow \text{So} \left( \hat{y} = \hat{y}_1 \omega_1 + \hat{y}_2 \omega_2 + \hat{y}_3 \omega_3 + - - - \right)$$

Both class/Reg $\Rightarrow$

Algo 1   $\omega_1 = \cdot 1 \rightsquigarrow$ Class 1 ✓

2   $\omega_2 = \cdot 4 \rightsquigarrow$ Class 0

3   $\omega_3 = \cdot 2 \rightsquigarrow$ Class 1 ✓

4   $\omega_4 = \cdot 7 \rightsquigarrow$ Class 0

5   $\omega_5 = \cdot 3 \rightsquigarrow$ Class 1 ✓

6   $\omega_6 = \cdot 4 \rightsquigarrow$ Class 1 ✓

$$\hat{y} = \left( \hat{y}_1 \omega_1 + \omega_2 \hat{y}_2 + - - - \right)$$

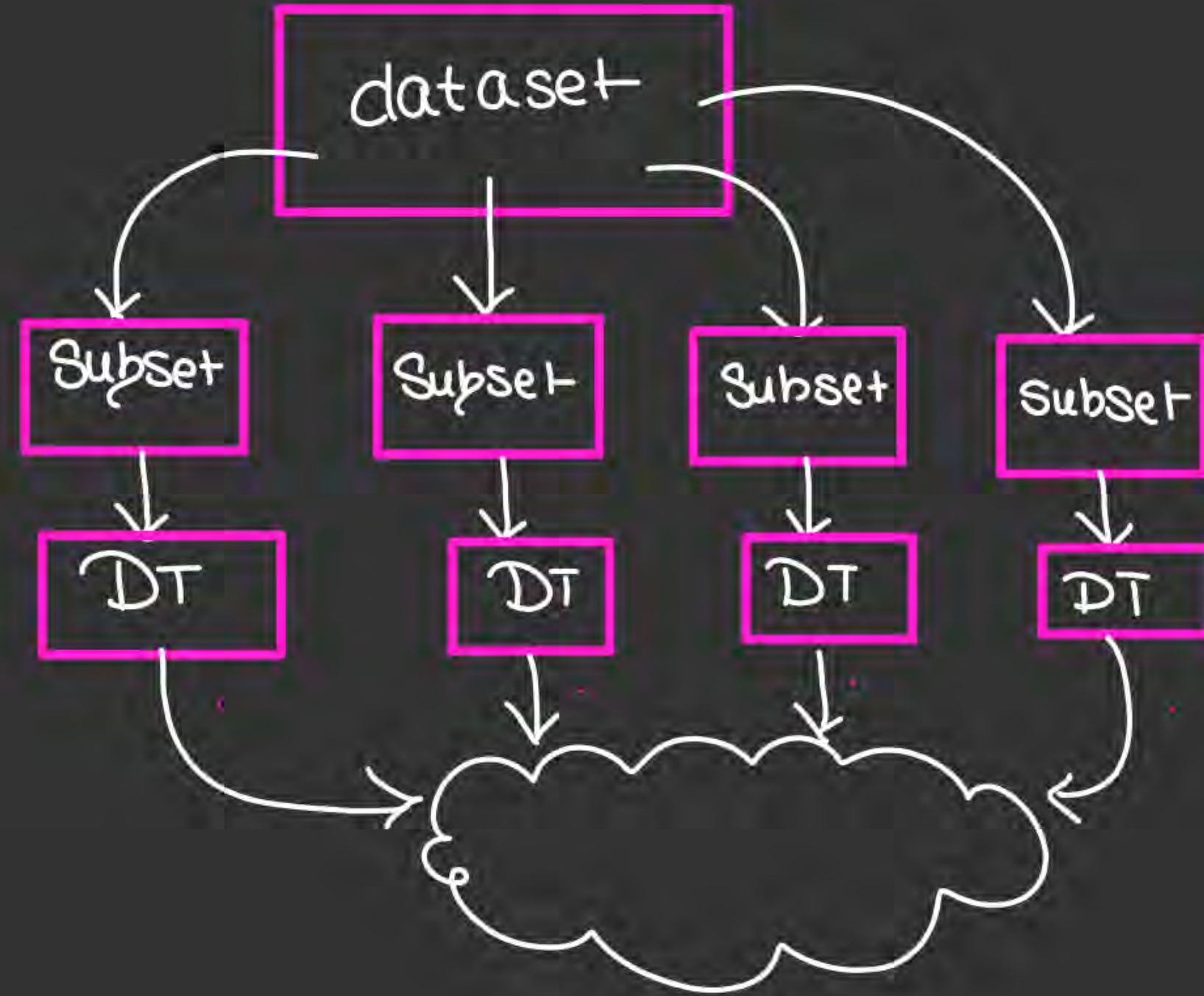$$\hat{y} = 1 \cdot 0 \text{ Class 1} + (1 \cdot 1) \text{ Class 0}$$

$\downarrow$

Class 0 assign.

⇒ Random forest ⇒
- If we have large dataset then, DT become computationally extensive and DT has property of overfitting
- So if we use ensemble learning technique in DT then we get Random forest.

→ *Power of feature Selection.*

→ *Computationally extensive*

**Steps Involved in Random Forest Algorithm**

- **Step 1:** In the Random forest model, a subset of data points and a subset of features is selected for constructing each decision tree. Simply put, n random records and m features are taken from the data set having k number of records.

- **Step 2:** Individual decision trees are constructed for each sample.

- **Step 3:** Each decision tree will generate an output.

- **Step 4:** Final output is considered based on *Majority Voting or Averaging* for Classification and regression, respectively.

Random Forest ⇒ Hyperparameters ⇒

1. No of subsets
2. No of dimension in subsets
3. Max depth of DT
4. Gini Impurity / IG / entropy.

**Important Hyperparameters in Random Forest**

- Hyperparameters are used in random forests to either enhance the performance and predictive power of models or to make the model faster.
- Number of decision trees to be constructed
- Maximum number of features a tree can use
- Splitting thresholds

Random forest $\rightarrow$

1) Reduce Risk of overfit
2) feature selection
3) low bias low Variance
4)

*no individual tree has whole data*

## Key Benefits

❏ Reduced risk of overfitting: when there's a robust number of decision trees in a random forest, the classifier won't overfit the model since the averaging of uncorrelated trees lowers the overall variance and prediction error.

❏ Provides flexibility: Since random forest can handle both regression and classification tasks with a high degree of accuracy.

❏ Easy to determine feature importance: Random forest makes it easy to evaluate variable importance, or contribution, to the model. There are a few ways to evaluate feature importance. Gini importance and decrease in impurity are usually used to measure how much the model's accuracy decreases when a given variable is excluded.

Problems ⟹

1) Computationally extensive

2) more memory

3) the subsets shd be good enough such that DT donot underfit

## Key Challenges

❑ Time-consuming process: Since random forest algorithms can handle large data sets, they can be provide more accurate predictions, but can be slow to process data as they are computing data for each individual decision tree.

❑ Requires more resources: Since random forests process larger data sets, they'll require more resources to store that data.

❑ More complex: The prediction of a single decision tree is easier to interpret when compared to a forest of them.

| Decision trees | Random Forest |
|---|---|
| 1. Decision trees normally suffer from the problem of overfitting if it's allowed to grow without any control. | 1. Random forests are created from subsets of data, and the final output is based on average or majority ranking; hence the problem of overfitting is taken care of. |
| 2. A single decision tree is faster in computation. *(Train & Test Fast)* | 2. It is comparatively slower. *(Train & Test both slow)* |
| 3. When a data set with features is taken as input by a decision tree, it will formulate some rules to make predictions. | 3. Random forest randomly selects observations, builds a decision tree, and takes the average result. It doesn't use any set of formulas. |

# Random Forest Algorithm

| Feature | Random Forest | Other ML Algorithms |
|---|---|---|
| Ensemble Approach | Utilizes an ensemble of decision trees, combining their outputs for predictions, fostering robustness and accuracy. | Typically relies on a single model (e.g., linear regression, support vector machine) without the ensemble approach, potentially leading to less resilience against noise. |
| Overfitting Resistance | Resistant to overfitting due to the aggregation of diverse decision trees, preventing memorization of training data. | Some algorithms may be prone to overfitting, especially when dealing with complex datasets, as they may excessively adapt to training noise. |
| Handling of Missing Data | Exhibits resilience in handling missing values by leveraging available features for predictions, contributing to practicality in real-world scenarios. | Other algorithms may require imputation or elimination of missing data, potentially impacting model training and performance. |
| Variable Importance | Provides a built-in mechanism for assessing variable importance, aiding in feature selection and interpretation of influential factors. | Many algorithms may lack an explicit feature importance assessment, making it challenging to identify crucial variables for predictions. |
| Parallelization Potential | Capitalizes on parallelization, enabling the simultaneous training of decision trees, resulting in faster computation for large datasets. | Some algorithms may have limited parallelization capabilities, potentially leading to longer training times for extensive datasets. |

# Maximum likelihood Estimation

## What is MLE (lets see an example)

Maximum Likelihood Estimation (MLE) is a statistical method used to estimate the parameters of a probability distribution that best describe a given dataset. The fundamental idea behind MLE is to find the values of the parameters that maximize the likelihood of the observed data, assuming that the data are generated by the specified distribution.

## What is MLE (lets see an example)

- So whenever analysis is done we take a small sample from data

- Now our task is to use this sample data to predict the parameters of the Probability distribution of whole data

election data    100 Crore

## What is MLE (lets see an example)

- data distribution $\longrightarrow$ data का PDF

  data is a RV
  and data distribution
  $\Rightarrow$ PDF of data

Gaussian

- $f_X(x) = \dfrac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$

  $\longrightarrow (\sigma^2, \mu)$

- Exponential dist
  $f_X(x) = \lambda e^{-\lambda x} \longrightarrow (\lambda)$

## What is MLE (lets see an example)

$ex \Rightarrow$ let data has gaussian distribution

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

- If we have a sample of data $\Rightarrow$
  Containining points of value $\begin{bmatrix} \text{Single D} \\ \text{data} \end{bmatrix}$
  $(x_1, x_2, x_3, x_4, x_5 - - - - -)$

- These points in sample of data are independent from each other.

Probability of getting these samples $\Rightarrow$

$$P\left[x_1, x_2, x_3, x_4 - - - - x_N / \mu, \sigma\right] \Rightarrow P(x_1) \, P(x_2) \, P(x_3) \, P(x_4) - - - P(x_N)$$

$$\left\{ \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_1-\mu)^2/2\sigma^2} \cdot \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_2-\mu)^2/2\sigma^2} \quad - - - - \right\}$$

data

$$P\left(x_1, x_2 - - - x_N / \mu, \sigma\right)$$

✓

Sample

$$= P(x_1 / \mu, \sigma) \, P(x_2 / \mu\sigma) - - -$$

This is the probability to get
this sampleset

\* So we always assume that
this sample represents whole
data

\* thus we want that the probability
of getting this sample is maximum.

So* we assume that if we sample the data the Probability of getting same sample values is maximum

⇒ the parameters of the distribution of data shape such that it maximises the probability of sample

It maximises the likelihood of sample.

So we have to
maximise $P_{x_1} P_{x_2} P_{x_3} - - - P_{x_N}$ ⇐

Probab of
Sample
$P(x_1, x_2 - - x_N / \mu, \sigma)$

$\Rightarrow \left\{ \left( \dfrac{1}{\sqrt{2\pi\sigma^2}} \right)^N e^{-(x_1-\mu)^2/2\sigma^2} e^{-(x_2-\mu)^2/2\sigma^2} - - - e^{-(x_N-\mu)^2/2\sigma^2} \right\}$

likelihood of
Sample

data



assumption

* Sample represent whole data

* The Parameters of PDF of data
Shd be such that it max Probab of
Sample.

So $\mu$ and $\sigma^2$ of the distribution of data shape
Such that the Probab of getting Sample is maximixed

$$\text{Likelihood of Sample} \Rightarrow \mathcal{L} = \left\{ \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^N e^{-\sum_{i=1}^{N}(x_i-\mu)^2/2\sigma^2} \right\}\Bigg|_{max}$$

$$\left(\log \alpha\right)_{max} \Rightarrow \left\{ \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^N + \log e^{-\sum_{i=1}^{N}(x_i-u)^2/2\sigma^2} \right\}$$

$$\Rightarrow \frac{\partial}{\partial u}\left\{ N\log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^{N}(x_i-u)^2/2\sigma^2 \right\} = 0$$

$$-2\sum_{i=1}^{N}(x_i^o-u)/2\sigma^2 = 0$$

$$\sum_{i=1}^{N}x_i^o = Nu$$

$$u = \frac{1}{N}\sum_{i=1}^{N}x_i^o \checkmark$$

$$\frac{\partial}{\partial \sigma} \left\{ N \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{2\sigma^2} \right\} = 0$$

$$\frac{\partial}{\partial \sigma} \left\{ -N \log\left(\sqrt{2\pi}\,\sigma\right) - \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{2\sigma^2} \right\} = 0$$

$$\log x = \frac{1}{x}$$

$$\left(\frac{-N\sqrt{2\pi}}{\sqrt{2\pi}\sigma}\right) - \frac{1}{2}\frac{(-2)}{\sigma^3} \sum_{i=1}^{N}(x_i - \mu)^2 = 0$$

$$; \qquad \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{\sigma^2} = N$$

$$\sigma^2 = \frac{1}{N}\sum_{i=1}^{N}(x_i - \mu)^2$$

## What is MLE (Logistic Regression)

- Using Sample we need to find PDF of data
  → maximizing Probab/likelihood of Sample.
    → to predict parameter of the PDF of data (MLE)

## Probability Density Estimation & Maximum Likelihood Estimation

### So what is Probability Density Estimation

❖ **Probability Density: Assume a random variable x that has a probability distribution p(x). The relationship between the outcomes of a random variable and its probability is referred to as the probability density.**

❖ **The problem is that we don't always know the full probability distribution for a random variable. This is because we only use a small subset of observations to derive the outcome. This problem is referred to as Probability Density Estimation as we use only a random sample of observations to find the general density of the whole sample space.**

# Maximum likelihood Estimation

## Probability Density Estimation & Maximum Likelihood Estimation

So what is Probability Density Estimation

❖ **Density Estimation: It is the process of finding out the density of the whole population by examining a random sample of data from that population.**

## Probability Density Estimation & Maximum Likelihood Estimation

### Definition

❖ **Maximum Likelihood Estimation**

❖ our primary job is to analyse the data that we have been presented with.
❖ First thing would be to identify the distribution from which we have obtained our data.
❖ Next, we need to use our data to find the parameters of our distribution.
❖ Normal distributions, as we know, have mean ($\mu$) & variance ($\sigma2$)
❖ Binomial distributions have the n and p.
❖ Exponential distributions have the inverse mean ($\lambda$).

## What is Maximum Likelihood Estimation (MLE)

### Definition

- ❖ The goal of the maximum likelihood principle is to fit an optimal statistical distribution to some data.
- ❖ This makes the data easier to work with, makes it more general, allows us to see if new data follows the same distribution as the previous data, and lastly, it allows us to classify unlabelled data points.
- ❖ MLE stands for Maximum Likelihood Estimation. It is a method for estimating the parameters of a statistical model, given a dataset. The goal of MLE is to find the parameter values that maximize the likelihood of the data, given the model. This is done by choosing the values of the parameters that make the observed data most probable.

## What is Maximum Likelihood Estimation (MLE)

❖ we want to do now is obtain the parameter set θ that maximises the joint density function of the data vector; the so-called Likelihood function L(θ).

❖ This likelihood function can also be expressed as P(X|θ), which can be read as the conditional probability of X given the parameter set θ.

$$L(\theta) = p(X \mid \theta) = p(X(1), X(2), \ldots X(n) \mid \theta)$$

X is the data matrix, and X(1) up to X(n) are each of the data points, and θ is the given parameter set for the distribution.

## What is Maximum Likelihood Estimation (MLE)

❖ **To obtain this optimal parameter set, we take derivatives with respect to θ in the likelihood function and search for the maximum: this maximum represents the values of the parameters that make observing the available data as likely as possible.**

$$\frac{\partial}{\partial \theta} p(X|\theta) = 0$$

Taking derivatives with respect to θ

## What is Maximum Likelihood Estimation (MLE)

❖ if the data points of X are independent of each other, the likelihood function can be expressed as the product of the individual probabilities of each data point given the parameter set:

$$L(\theta) = p(X \mid \theta) = \prod p(X(j) \mid \theta)$$

Taking the derivatives with respect to this equation for each parameter (mean, variance, etc...) keeping the others constant, gives us the **relationship between the value of the data points, the number of data points, and each parameter.**

## What is Maximum Likelihood Estimation (MLE)

From the likelihood function we take log likelihood function

## What is Maximum Likelihood Estimation (MLE)

The goal of MLE is to infer $\Theta$ in the likelihood function $p(X|\Theta)$.

$$\theta_{MLE}$$
$$= arg\ max\ p(X|\theta)$$
$$= arg\ max\ \prod_i p(x_i|\theta)$$
$$= arg\ max\ log \prod_i p(x_i|\theta)$$
$$= arg\ max\ \sum_i log\ p(x_i|\theta)$$

Lets see some examples of MLE

# Maximum Aposterori Probability Rule

## What is Maximum Aposteriori Probability Rule(MAP)

Here we maximize ...

- MAP stands for Maximum A Posteriori probability. It is a method for estimating the parameters of a statistical model, given a dataset and some prior knowledge about the model. The goal of MAP is to find the parameter values that maximize the posterior probability of the data, given the model and the prior knowledge. This is done by choosing the values of the parameters that make the observed data most probable, given the prior knowledge.

## What is Maximum Aposteriori Probability Rule(MAP)

**Here we maximize ...**

- The goal of MAP is to find the parameter values that maximize the posterior probability of the data, given the model and the prior knowledge.
- MAP is similar to MLE (Maximum Likelihood Estimation), but it incorporates prior knowledge about the model into the estimation process. This can be useful in cases where the data is limited or noisy, or where there is a need to incorporate domain-specific knowledge into the model.

# Maximum Aposterori Probability Rule

## What is Maximum Aposteriori Probability Rule(MAP)

**Here we maximize ...**

- MAP is similar to MLE (Maximum Likelihood Estimation), but it incorporates prior knowledge about the model into the estimation process. This can be useful in cases where the data is limited or noisy, or where there is a need to incorporate domain-specific knowledge into the model.

$$p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)} \propto p(X|\theta)p(\theta)$$

$$
\begin{aligned}
\theta_{MAP} \\
&= arg\ max\ p(X|\theta)p(\theta) \\
&= arg\ max\ log[p(X|\theta)] + log(p(\theta)) \\
&= arg\ max\ log \prod_i p(x_i|\theta) + log(p(\theta)) \\
&= arg\ max \sum_i log\ p(x_i|\theta) + log(p(\theta))
\end{aligned}
$$

Comparing the equation of MAP with MLE, we can see that the only difference is that MAP includes prior in the formula, which means that the likelihood is weighted by the prior in MAP.

THANK - YOU