

# **Introduction to** **Machine Learning**

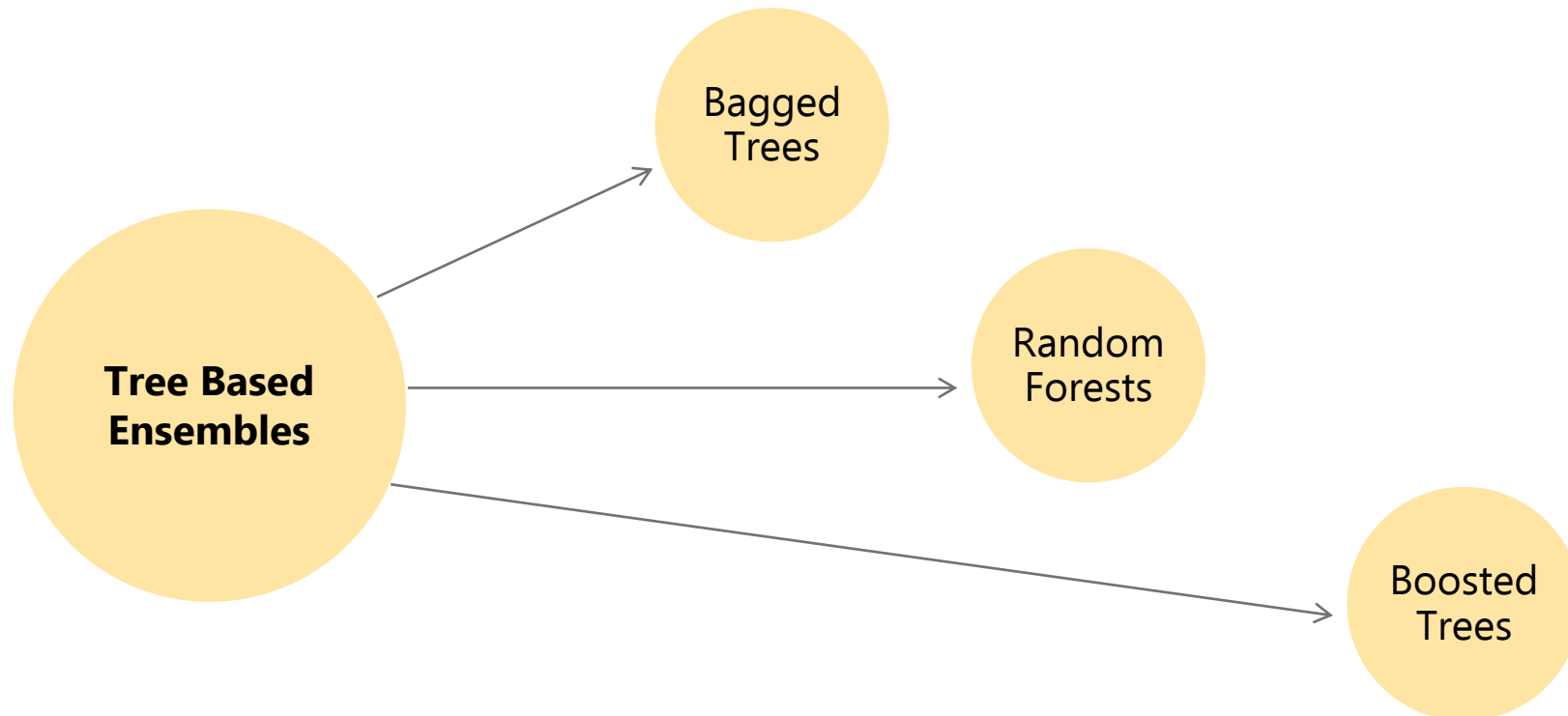


Class  
**Tree Based Model**

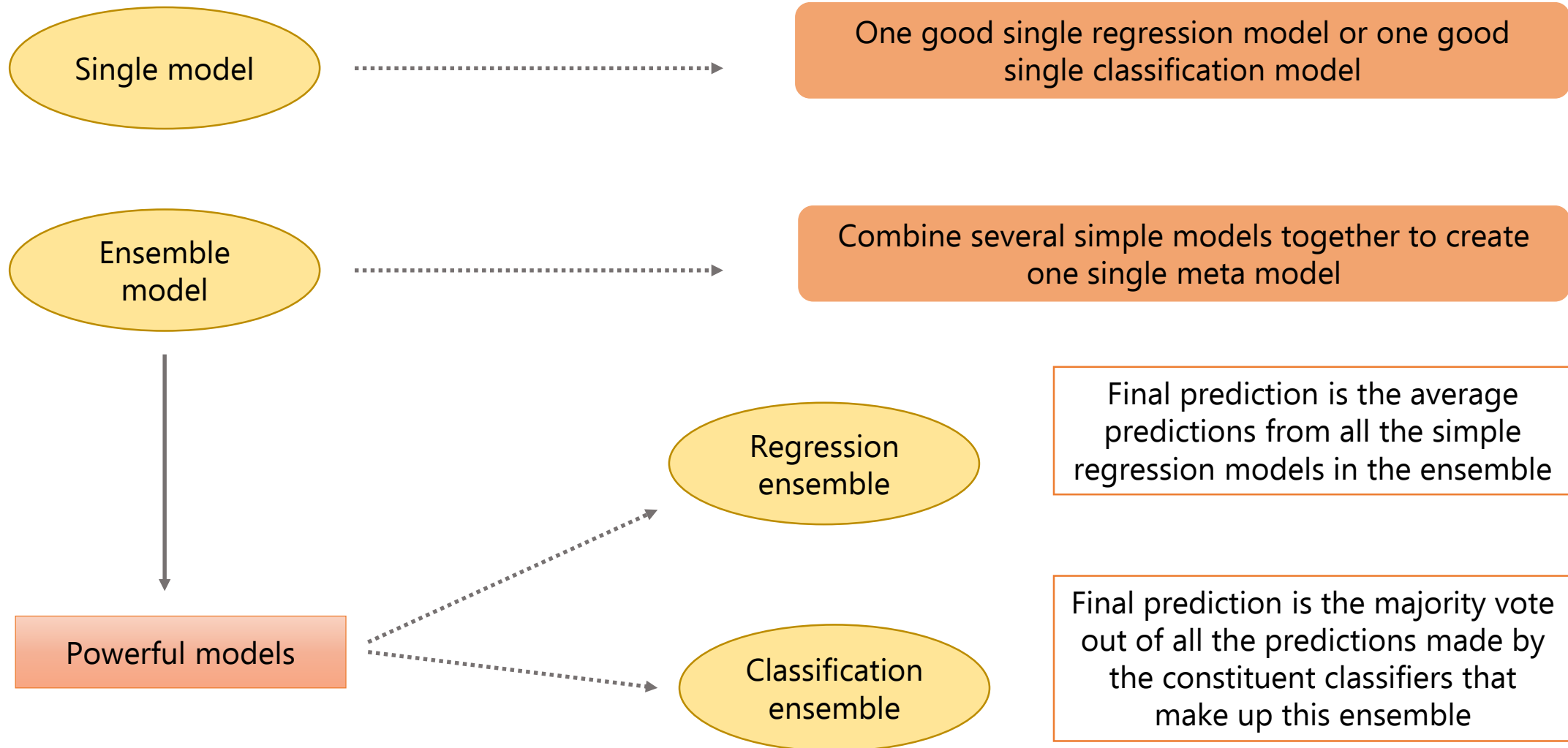


Topic  
**Tree Based Ensembles: Bagged Trees and  
Random Forests**

# Agenda

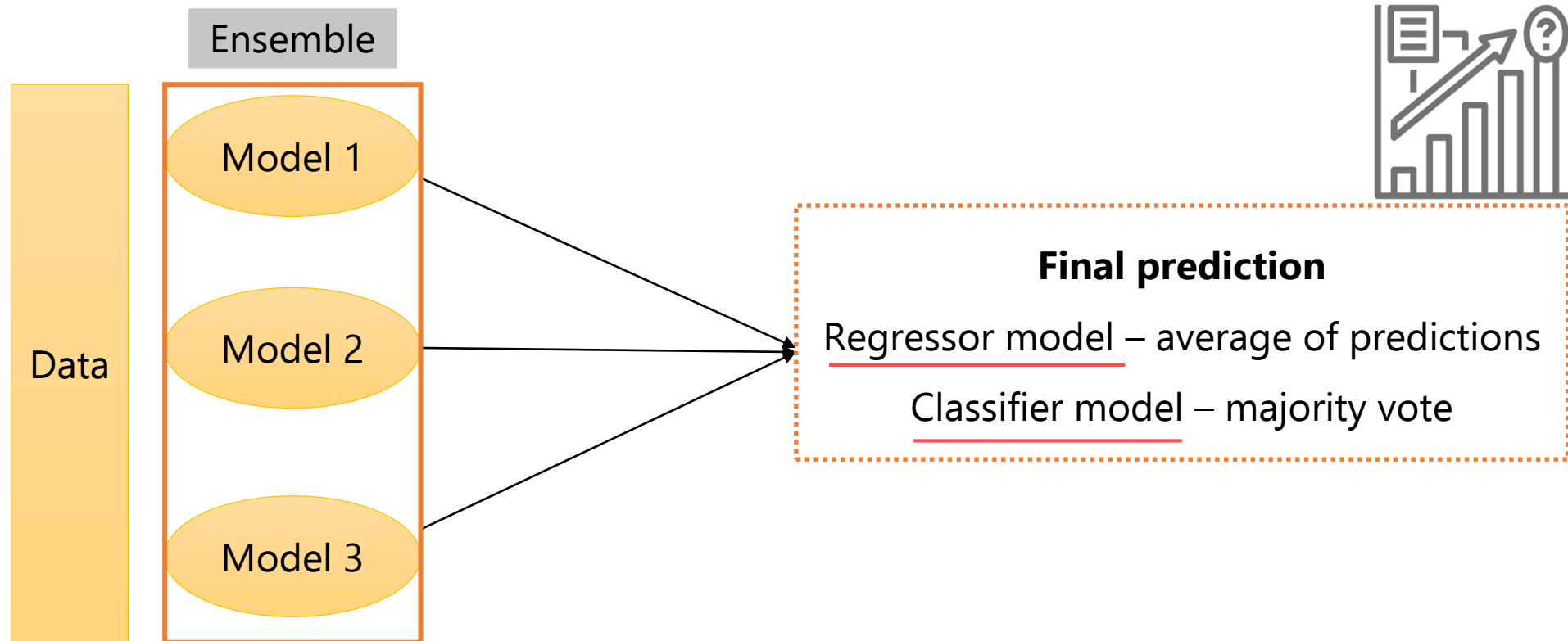


# Tree Based Ensembles Overview

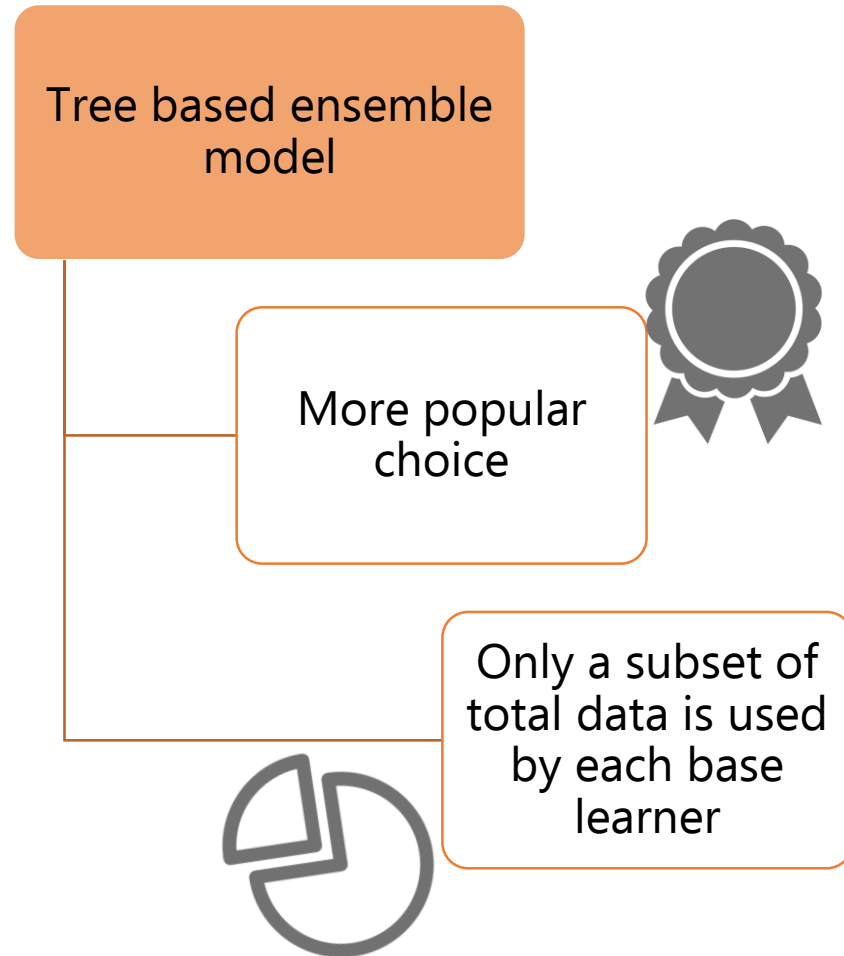


# Tree Based Ensembles Overview

Schematic working of ensembles



# Tree Based Ensembles Overview

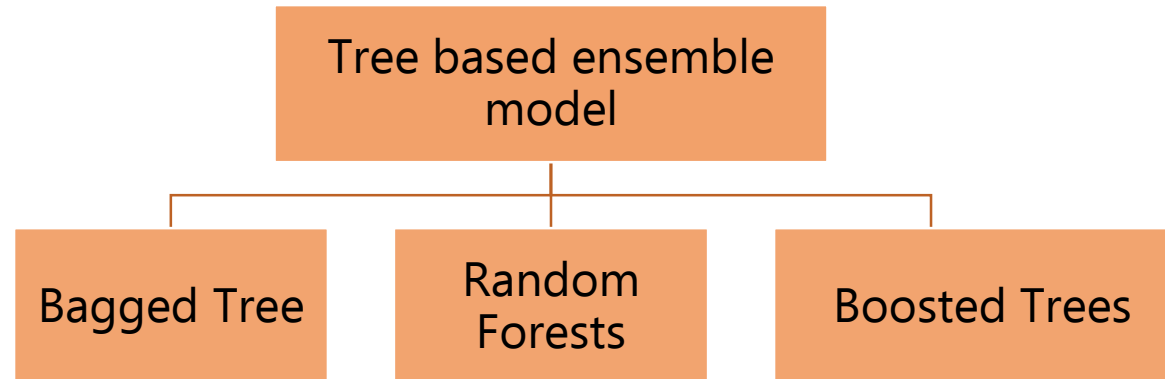


The way this data set is fed into each of the base learners is based on a **data sampling scheme**

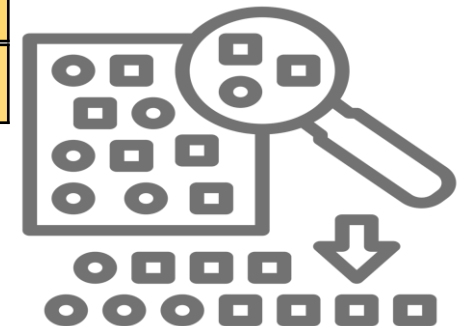
Different sampling schemes give rise to different types of tree based ensembles



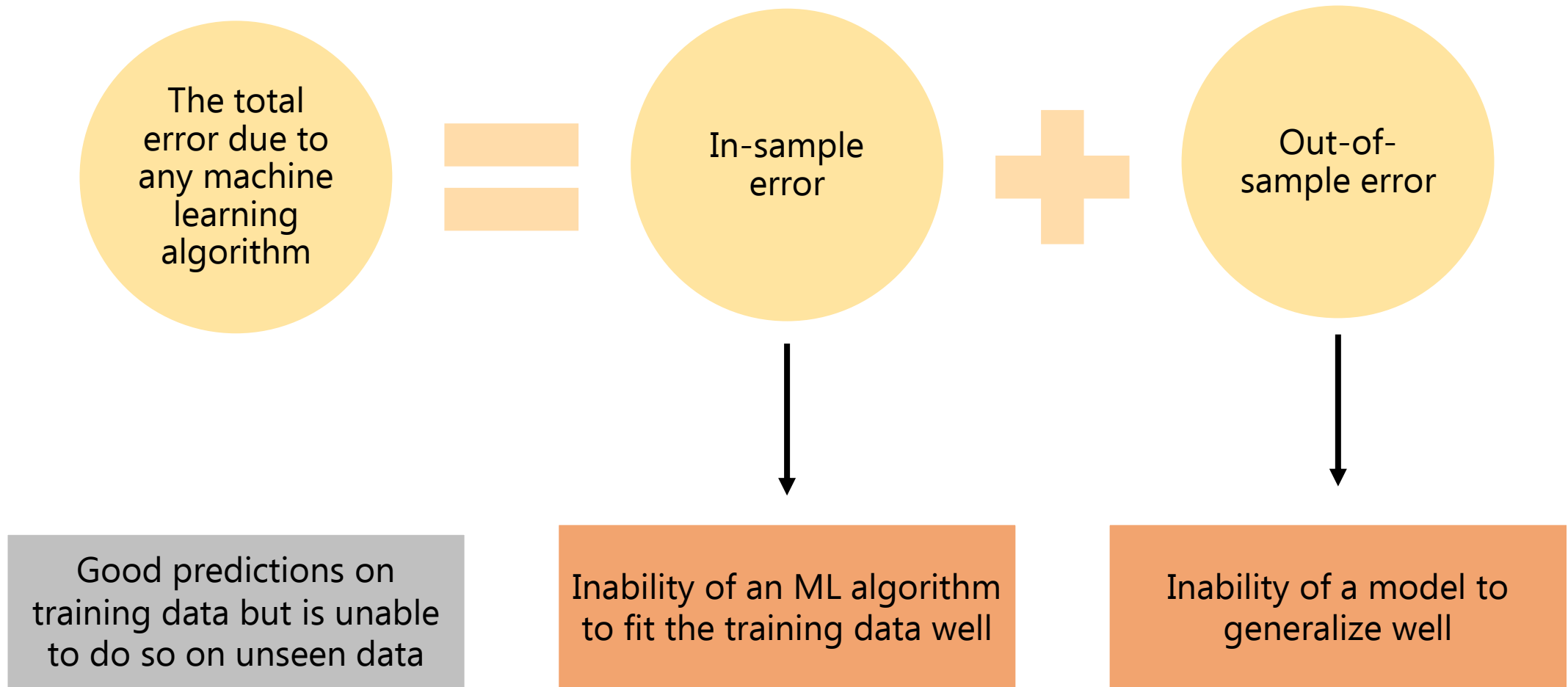
# Tree Based Ensemble Models



Sampling Scheme	Bootstrap Sampling	Bootstrap Sampling + Feature Sampling	Data Reweighing
Base Learner	Tree	Tree	Tree

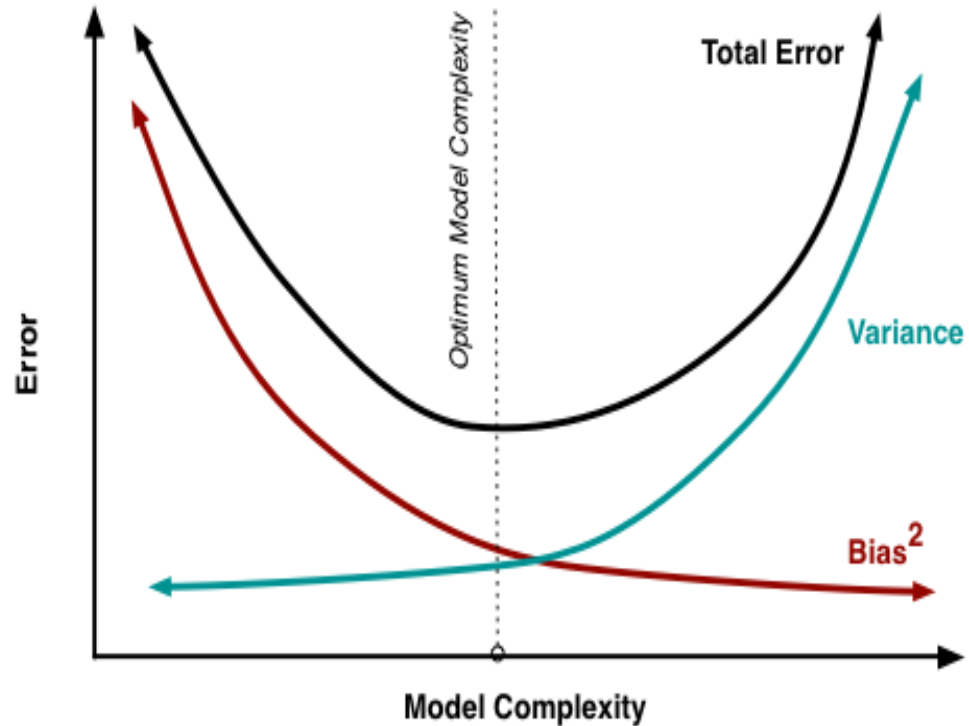


# Bagged Trees





# Bagged Trees



$$\text{Error} = \text{Bias} + \text{Variance} = \text{In-sample Error} + \text{Out-of-sample Error}$$

Trade offs between the model complexity and the error in models

More complicated models have very **low in-sample error but have a high out-of-sample error**

Simpler models have **low out-of-sample error but high in-sample error**

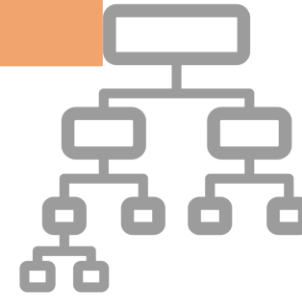
Theoretically, there is a limit to minimum error that can be achieved

Reduce error further by decreasing in-sample error and out-of-sample error simultaneously

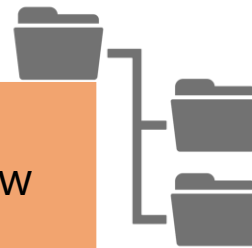


# Bagged Trees

Use tree based models as base learners to reduce in-sample error



While training a tree based ensemble the constituent tree models are allowed to grow **many levels deep**

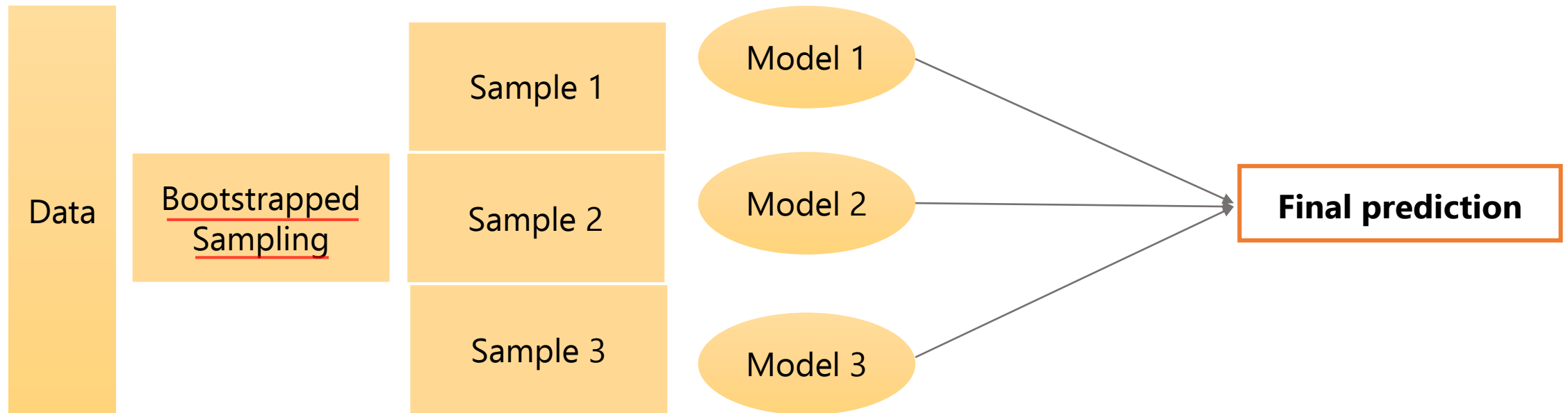


The intricacies of training data are captured intimately, thereby reducing the in-sample error



# Bagged Trees

In the case of **Bagged trees** each of the unpruned trees are fed bootstrapped samples of original data set



# Bootstrapped Sampling

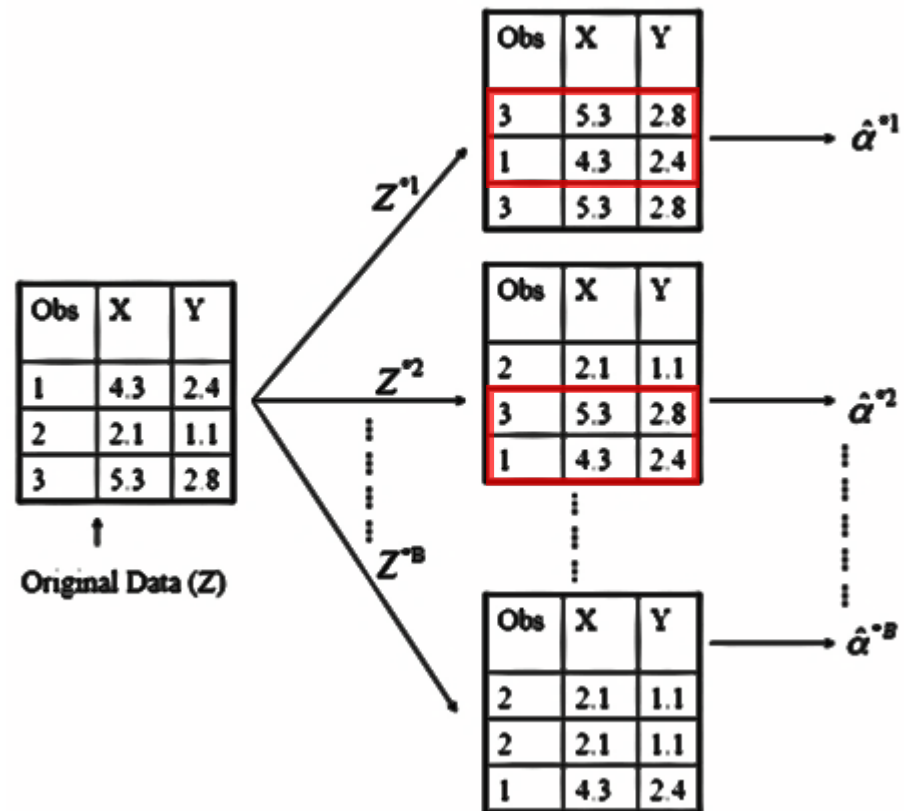
Bootstrapped sampling simply refers to **sampling by replacement**



Blue and red dots are repeated more often in the samples than they are present in the original data

# Bootstrapped Sampling

Bootstrapped Sampling at the data level

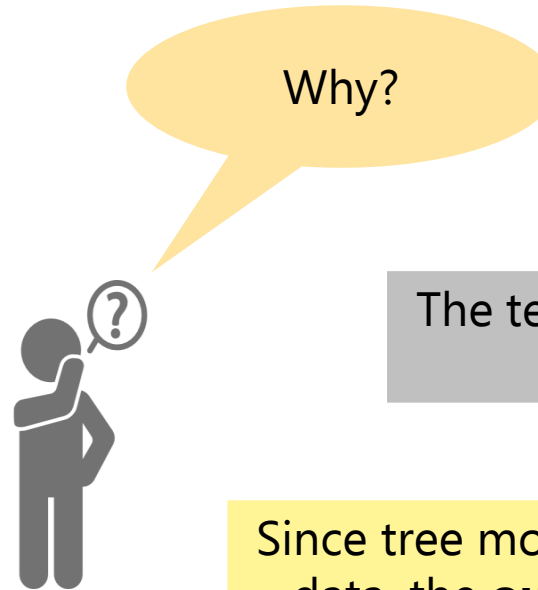


# Bootstrapped Sampling

Bootstrapped Sampling helps in reducing out of sample error

If an unpruned decision tree is fit into any data set then the model will have **very high out-of-sample error**

**Hypothetically**, if unpruned tree model is being fitted on **total population data** the error will be very low

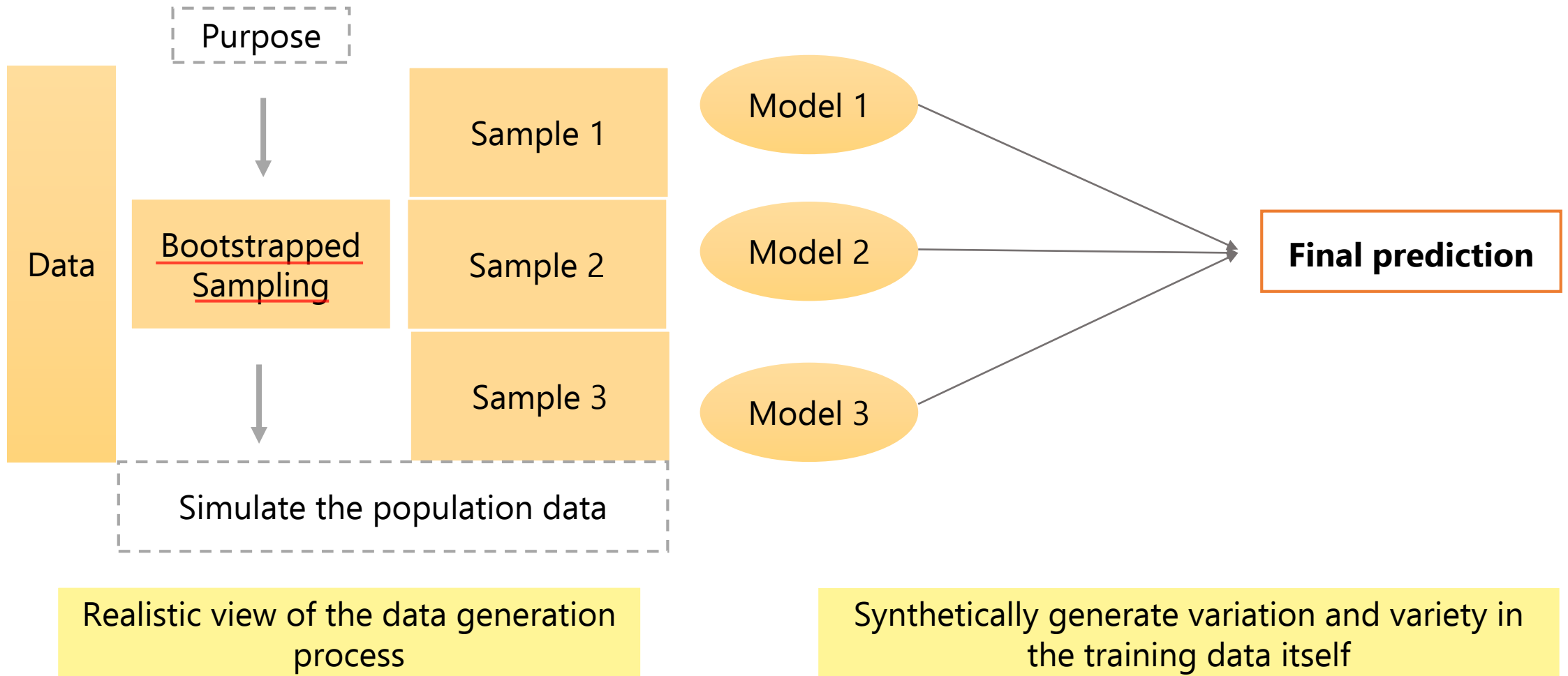


The test data can be very different from training data

Since tree model is being overfitted on training data, the **out-of-sample error will be high**

Low error due to all the variation and variety in data has been already seen by the model as the population data has been used to train the model

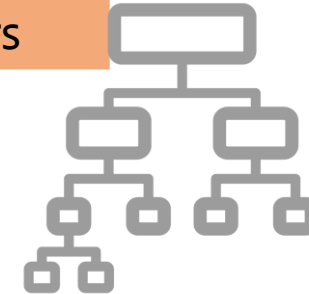
# Bootstrapped Sampling



# Bagged Trees

## Characteristics of Bagged Model

Using **unpruned decision trees** as base learners



Using **Bootstrapped Sampling** to create samples that are fed to each of the base learners

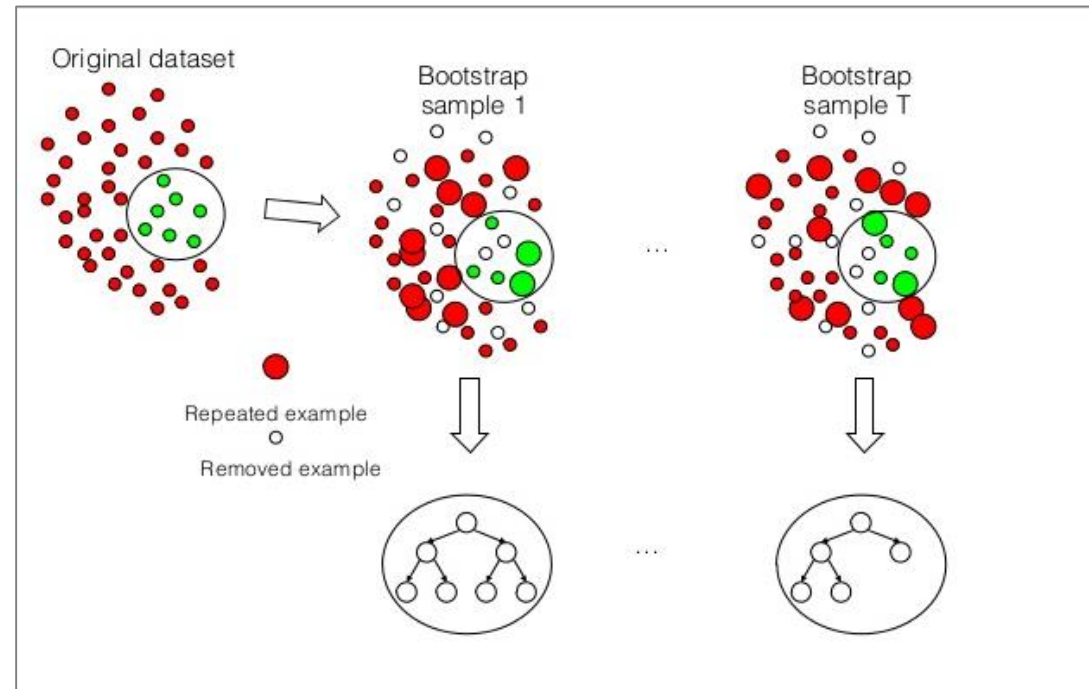




# Peculiarities of Bagged Trees

Bagged tree ensemble is comprised of multiple decision trees

Not interpretable as a linear model or simple decision tree



Qualitative statement on ensemble models

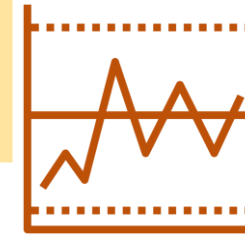


Identify the important predictors by looking at **Variable Importance**



# Variable Importance

Variable Importance - Averaging or summing the improvement in **Gini or Entropy for a classification model** and **RSS for a regression model** for all the variables



Bagged tree ensemble model contains many tree models

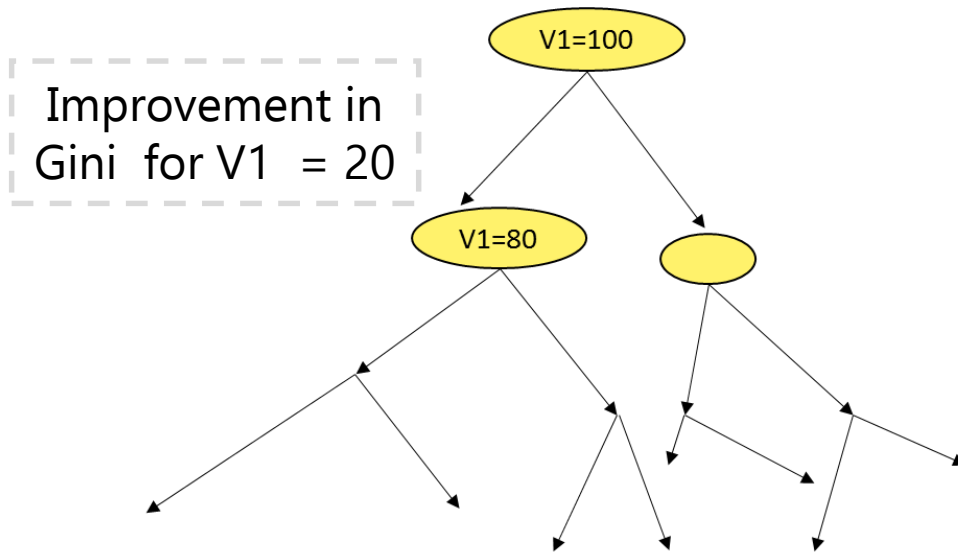
**Feature importance** of each variable in each of the constituent trees

Tracking the decrease in Gini metric and weighing this decrease appropriately

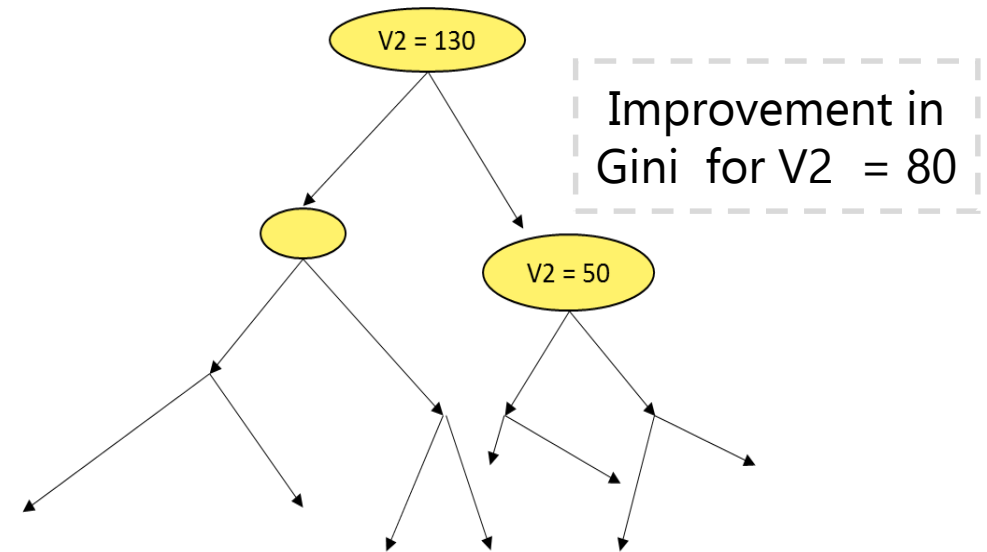


# Variable Importance

## Example



Tree 1  
Gini Measure for each split



Tree 2  
Gini Measure for each split

Computing variable importance for all the variables used in the split

# Variable Importance

Ensemble has N trees

Improvement in Gini/RSS Across Splits

Variable	Tree 1	Tree 2	Tree 3	.....	Tree N
V1	300	30	12	.....	0
V2	600	0	200	.....	150
...	...	...	...	.....	...
$V_k$	120	450	30	.....	19

Variable	Variable Importance
V1	$\frac{(300 + 30 + 12 + \dots + 0)}{N}$
V2	$\frac{(600 + 0 + 200 + \dots + 150)}{N}$
...	....
$V_k$	$\frac{(120 + 450 + 30 + \dots + 19)}{N}$

The average values of importance measures per variable will produce a consolidated number



# Parameters of Bagged Trees

What could be the user specified parameters while building a bagged tree model?



User specified parameters or  
**Hyperparameters**

Number of tree used  
to build an ensemble

Depth of the tree

Number of  
observations per node  
of a tree



# Parameters of Bagged Trees

User specified parameters have an implication

Different ensemble model depending on different parameters



Which among the three model?

**K-Fold CV** to get an estimate of out of sample error

Expensive

**Out of Bag Error** is generally used in most tree based models

Model 1:  
Trees = 100  
Depth of Tree = 4

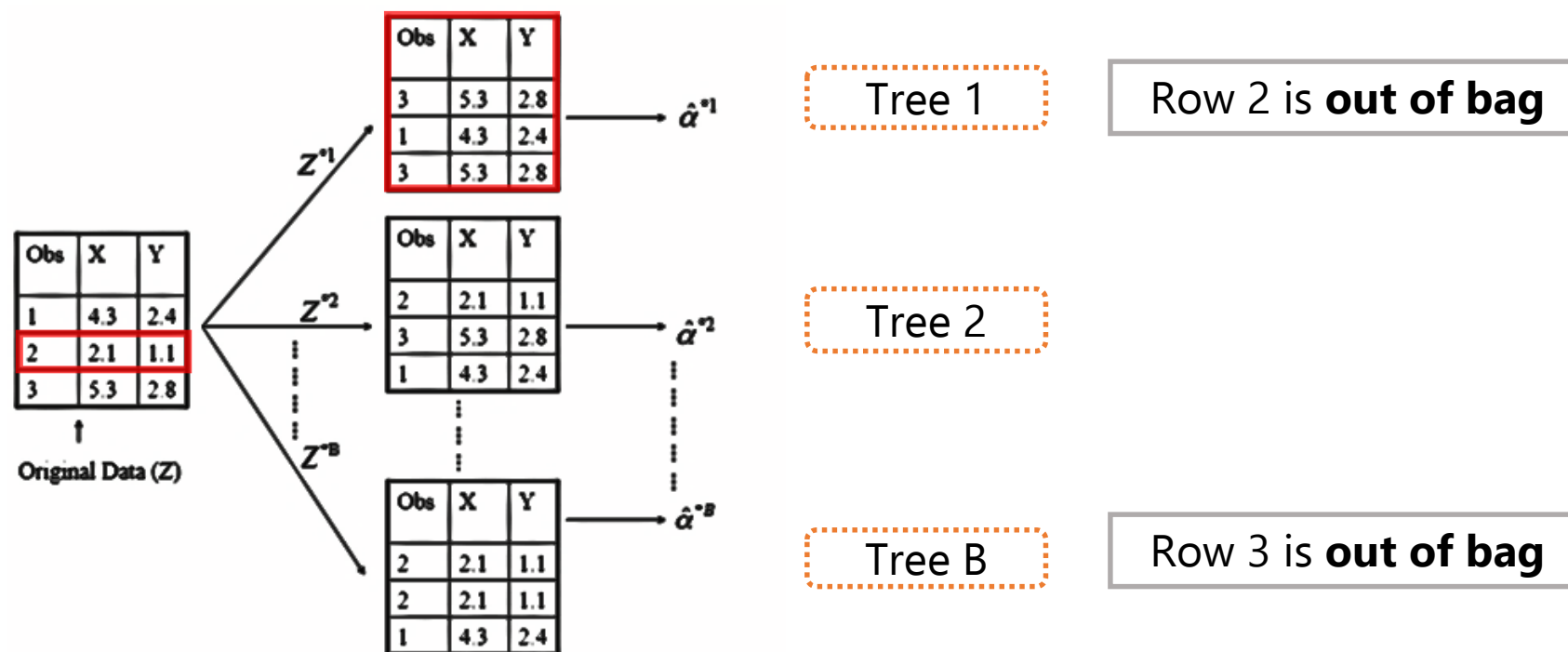
Model 2:  
Trees = 150  
Depth of Tree = 3

Model 3:  
Trees = 500  
Depth of Tree = 4



# Out Of Bag Error (OOB)

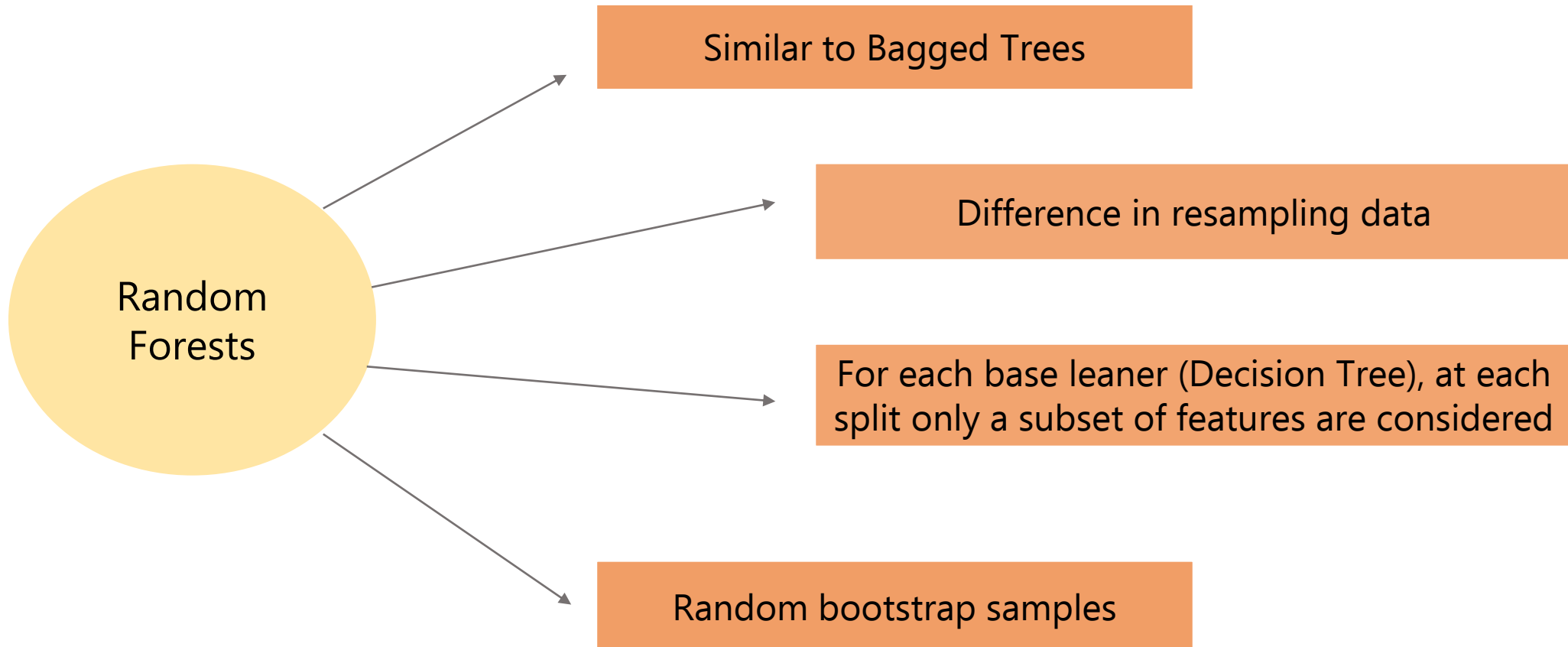
In Bootstrap Sampling, some observation gets left out from the original data



Average Out of Bag observations in Bootstrapped Sampling is around **33%**

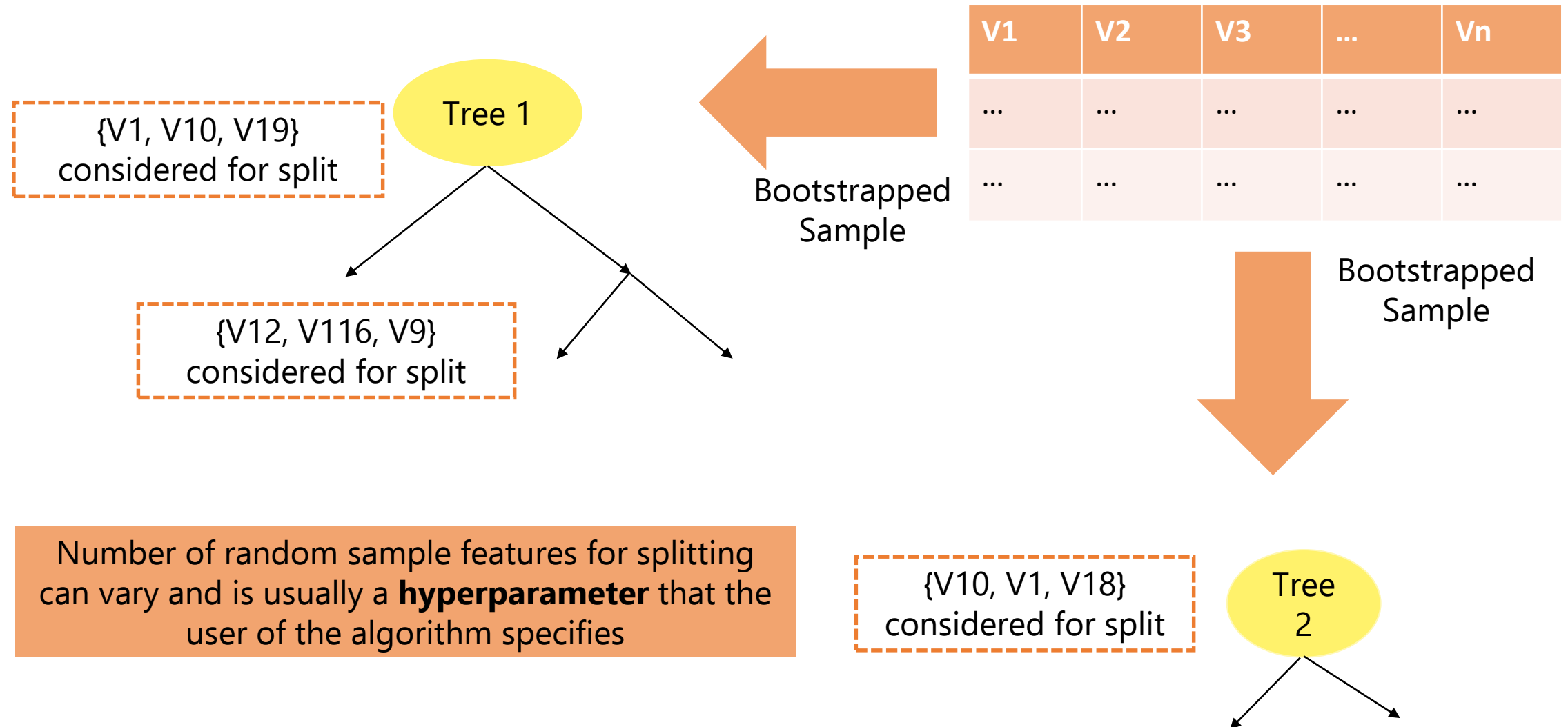


# Random Forests





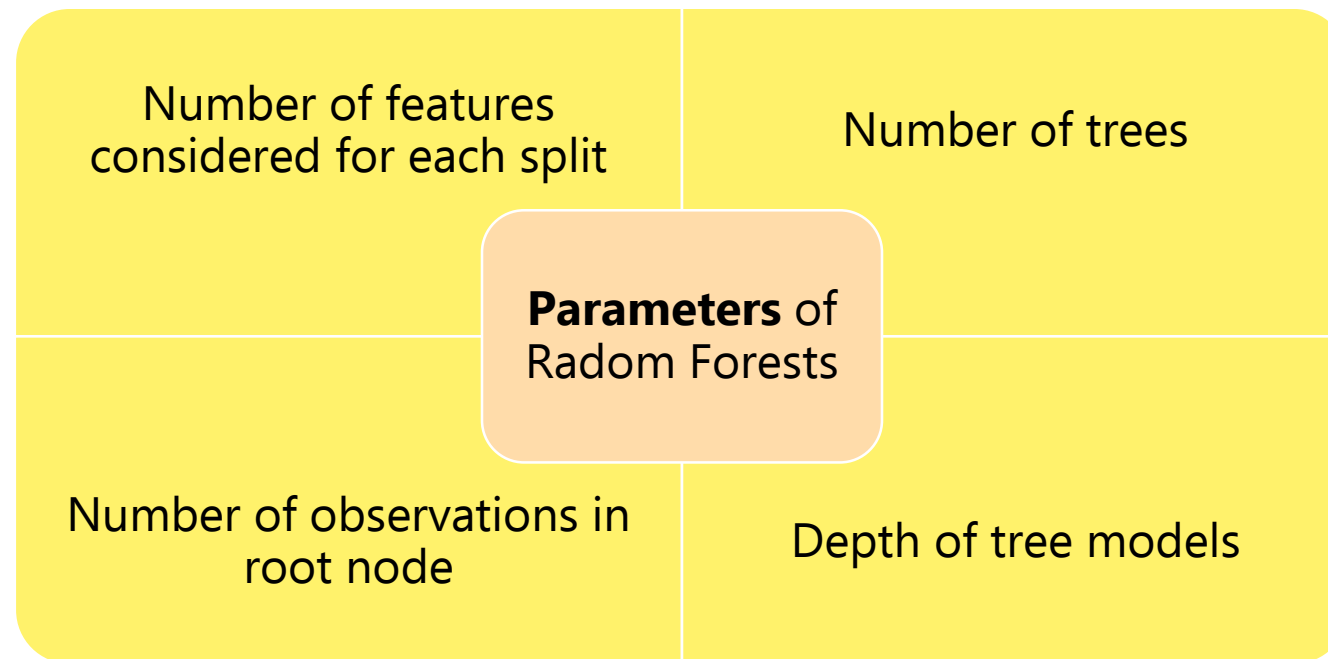
# Random Forests



# Random Forests

Random Forests uses tree models as base learner

Extract **Variable Importance** or compute **Out-of-bag Error** to get an estimate of out of sample model performance for parameter tuning



# Recap

- Tree based ensembles overview
- Tree based ensembles models – Bagged Tree and Random Forests
- Bootstrapped sampling
- Variable importance
- Out of bag error (OOB error)
- Random Forests
- Code Demo

