

Random Forest Ensemble Method

Session-13



SUMMARY of PREVIOUS CLASS

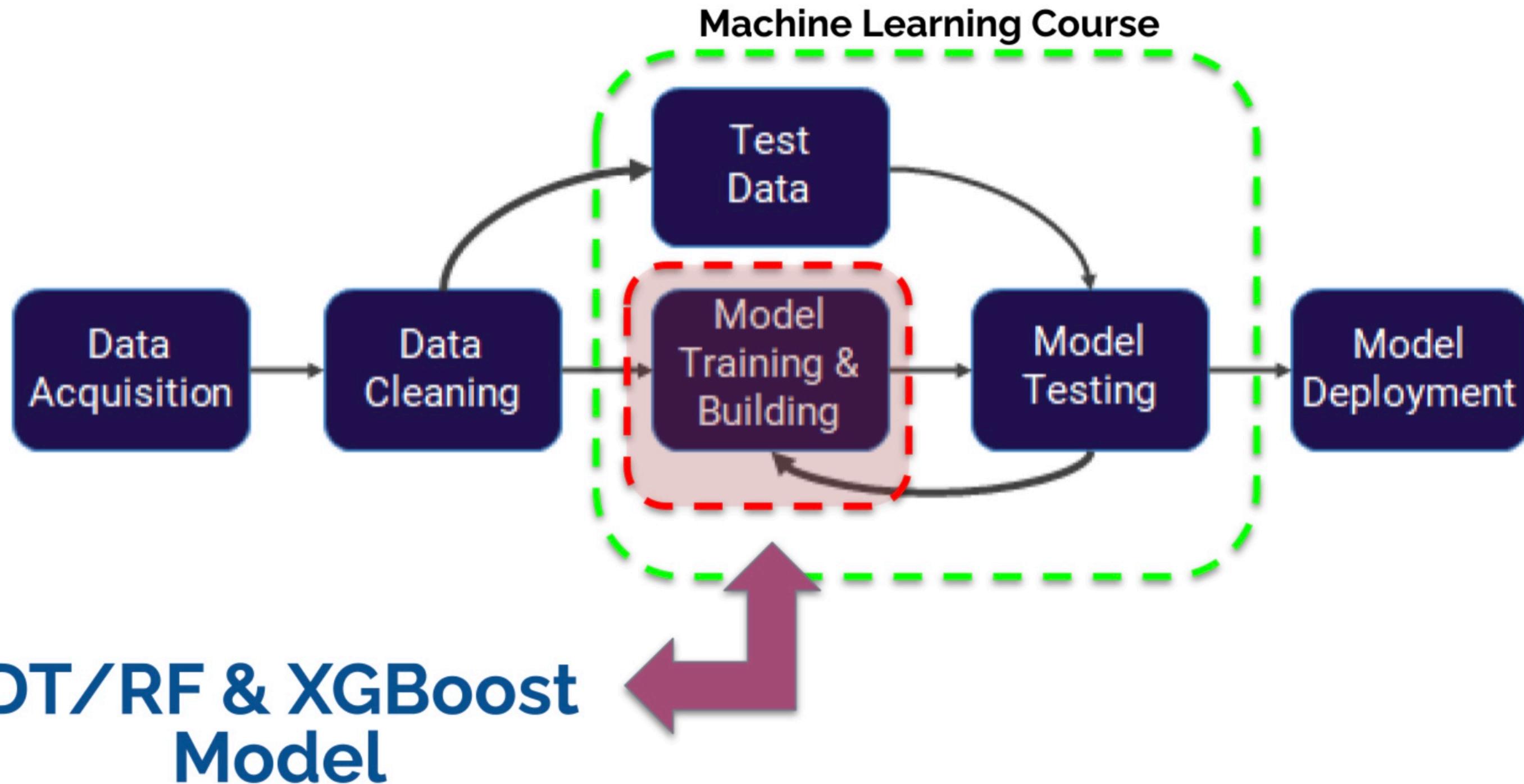
- **Decision Tree Theory:** Root node, Decision node, Leaf Node
- **Gini Index:** minimize the probability of misclassification.
- **Information Gain-Entropy:** obtain the largest IG by minimizing the entropy
- **Regression:** Minimize variance, calculate mean



- Ensemble Methods
- Bagging
(Bootstrap AGGRegatING)
- Random Forest Theory
- Random Forest /w Python

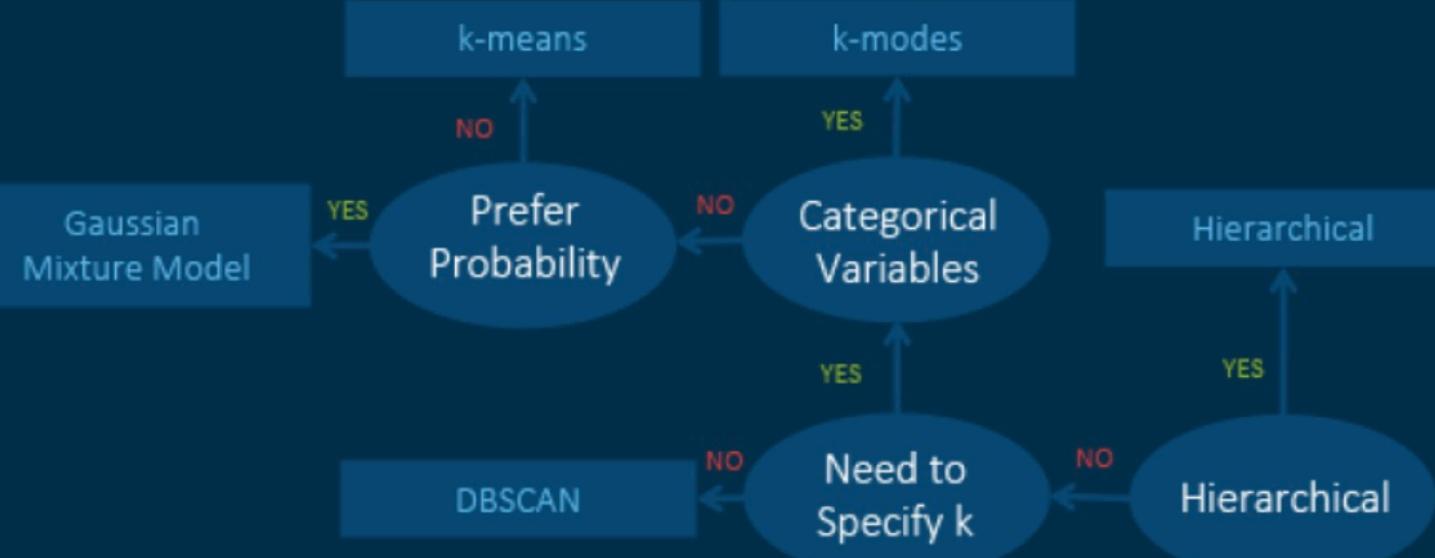


Where are we?

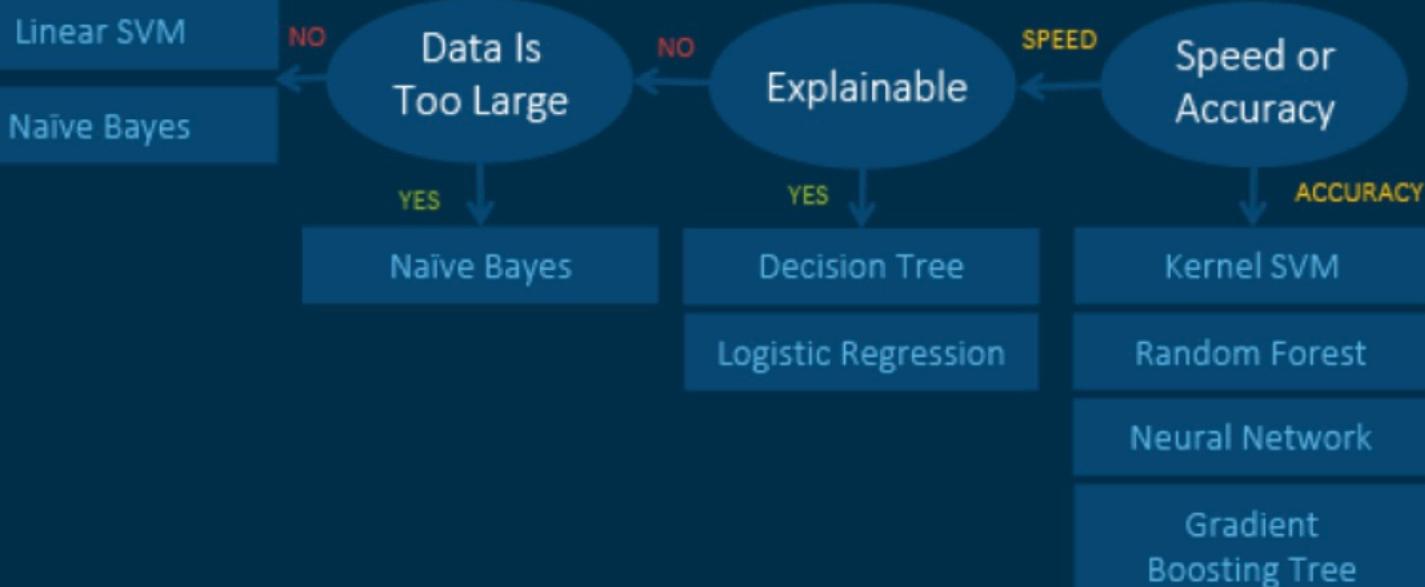


Machine Learning Algorithms Cheat Sheet

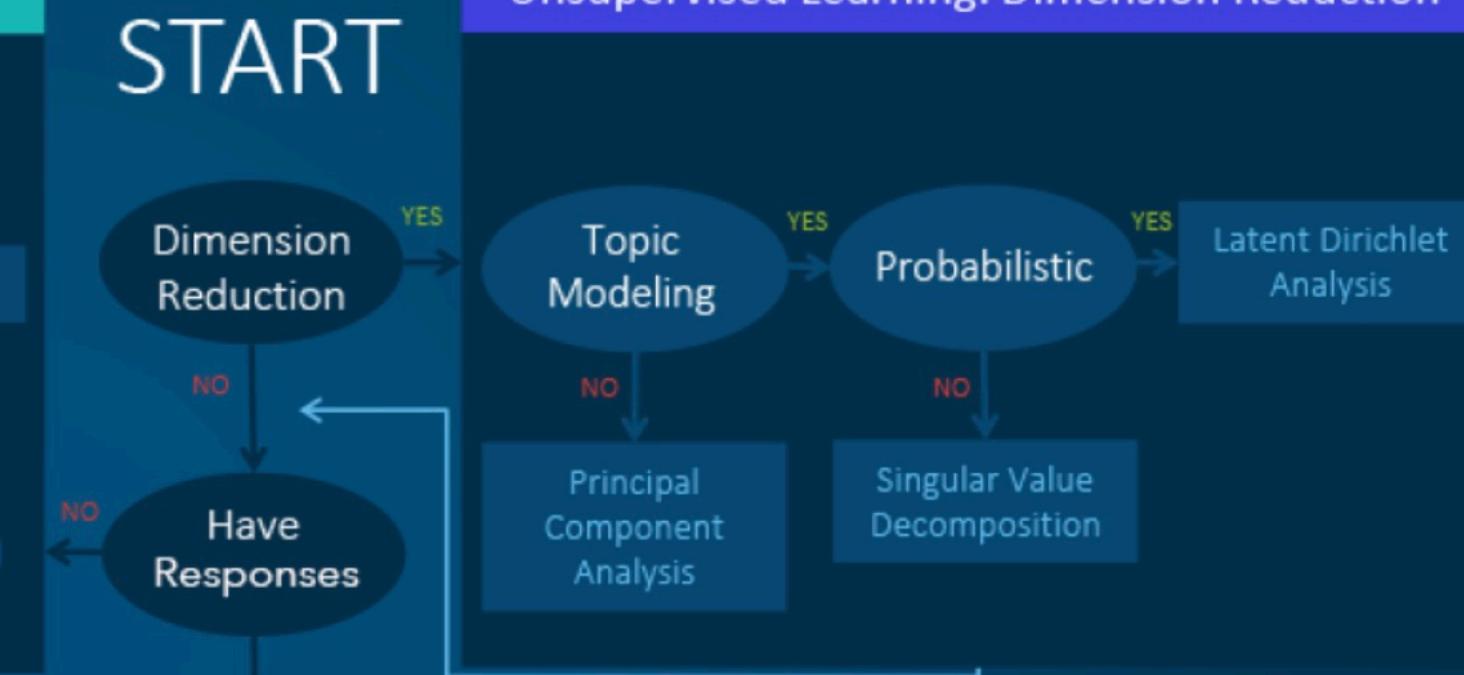
Unsupervised Learning: Clustering



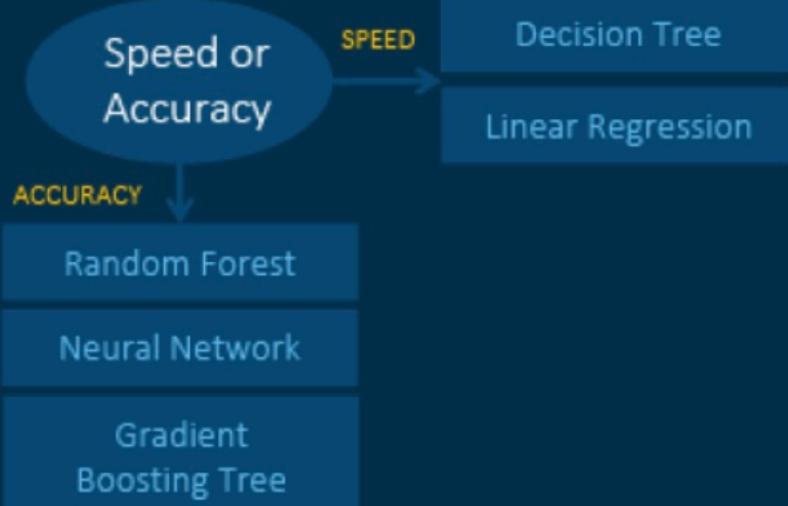
Supervised Learning: Classification



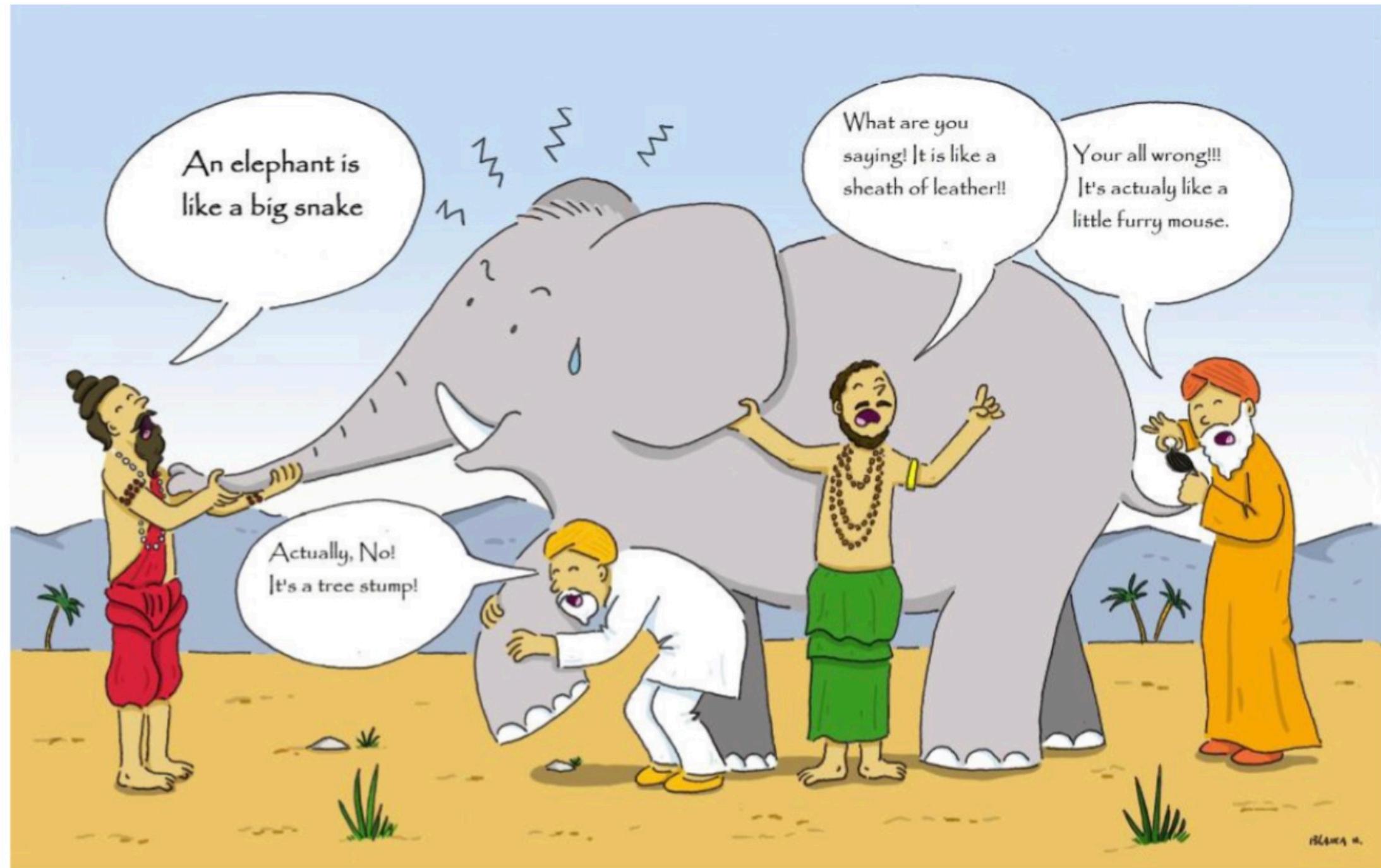
Unsupervised Learning: Dimension Reduction



Supervised Learning: Regression



Ensemble Methods



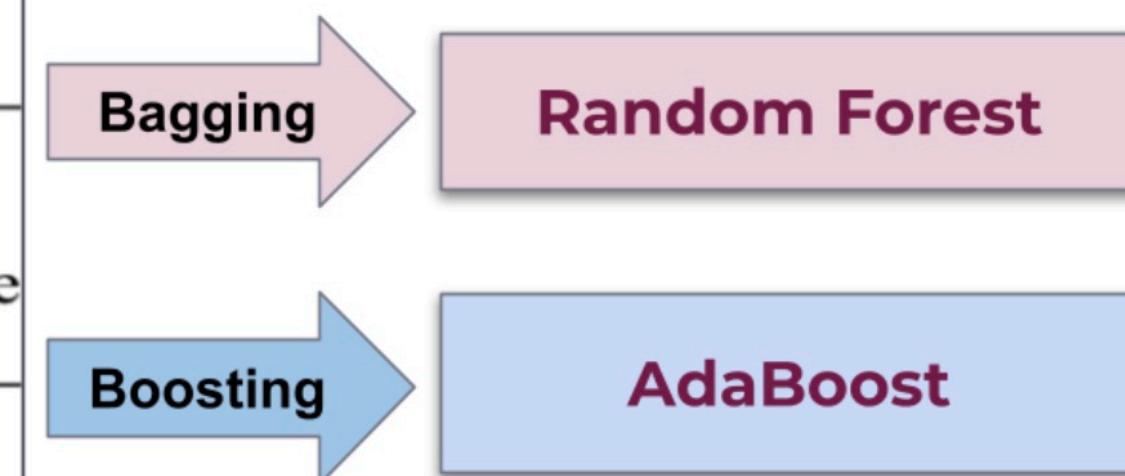
Ensemble Methods



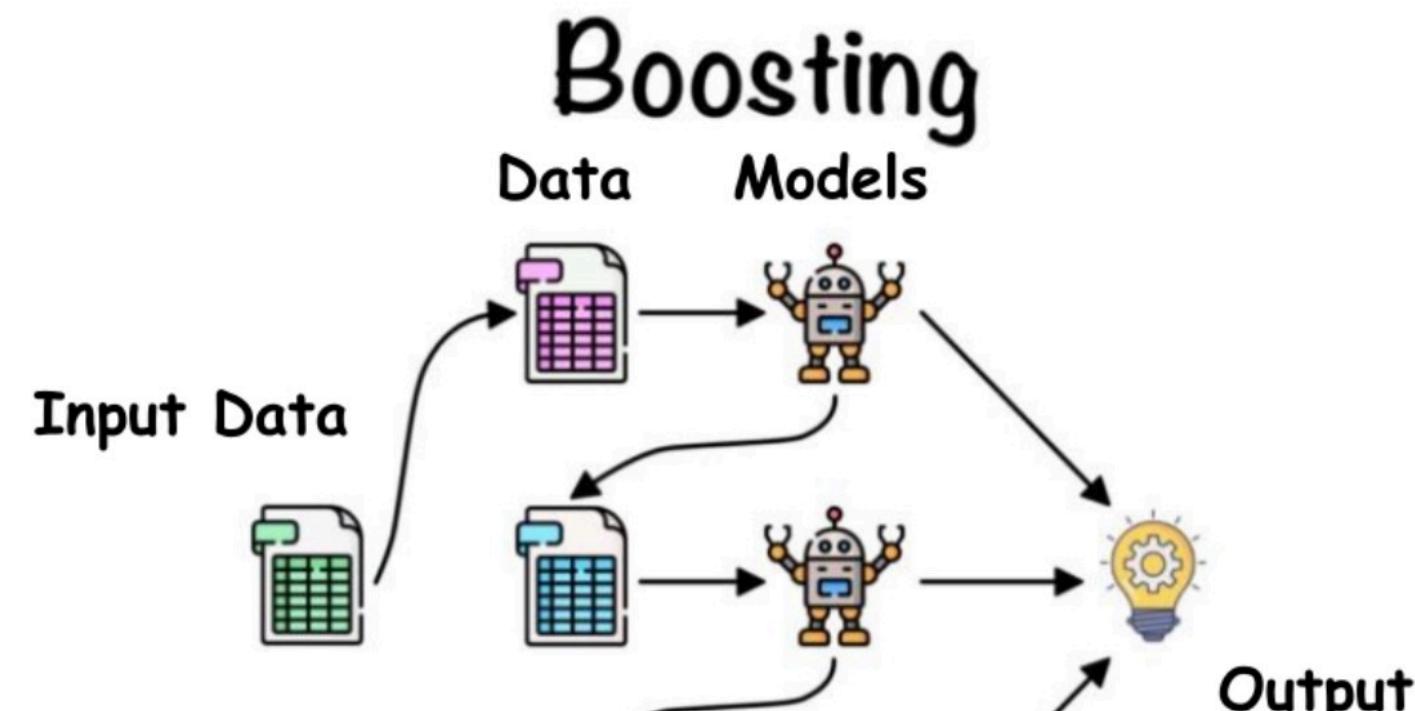
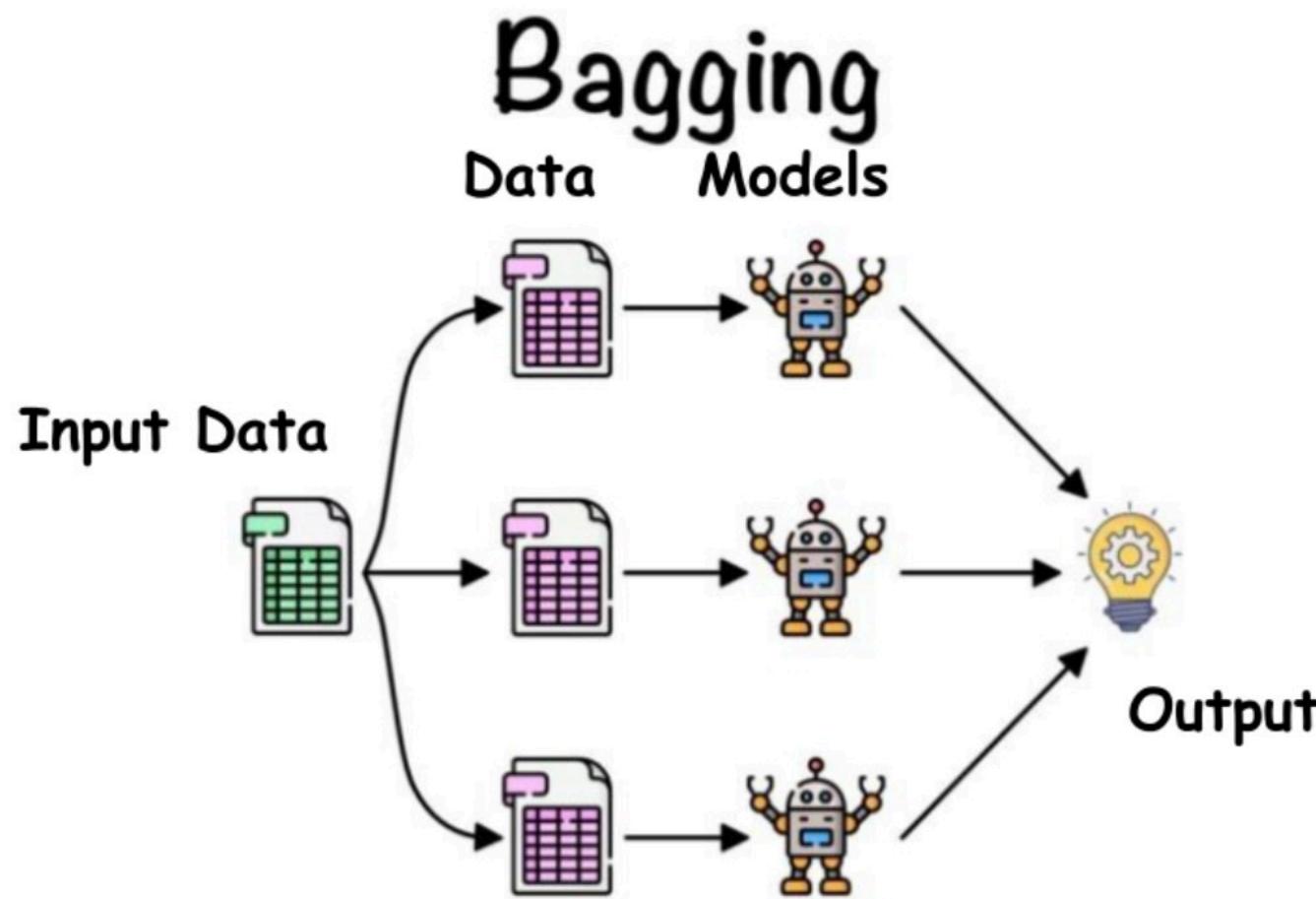
General info about Ensemble Methods: Bagging & Boosting

Ensemble methods provide training of **different models on the same data set**. Each model makes its own prediction and a **metamodel** is formed from the **combination of predictions**. Thus, a metamodel is created that makes the most successful prediction. **Bagging-method** prefers strong **sub models**, in contrast **boosting-method** works on **weak learners**.

	Bagging	Boosting
Similarities	<ul style="list-style-type: none">• Uses voting• Combines models of the same type	
Differences	Individual models are built separately	Each new model is influenced by the performance of those built previously
	Equal weight is given to all models	Weights a model's contribution by its performance



Ensemble Methods



Parallel

Sequential

Ensemble Methods

General info about Ensemble Methods: Boosting & Bagging

Bootstrap aggregating or Bagging is a ensemble meta-algorithm combining predictions from multiple decision trees through a majority voting mechanism

Models are built sequentially by minimizing the errors from previous models while increasing (or boosting) influence of high-performing models

Optimized Gradient Boosting algorithm through parallel processing, tree-pruning, handling missing values and regularization to avoid overfitting/bias

Bagging



Decision Trees

A graphical representation of possible solutions to a decision based on certain conditions



Bagging-based algorithm where only a subset of features are selected at random to build a forest or collection of decision trees



Gradient Boosting

Gradient Boosting employs gradient descent algorithm to minimize errors in sequential models

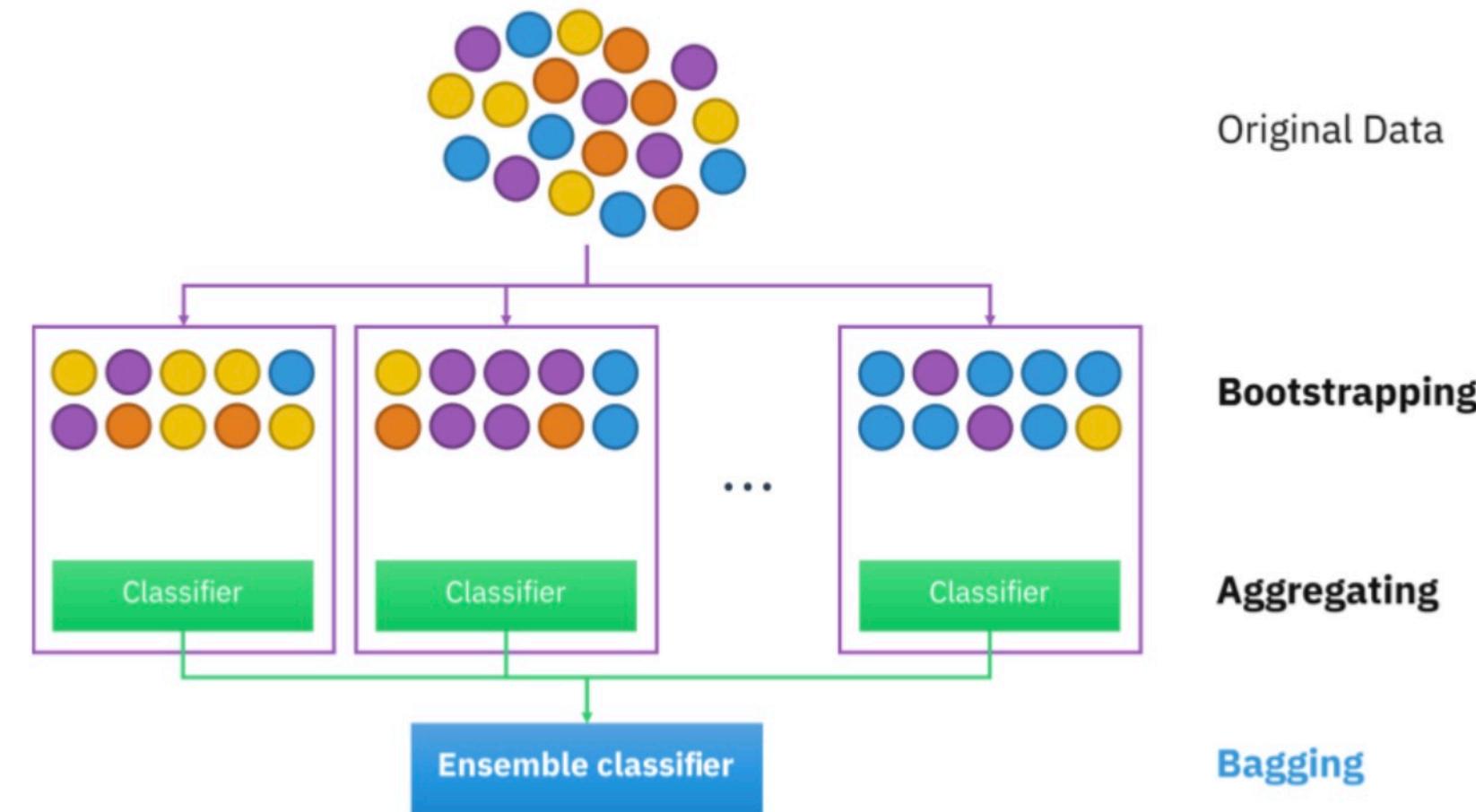


Boosting

Ensemble Methods

Bagging (Bootstrap AGGRegatING) :

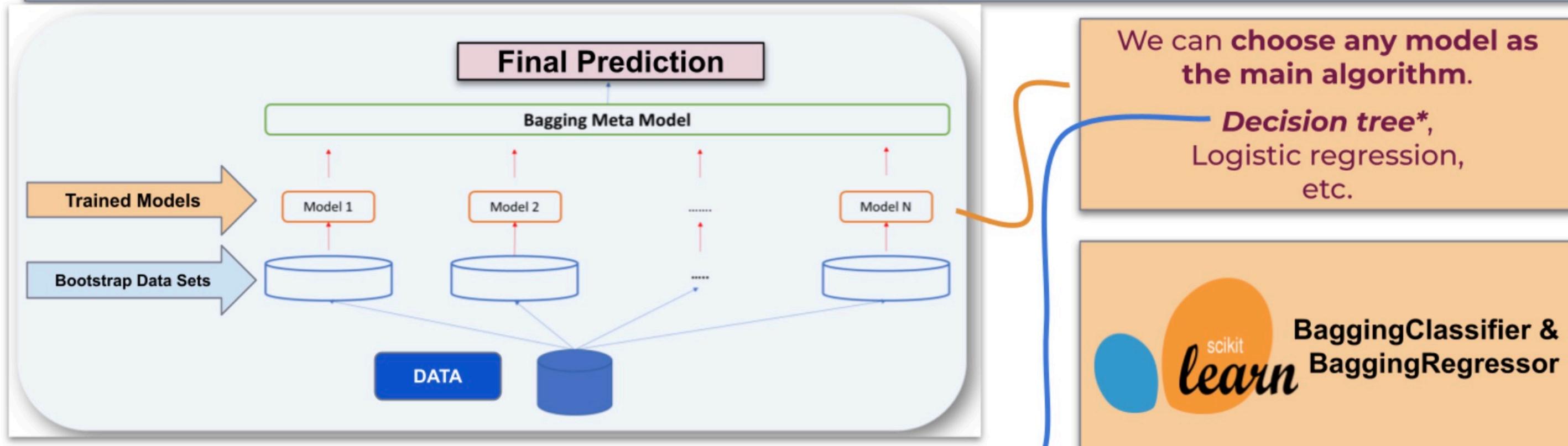
Bagging is used instead of "**bootstrap aggregation**". We can define **Bootstraps** as creating a **new sample data set** using values from the original data set. The clue is that each sub-sample has the same size but different parts of the observations of the original dataset.



Ensemble Methods

Bagging (Bootstrap AGGregatING) :

In bagging methods, **more than one sub dataset created from the data set is trained with the same algorithm**. With these different data sets trained, a final model is created. While creating this model, **voting (mode)** is made for **classification** problems, and **averages (mean)** are taken for **regression** problems.

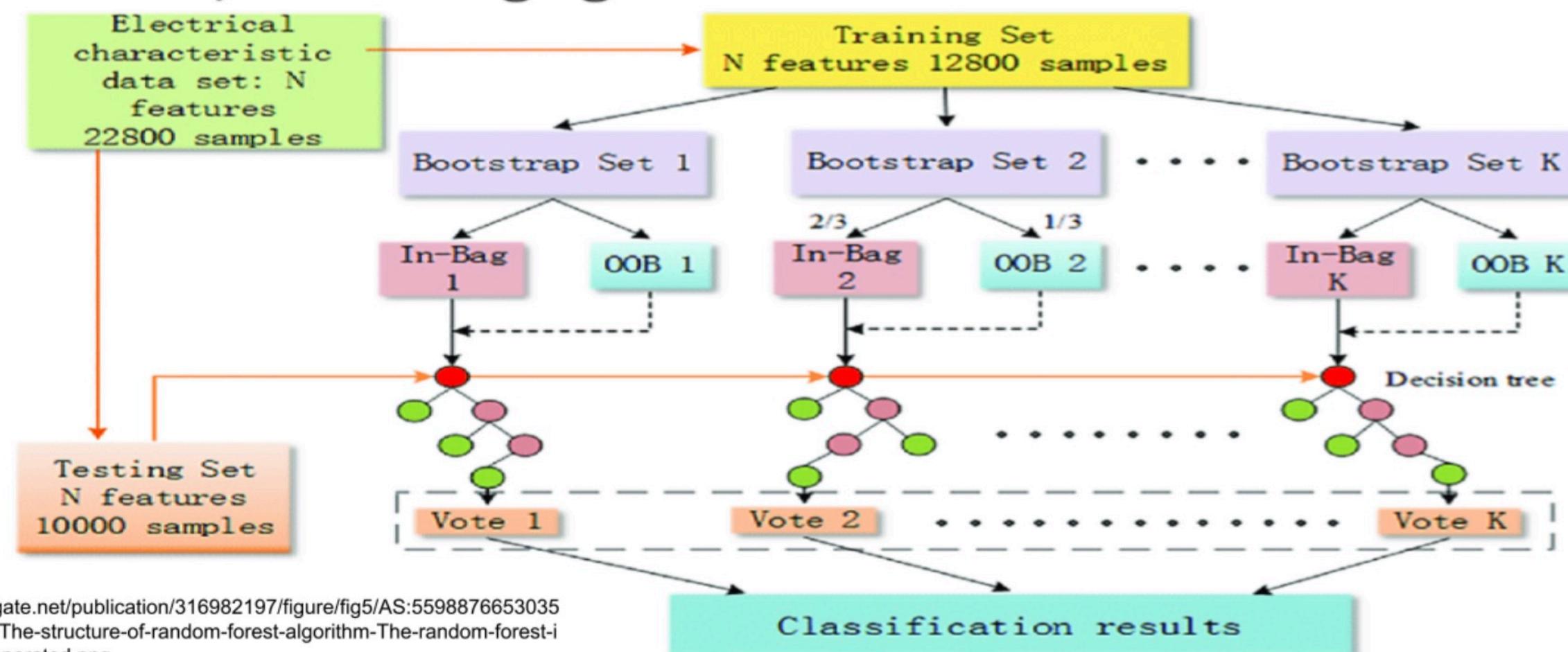


*Random Forest is a bagging method using Decision Tree as the main algorithm.

Random Forest Algorithm

Random Forest Algorithm:

- Random Forests (RF) is an ensemble bagging algorithm composed of many individual trees (Breiman, 2001)
- RF reduces the variance of individual trees by randomly selecting (and thus de-correlating) many features from the dataset, and averaging them



Random Forest Algorithm



Random Forest Algorithm:



Prevent Overfitting

In addition to taking the random subset of data, as in bagging, Random Forest uses a random subset of features rather than using all features to grow trees.

Hyperparameters:

“n_estimators” parameter: (default=100)

The number of trees in the forest.

The more # of trees, the better accuracy. But CPU intensive.

Random Forest Algorithm

Hyperparameters:

“max_depth” parameter: (default=None)

The maximum depth of the tree.

If None, then nodes are expanded until all leaves are pure.

“max_features” parameter: (default="sqrt")

Number of features to consider when looking for the best split.

Increase will improve the performance, but results in a correlation between the trees.

Random Forest Algorithm

Hyperparameters:

“bootstrapbool” parameter: (default=True)

Whether bootstrap samples are used when building trees. If False, the whole dataset is used to build each tree.

“oob_scorebool” parameter: (default=False)

This score can be used as validation score, in case of you have small data set

Whether to use out-of-bag samples to estimate the generalization score. Only available if bootstrap=True.

Random Forest Algorithm



Pros & Cons

Pros:

- Good for multi-class problems
- Excellent Predictive Powers
- Easy Data Preparation
- Ranking for feature importance
- Suitable for large data
- Can handle missing values (but not in Scikit-learn)
- Compare with DT robust to overfitting

Cons:

- Not good on smaller data sets
- Training could be computationally expensive for numerous trees
- For Regression, not as good option as Classification
- Lack of interpretability
- Biased towards high cardinality categorical features

RF Algorithm



Be ready for
RF
Python 1 & 2
Session