

© 2020 Trilogy Education Services, a 2U, Inc. brand. All Rights Reserved.

## Class Objectives

In today's class we'll learn about a new ML algorithms family: Tree based algorithms

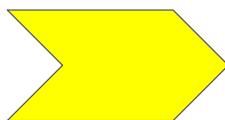
-  Decision trees
-  Random forest
-  Weak learners
-  Ensemble methods

## Categorical Data

---

Also we'll learn how to deal with categorical data

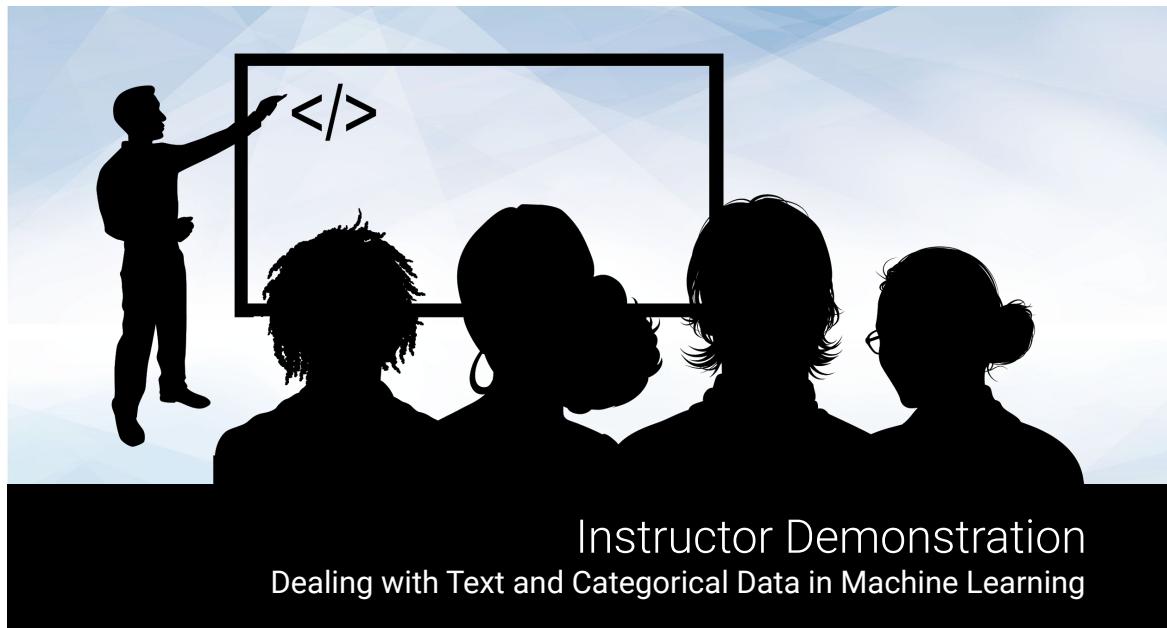
Color
Red
Red
Yellow
Green
Yellow



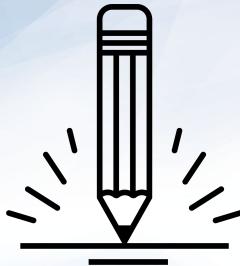
Red	Yellow	Green
1	0	0
1	0	0
0	1	0
0	0	1

---

3



4



## **Activity:** Encoding Categorical Data for Machine Learning

In this activity, you will be tasked with encoding categorical and text features of a dataset that contains 2,097 loan applications.

Suggested Time:  
10 minutes



5



## **Time's Up! Let's Review.**

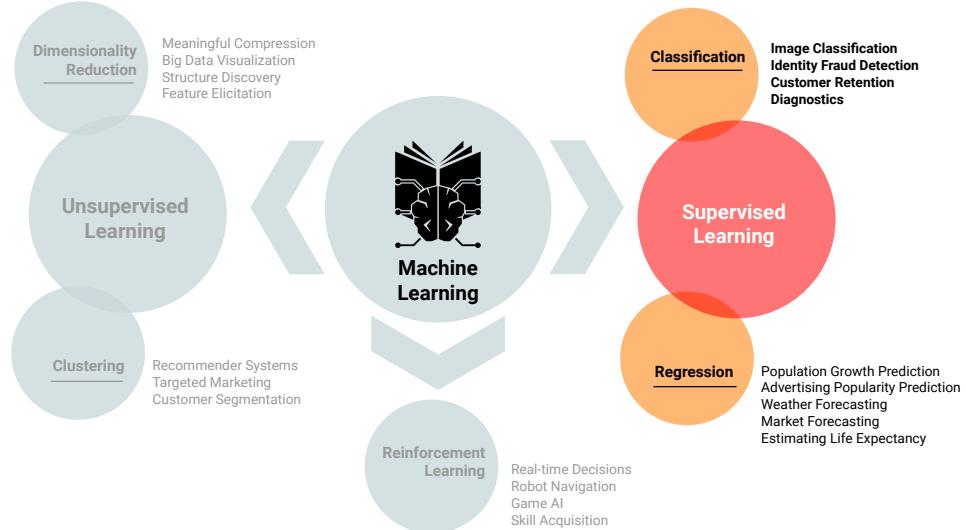
6

# Walking into the Algorithms Forest

7

## Tree based algorithms

Tree based algorithms, are part of the supervised machine learning methods.



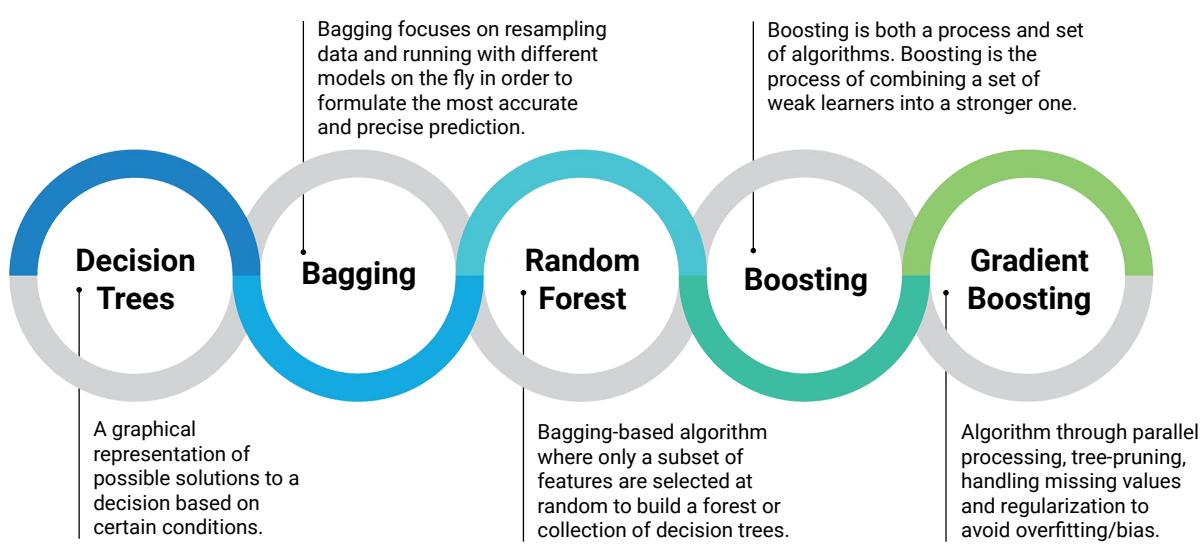
8



**Tree-based algorithms** are supervised learning methods that are mostly used for classifications and regression problems.

9

## Tree Based Algorithms at a Glance



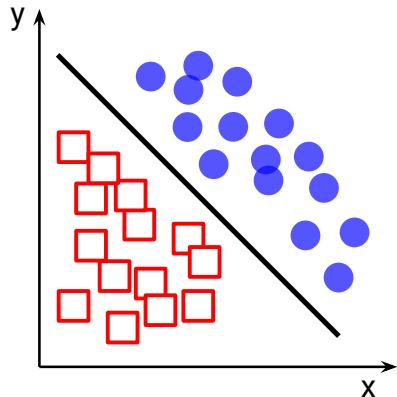
10

# Algorithms

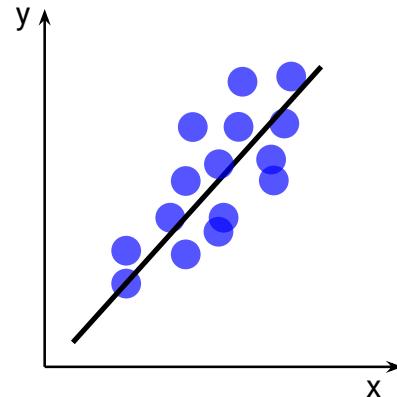
---

These algorithms can be used to solve classification or regression problems.

## Classification



## Regression

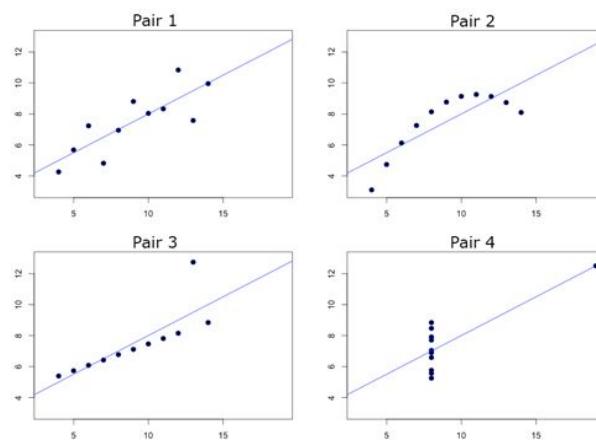


11

## Linear vs. Non-Linear Models

---

In linear models, the relationship among input variables can be represented as a straight line, while non-linear models have a different shape.



12

## Linear Models

---

Predicting the price of a house based on its size is an example of a linear problem. This is because, as a general rule, the size of the house is directly proportional to the price of the house.

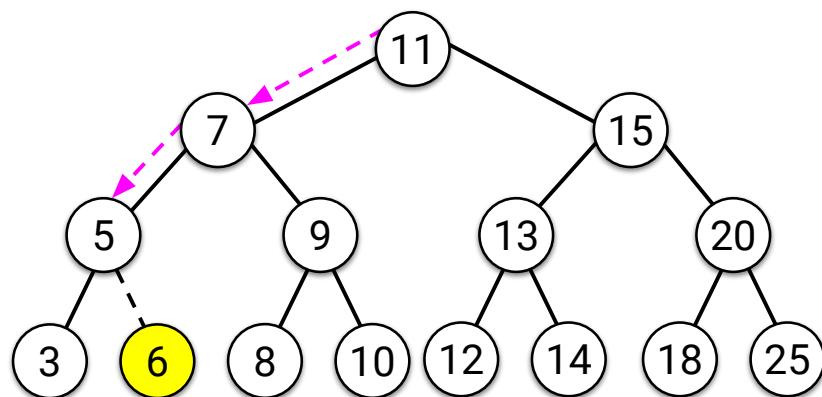


13

## Non-Linear Models

---

Tree-based algorithms can map non-linear relationships in data.



14

## Non-Linear Models

Predicting if a credit application is going to be fraudulent or not may be an example of a non-linear problem due to the complex relationship between the input features and the output prediction.



15

## Tree-Based Algorithms

These algorithms are quite often used in finance for assessing risk, preventing fraud, or fighting money laundering.



16

## sklearn

sklearn has two modules that implement tree-based algorithms that we will be covering Today.

01    `sklearn.tree`

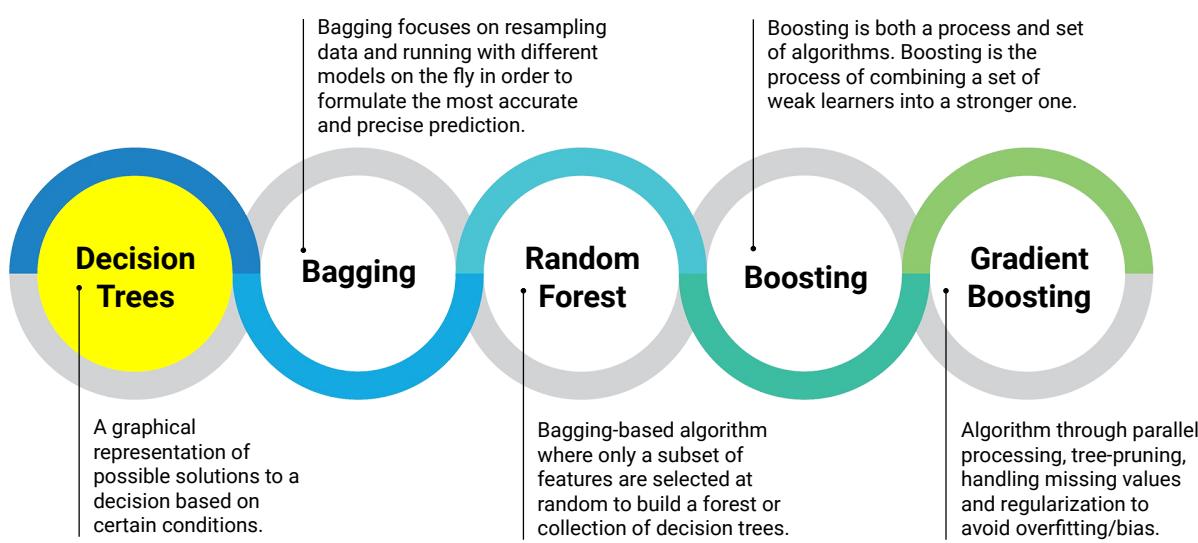
`sklearn.tree` implements decision trees.

02    `sklearn.ensemble`

`sklearn.ensemble` offers implementations for random forest, gradient boosting, boosting and bagging algorithms.

17

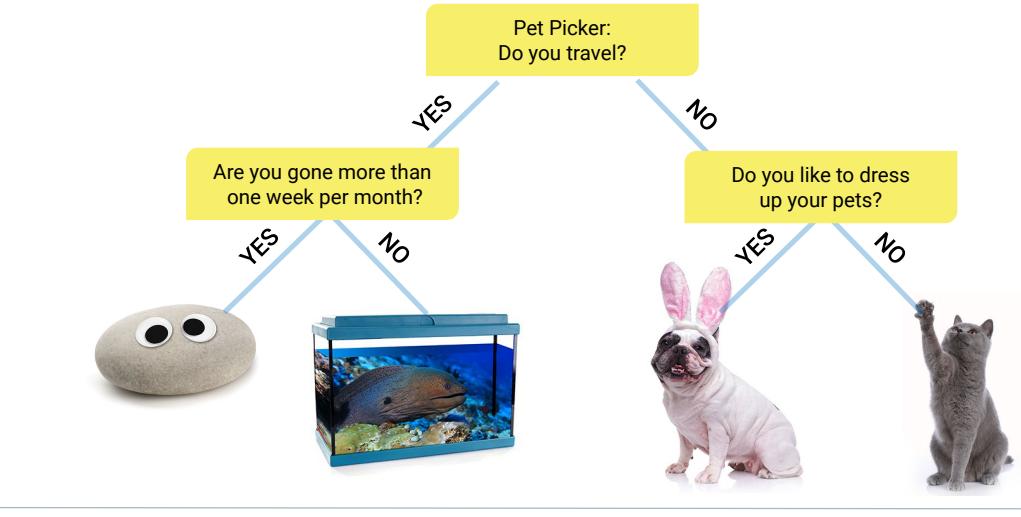
## Decision Trees



18

## Decision Trees

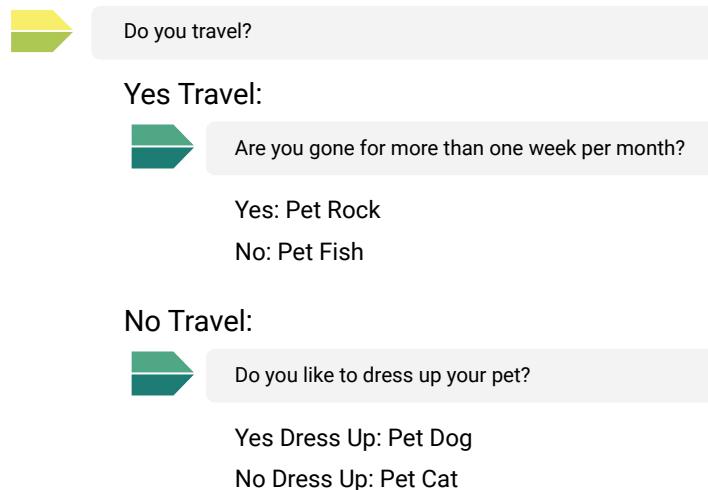
Decision trees encode a series of true/false questions.



19

## Decision Trees

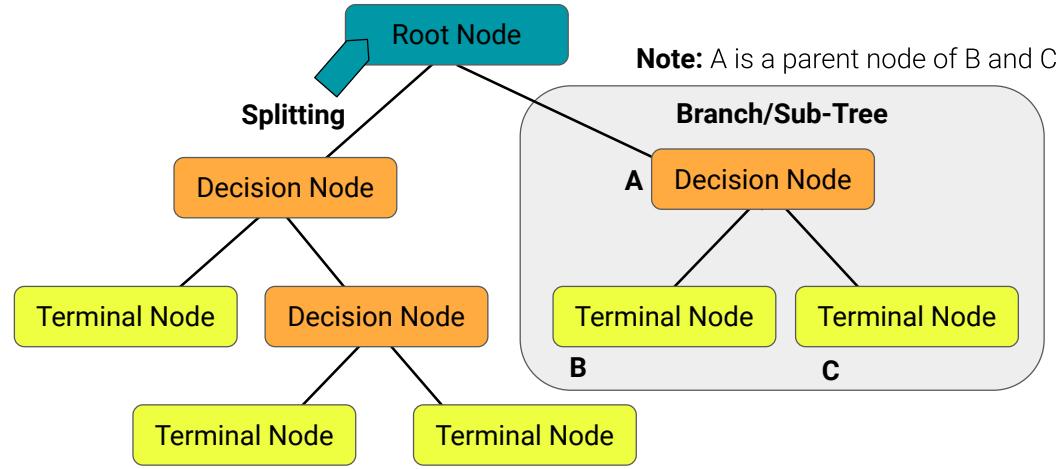
These true/false questions can be represented with a series of if/else statements



20

## Root Node

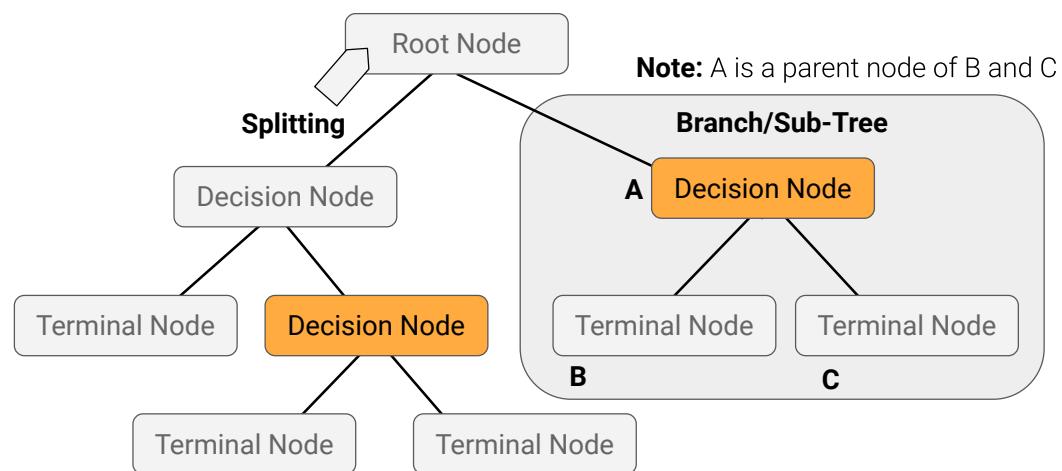
Represents the entire population or sample data, this node gets divided into two or more homogeneous sets.



21

## Parent Node

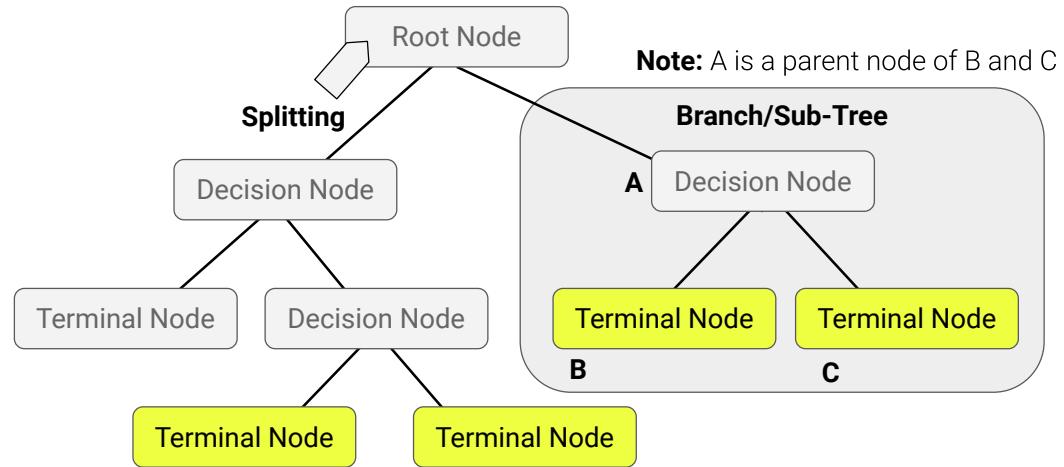
A node that is divided into sub-nodes.



22

## Child Node

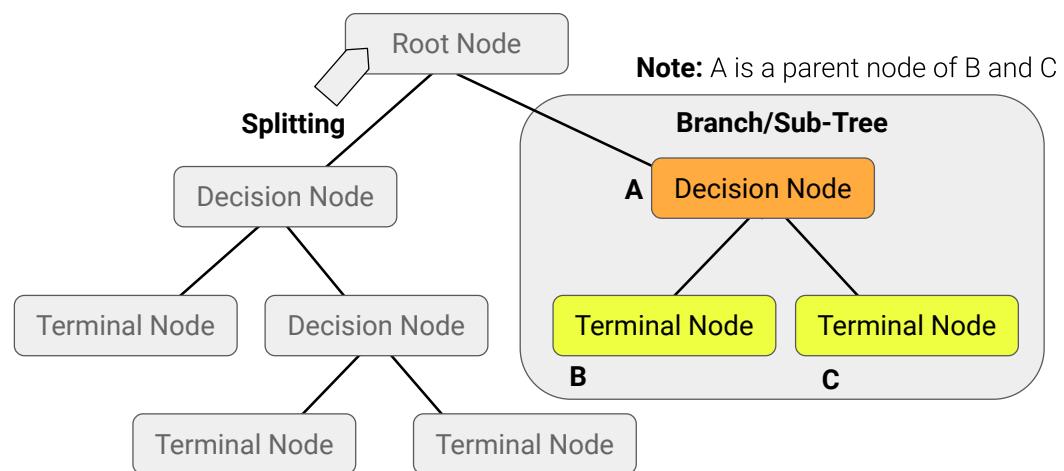
Sub-nodes of a parent node.



23

## Decision Node

A sub-node that is split into further sub-nodes.



24

## Decision Trees

---

Key concepts to know while working with decision trees:

<b>Leaf or Terminal Node</b>	Nodes that do not split.
<b>Branch or Sub-Tree</b>	A subsection of entire tree.
<b>Splitting</b>	Process of dividing a node into two or more sub-nodes.
<b>Pruning</b>	Process of removing sub-nodes of a decision node.
<b>Tree's Depth</b>	The number of decision nodes encountered before making a decision.

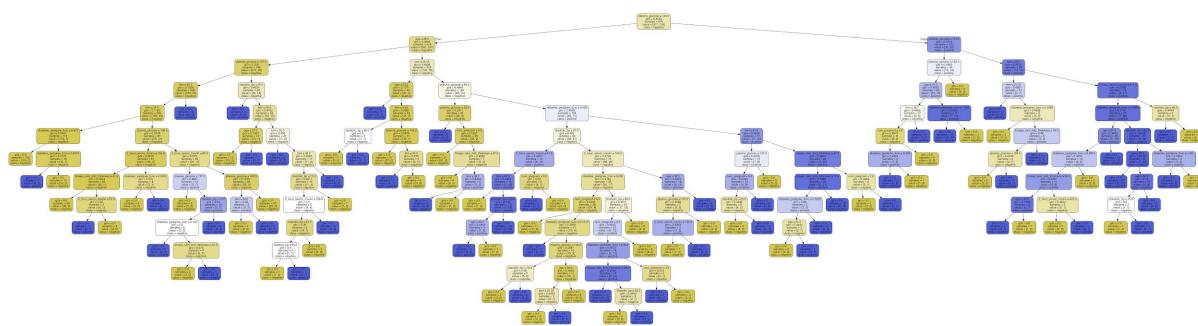
---

25

## Decision Trees

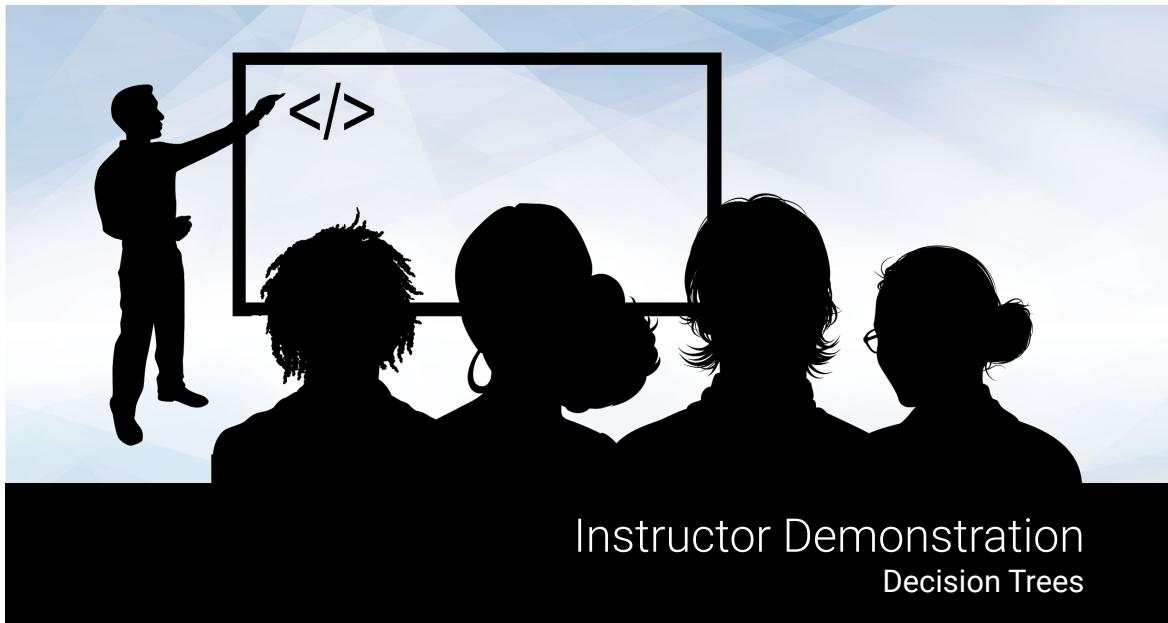
---

Decision trees can become very complex and may not generalize well.



---

26



## Instructor Demonstration Decision Trees

27

A graphic icon of a pencil with several short lines radiating from its tip, set against a light blue background with a geometric pattern.

## **Activity:** Predicting Fraudulent Loans Applications

In this activity, you will create a decision tree model to predict fraudulent loan applications.

Suggested Time:  
10 minutes

A small icon of a clock face with a circular arrow or hand pointing around it, located next to the suggested time text.

28



**Time's Up! Let's Review.**

29

**Introduction to Ensemble Learning**

30

## The Classification Algorithm Race

---

If we compare the performance of classification algorithms, we'll find that some algorithms performed better than others



---

31

## Weak Learners

---

- Algorithms that actually fail at learning in an adequate fashion.
- They are a consequence of limited data to learn from.
- Their predictions are only a little better than random chance.



---

32

## Weak Learners are still valuable in Machine Learning

---

They can be combined with other classifiers in order to make a more accurate and robust prediction engine.

Combined weak learners are an example of ensemble learning:



33

## Ensemble Learners

---

Ensemble learners improves accuracy and robustness, as well as decrease variance.

Combined, weak learners can perform as well as strong learners.



34

## Combining Weak Learners

Weak learners have to be combined using specific algorithms like:



GradientBoostingTree



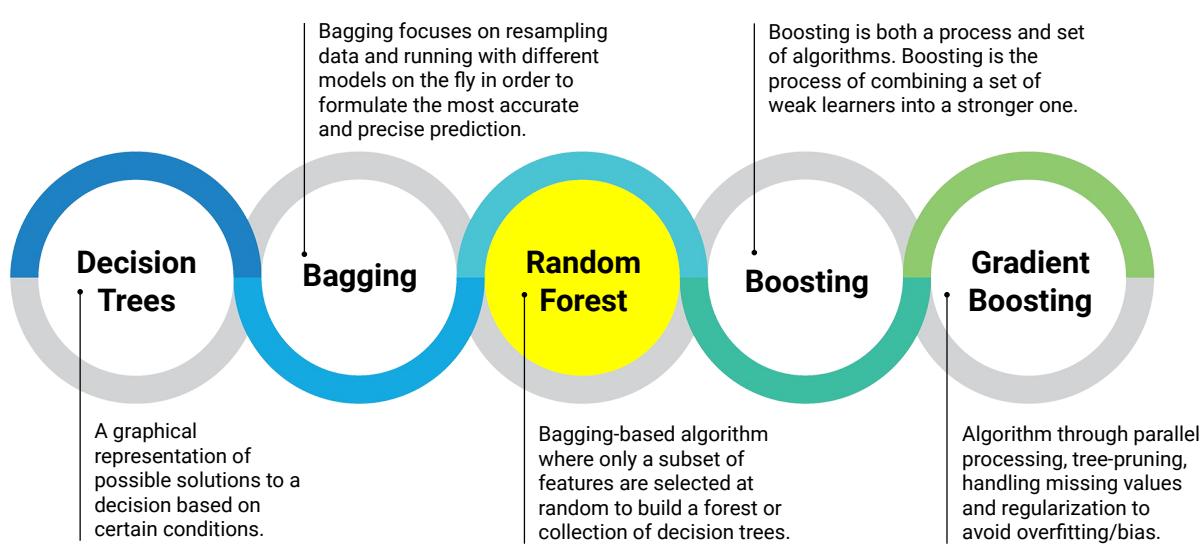
XGBoost



Random Forest

35

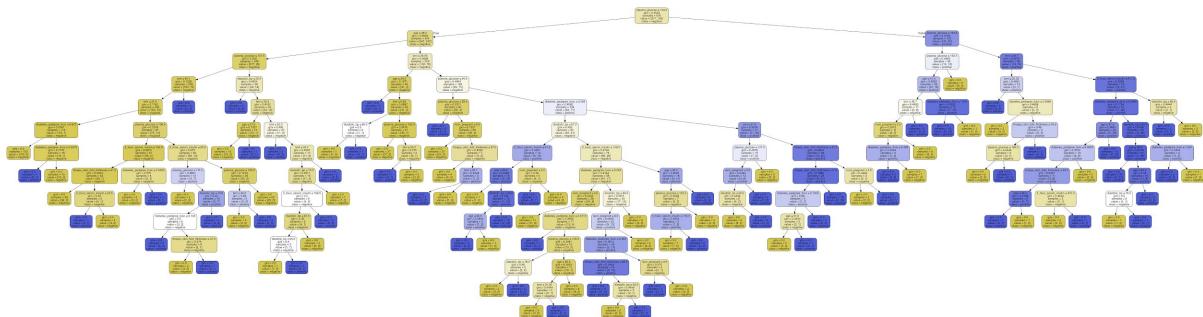
## Random Forest



36

## Random Forest

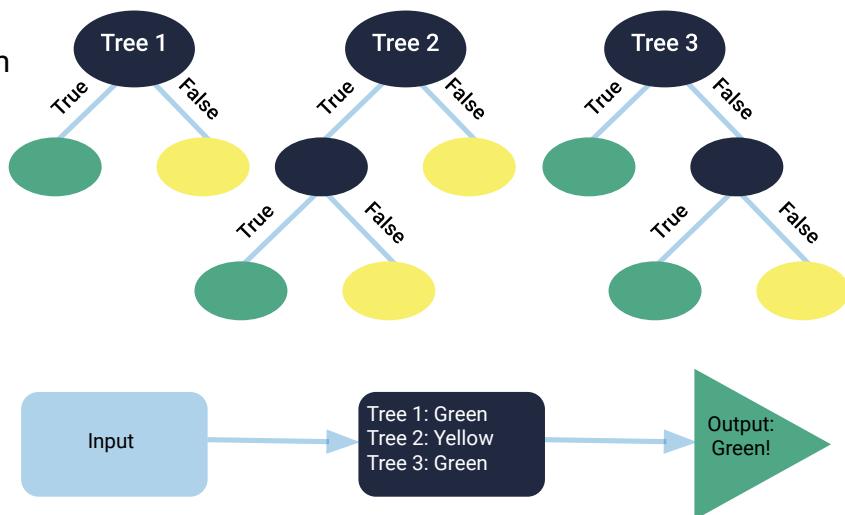
Instead of having single, complex tree like the ones created by decision trees, a random forest algorithm will sample the data and build several smaller, simpler decisions trees.



37

## A Forest of Trees

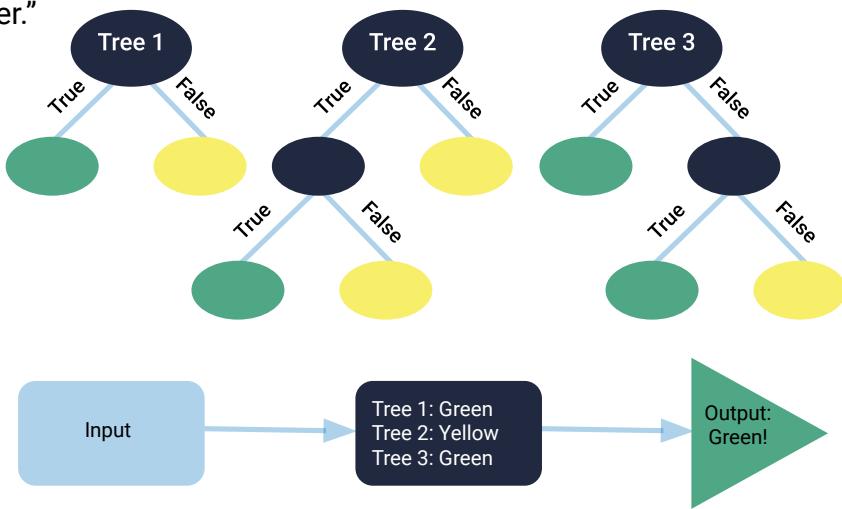
In a random forest, each tree is much simpler because it is built from a subset of the data.



38

## A Forest of Trees

Each tree is considered a “weak classifier” but when you combine them, they form a “strong classifier.”

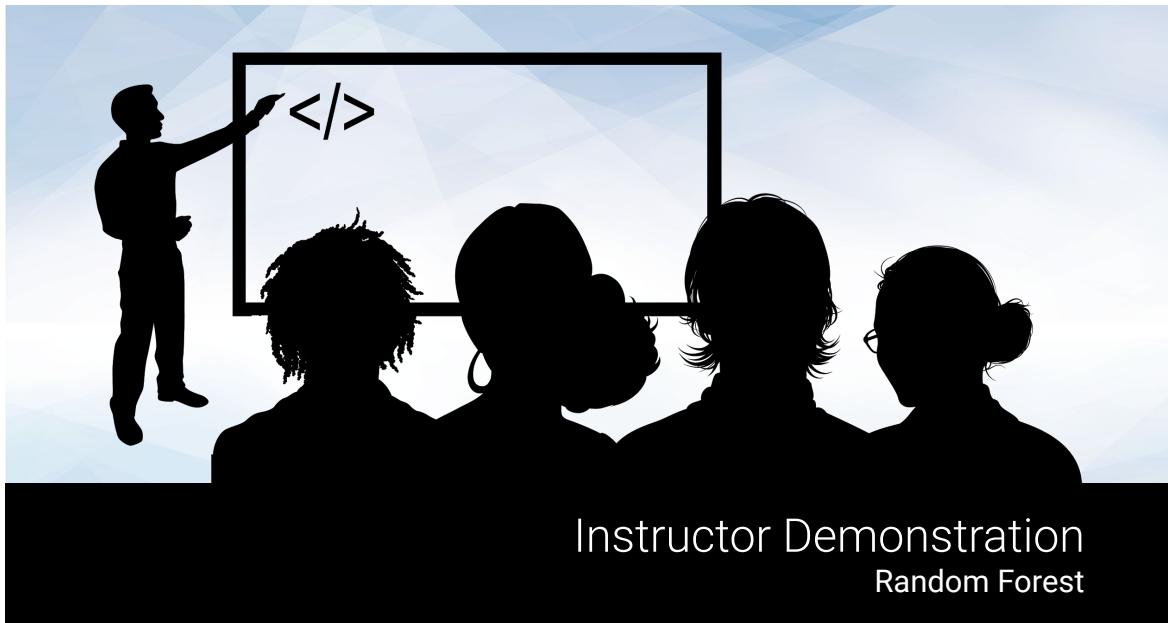


39

## Benefits of Random Forest Algorithm

- It's robust against overfitting.
- It can be used to rank the importance of input variables in a natural way.
- It can handle thousands of input variables without variable deletion.
- It's robust to outliers and non-linear data.
- It runs efficiently on large databases.

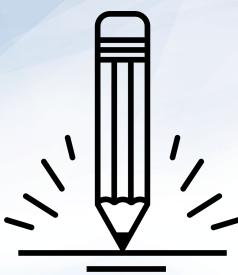
40



## Instructor Demonstration

### Random Forest

41



## Activity: Predicting Fraud with Random Forests

In this activity, you will explore how the random forest algorithm can be used to identify fraudulent loan applications. You will use the `sba_loans_encoded.csv` file that they created before to train the model.

Suggested Time:  
10 minutes



42



## Time's Up! Let's Review.

43

### Review: Predicting Fraud with Random Forests

---



Would you trust in this model to deploy a fraud detection solution in a bank?



What are your insights about the top 10 most importance features?

44

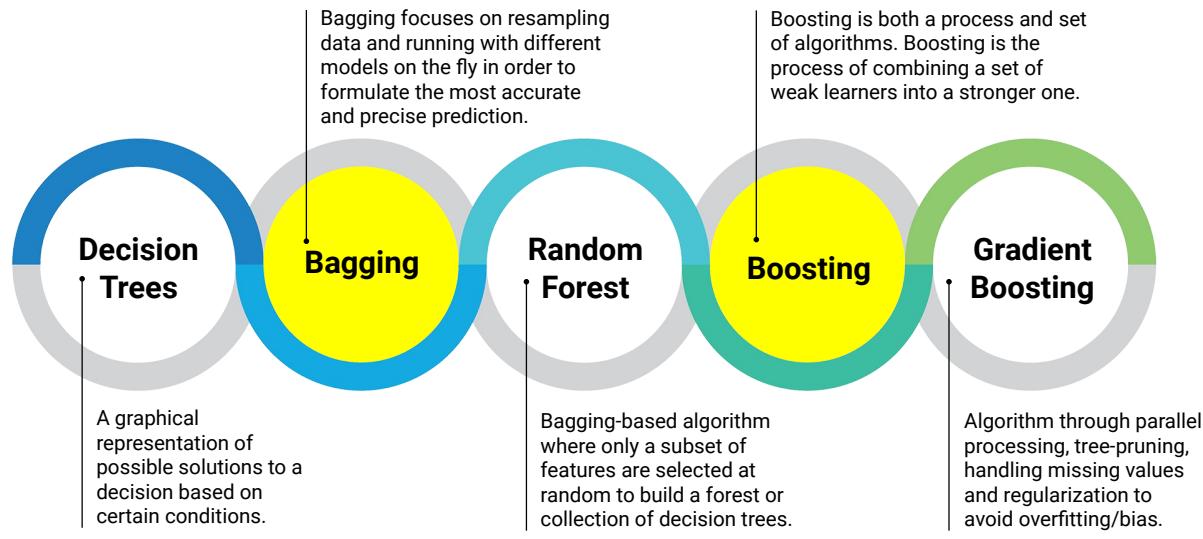


45

Boosting and Bagging

46

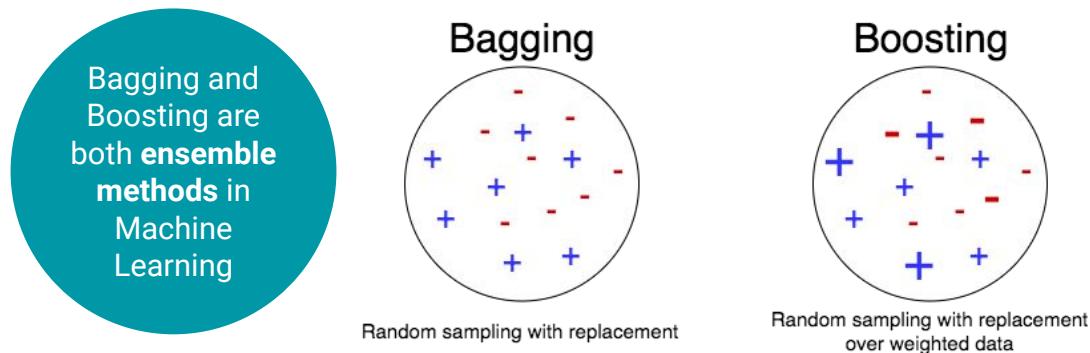
## Boosting and Bagging



47

## Boosting and Bagging

Boosting and bagging algorithms are used to improve the robustness and reliability of machine-learning models



48

## Boosting and Bagging

Boosting and bagging algorithms like XGBoost are often the best performing in Kaggle machine-learning contests. Their ability to make accurate predictions with precision and substantial recall is almost unparalleled

The screenshot shows a list of three competitions from Kaggle:

Competition	Description	Prize	Teams
Santander Value Prediction Challenge	Predict the value of transactions for potential customers.	\$60,000	4,477 teams
Two Sigma Financial Modeling Challenge	Can you uncover predictive value in an uncertain world?	\$100,000	2,070 teams
The Winton Stock Market Challenge	Join a multi-disciplinary team of research scientists	\$50,000	832 teams

49

## Boosting vs. Bagging

### Boosting

Boosting takes multiple algorithms and coordinates them as an ensemble and runs the algorithms in tandem to identify the best prediction

VS.

### Bagging

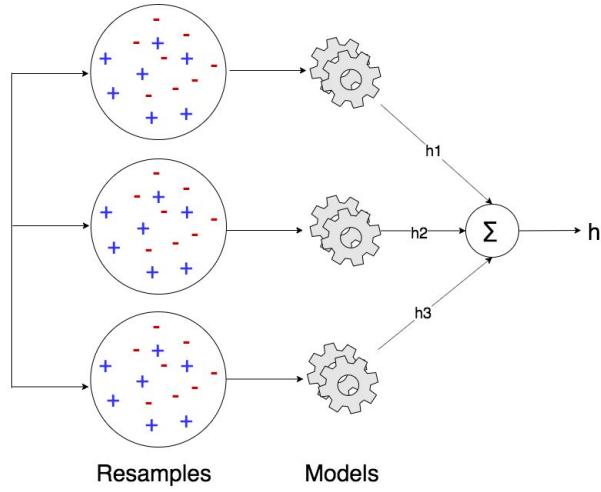
Bagging focuses on resampling data and running with different models on the fly to formulate the most accurate and precise prediction.

50

## Bagging

---

Bagging averages predictions from multiple samples and or models.

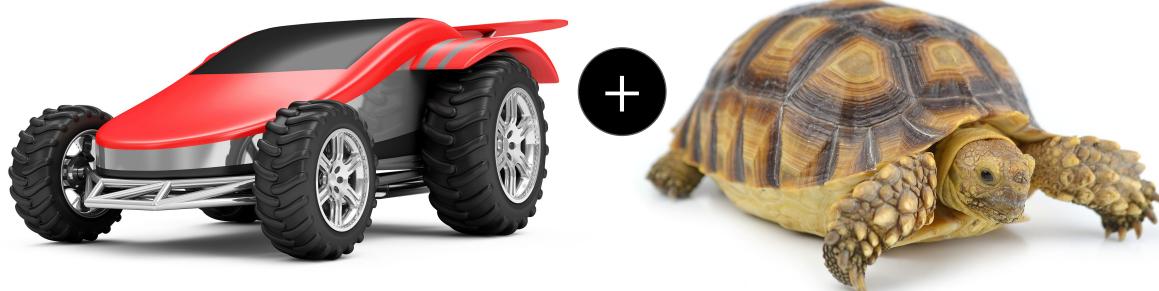


51

## Boosting

---

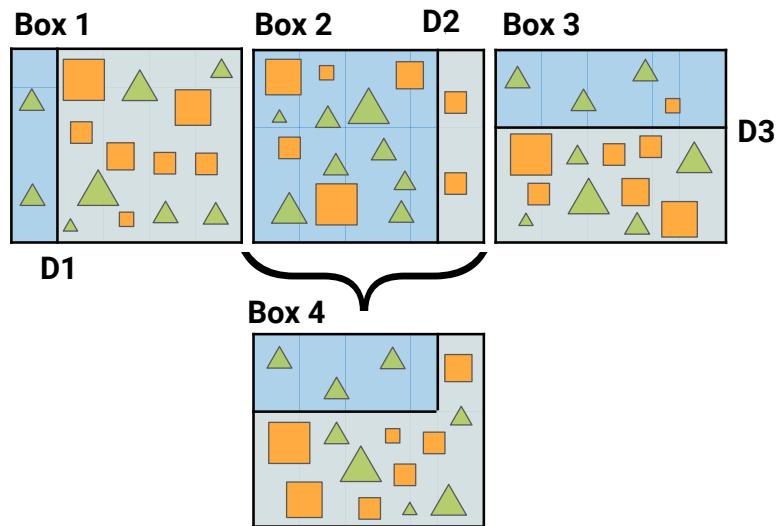
Boosting algorithms work iteratively, by sampling more heavily the observations with worst predictions. Subsequent weak learners then learn these more. Weak learners are then aggregated to produce a more accurate and precise prediction. The goal of a boosting algorithm is to combine weak learners into ensemble learners.



52

## Boosting

For this reason, boosting algorithms are considered meta-algorithms. Instead of working with and affecting data, boosting algorithms work with and affect other algorithms.

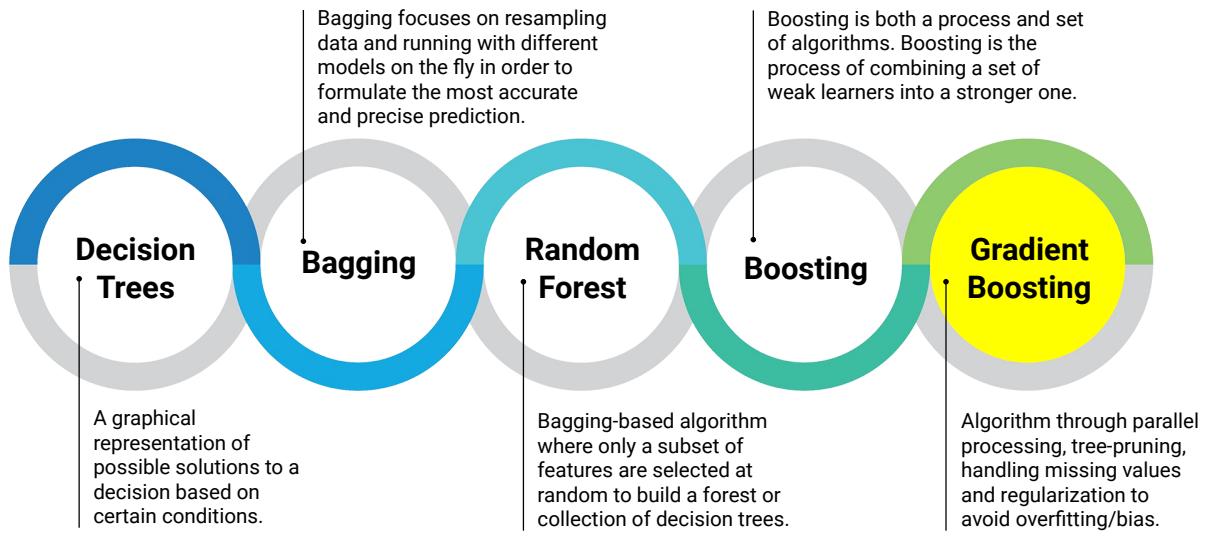


53

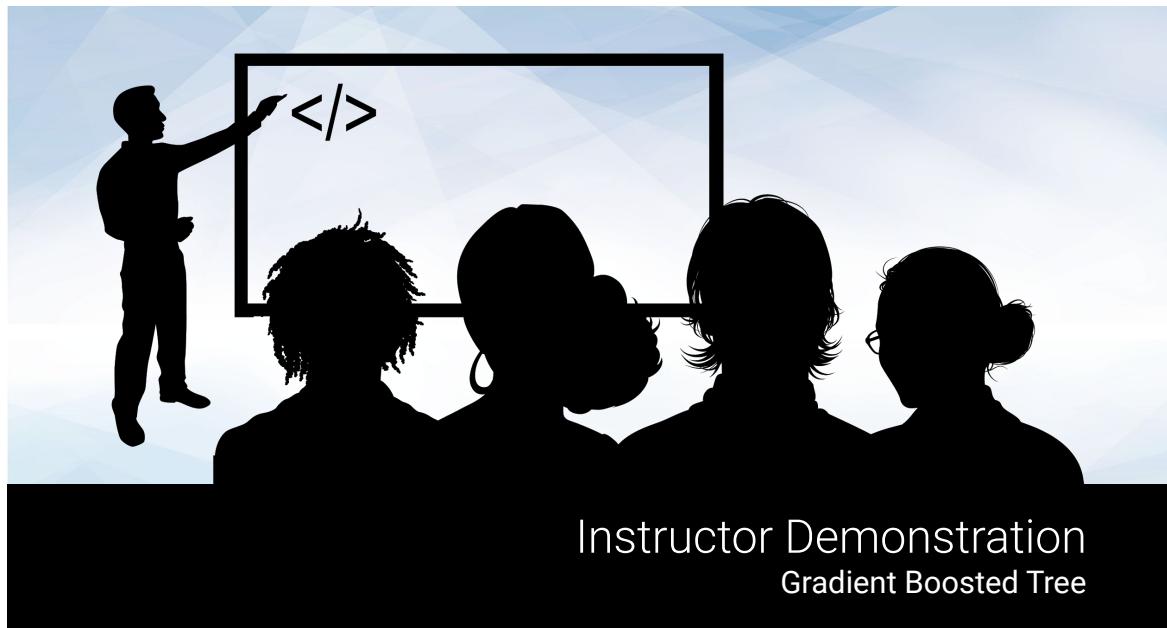
Gradient Boosted Tree

54

## Gradient Boosting



55



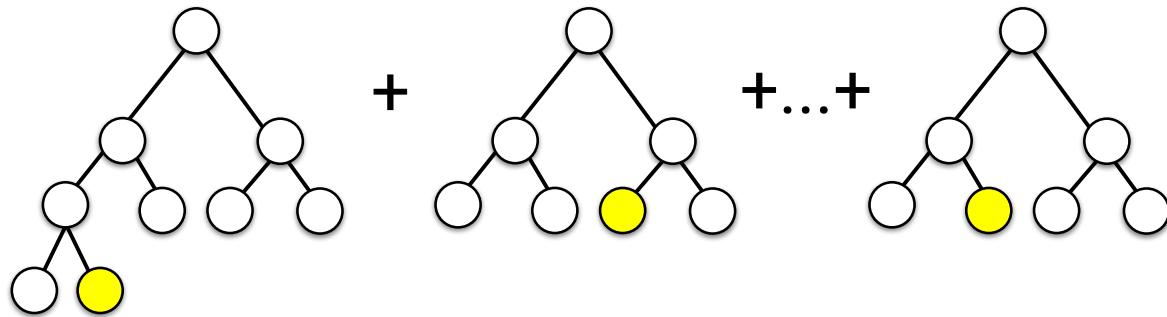
56

## Gradient Boosted Tree

---

Gradient Boosted Trees can be created using the GradientBoostingClassifier module from the ensemble package

```
From sklearn.ensemble import GradientBoostingClassifier
```



57

## Gradient Boosted Tree

---

### Arguments

GradientBoostingClassifier has three main arguments:

`N_estimators`

`Learning_rate`

`Max_depth`

### Definitions

The `n_estimators` parameter configures the number of weak learners being used with the boosting algorithm.

58

# Gradient Boosted Tree

---

## Learning\_rate

`Learning_rate` controls overfitting.  
Smaller values should be used when setting `learning_rate`.

## max\_depth

The `max_depth` argument identifies the size/depth of each decision tree being used.

---

59



## Activity: Turbo Boost

In this activity you will use the `sklearn` `GradientBoostingClassifier` boosting algorithm to detect fraudulent loan applications using ensemble learning.

Suggested Time:  
10 minutes



60



**Time's Up! Let's Review.**

61

The Trees vs. The World

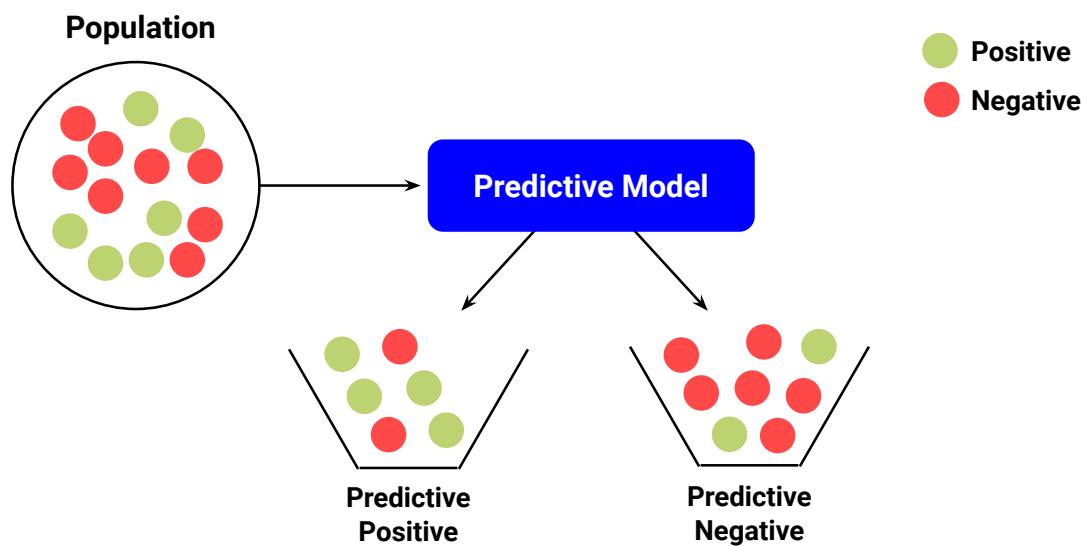
62

## Classification: A multidisciplinary challenge

<b>Finance and Banking</b>	Fraud detection, money laundering, credit risk assessment.
<b>Retail and Marketing</b>	Customized product offers, product recommendation, direct-marketing optimization.
<b>Politics</b>	Vote intention, party affinity.
<b>Health</b>	Trials tests, ills diagnosis.
<b>Security</b>	Intruders detection, predictive maintenance.
<b>Education</b>	Programs affinity, customized curricula, desertion prevention.

63

## Classification: A multidisciplinary challenge



64



## Are tree-based algorithms the strongest for classification?

65

### Tree-based algorithms



Are easy to represent, making a complex model much easier to interpret.



Can be used for any type of data: Numerical (e.g., loan's amount) or categorical (e.g., name of bank that issues a loan).



Require little data preparation.



Can handle data that are not normally distributed.



Can avoid overfitting.

66

## Trees vs. Classical Classifiers

---



Generally speaking, classical classifiers may be faster.



Logistic regression may outperform decision trees or random forests having a large number of features with low noise.



SVM also support linear and non-linear models.



SVM handles outliers better.



KNN naturally supports incremental learning (data streams).

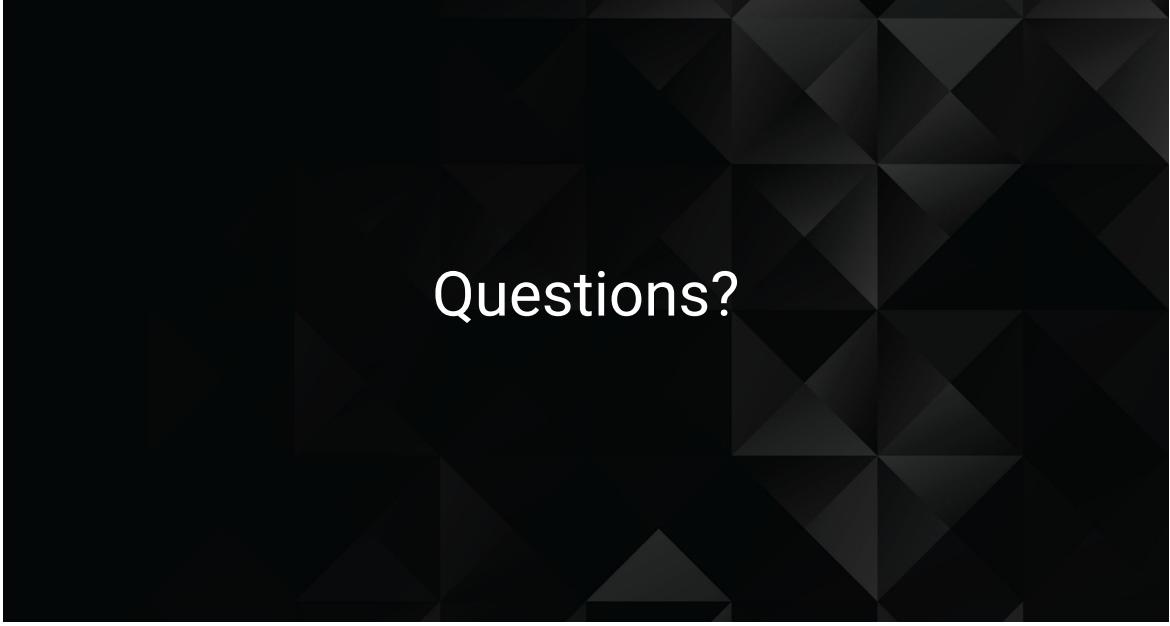
---

67



**Which algorithm should  
I use for classification?**

68



Questions?