

# Boosting a Signal

Using Trees and Forests for Classification Tasks

Evan Simpson

[p.evansimpson@gmail.com](mailto:p.evansimpson@gmail.com)

[EducatorsRLearners](#)

# Agenda

- Introduction
- Review of Machine Learning Types
- Decision Trees
- Random Forest
- XGBoost
- Future Directions

A stylized illustration of a person from the chest up, wearing a dark blue suit jacket, a white dress shirt, and a red tie. The person's head is not visible. On the right side of the chest, there is a white rectangular name tag with a yellow clip at the top. The background is a solid yellow color.

**HELLO,**  
My Name Is..





**HELLO,**  
**My Name Is..**

# Activity Time!

Match the terms with their definitions

Check your Answers

F, D, C, B, E, A

# Types of Machine Learning

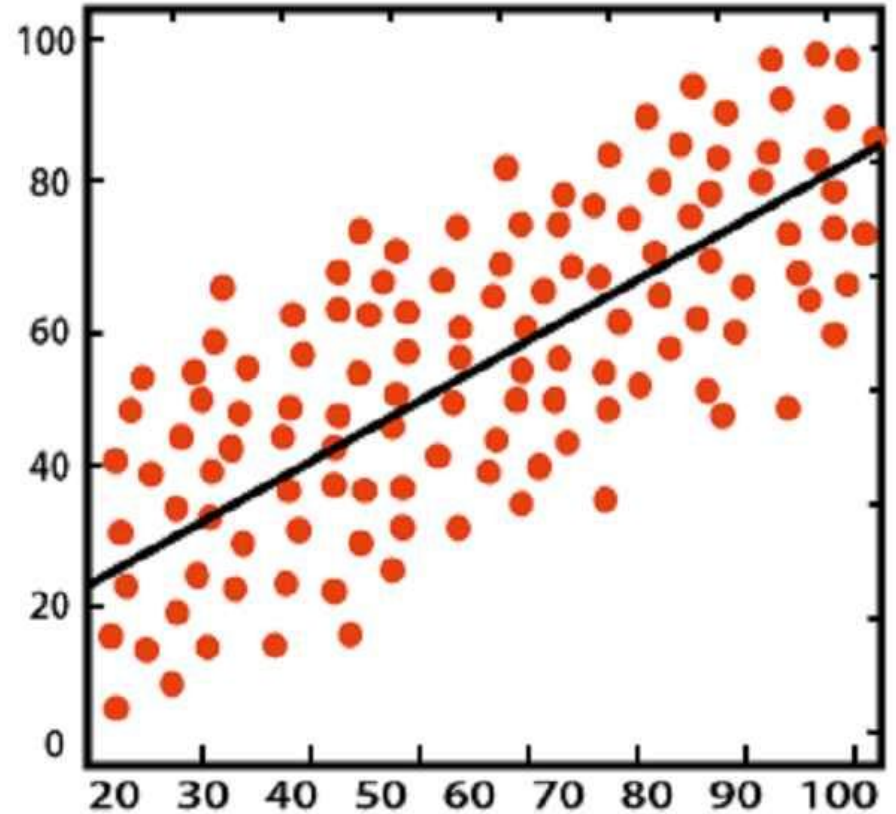
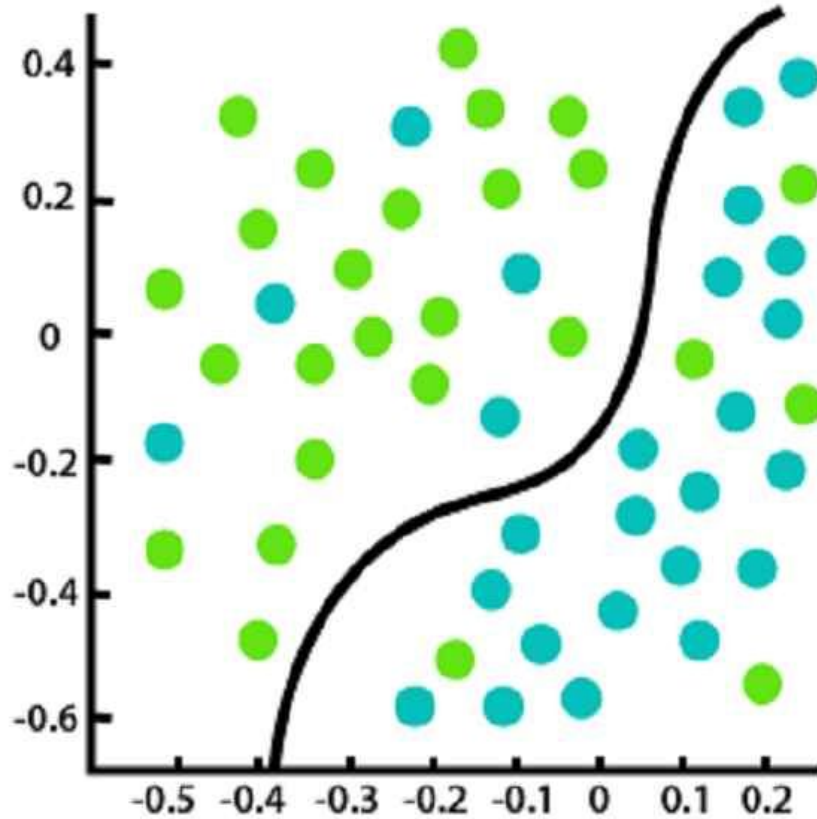


- Unsupervised?
- Supervised?

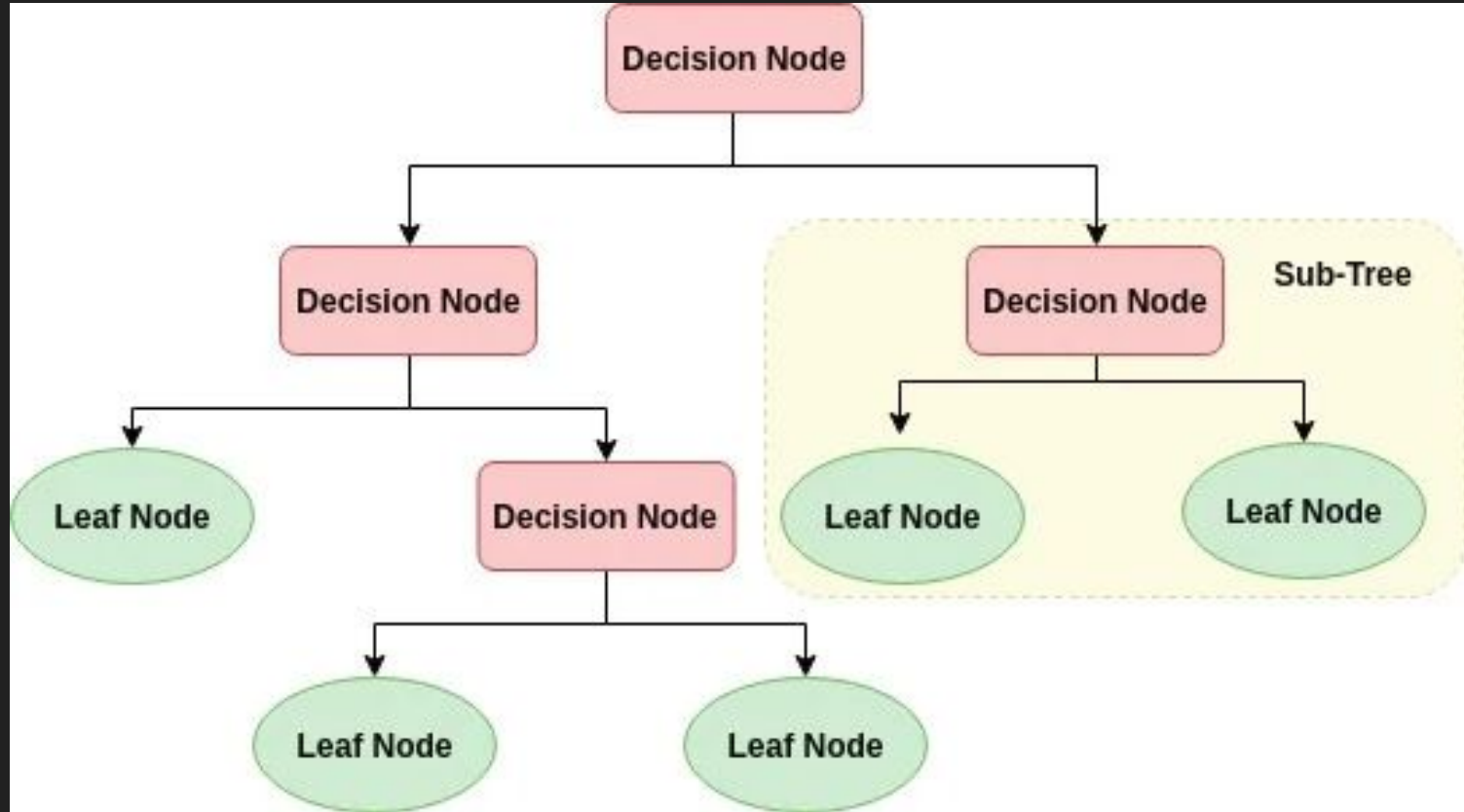
Why?



# Types of Machine Learning



# Decision Trees

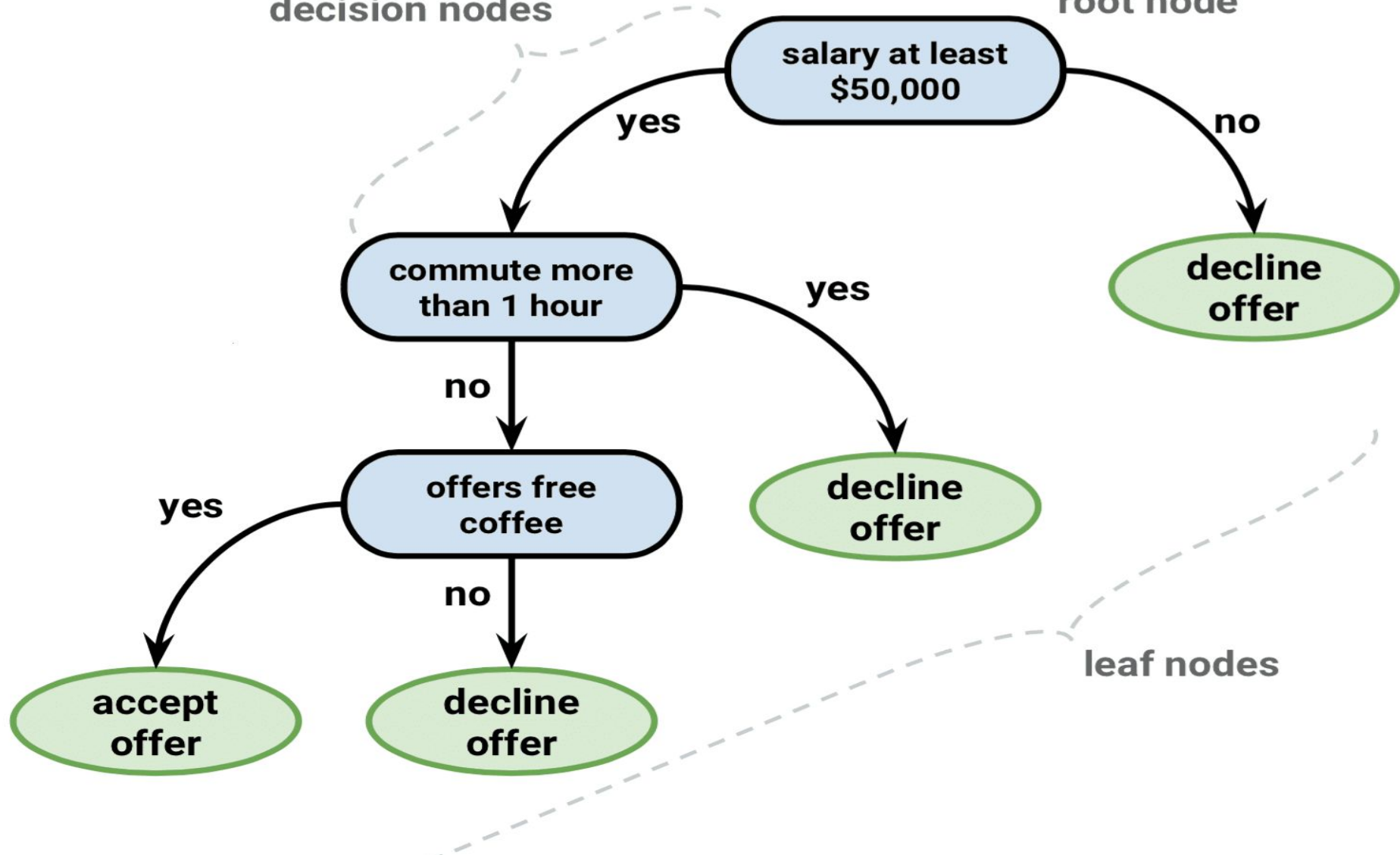


# Decision Trees

Example: Job Offer

decision nodes

root node



# Decision Trees

## Advantage

- Easy to interpret and use

## Disadvantage

- Easy to overfit

# Ensemble Methods

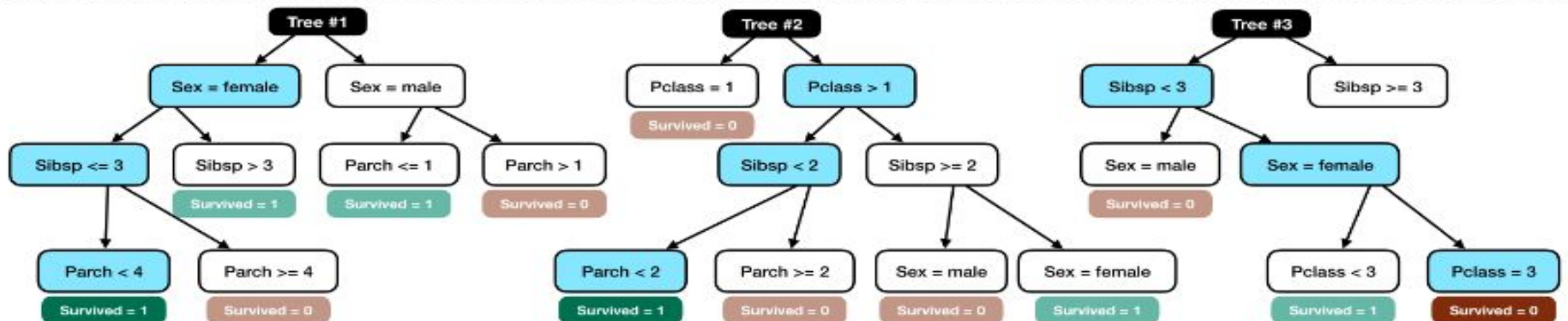




# Random Forest

Did the passenger survive?

PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
893	3	Wilkes, Mrs. James (Ellen Needs)	female	47	1	0	363272	7		S



Tree #1 votes Survived = 1

Tree #2 votes Survived = 1

Tree #3 votes Survived = 0



Random forest predicts Survived = 1

# Random Forest

- + Incredibly Powerful

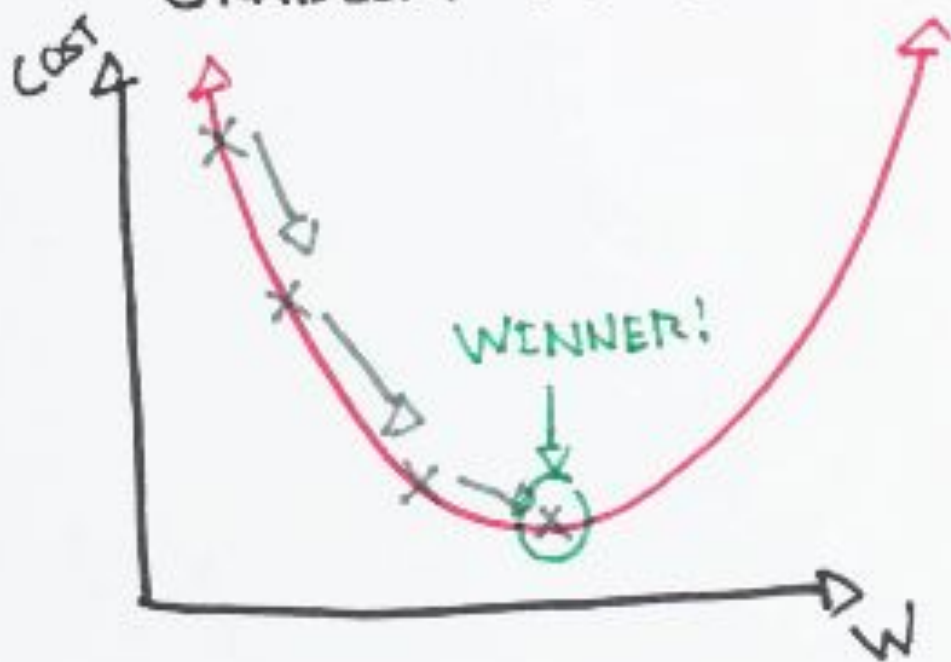
- Trees are built randomly



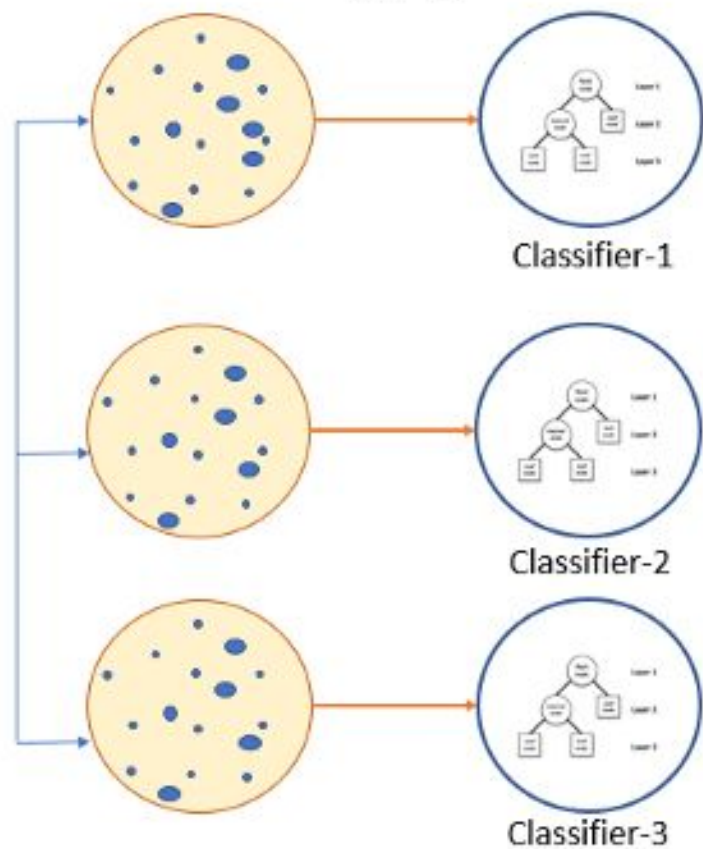
*dmlc*

***XGBoost***

# GRADIENT DESCENT

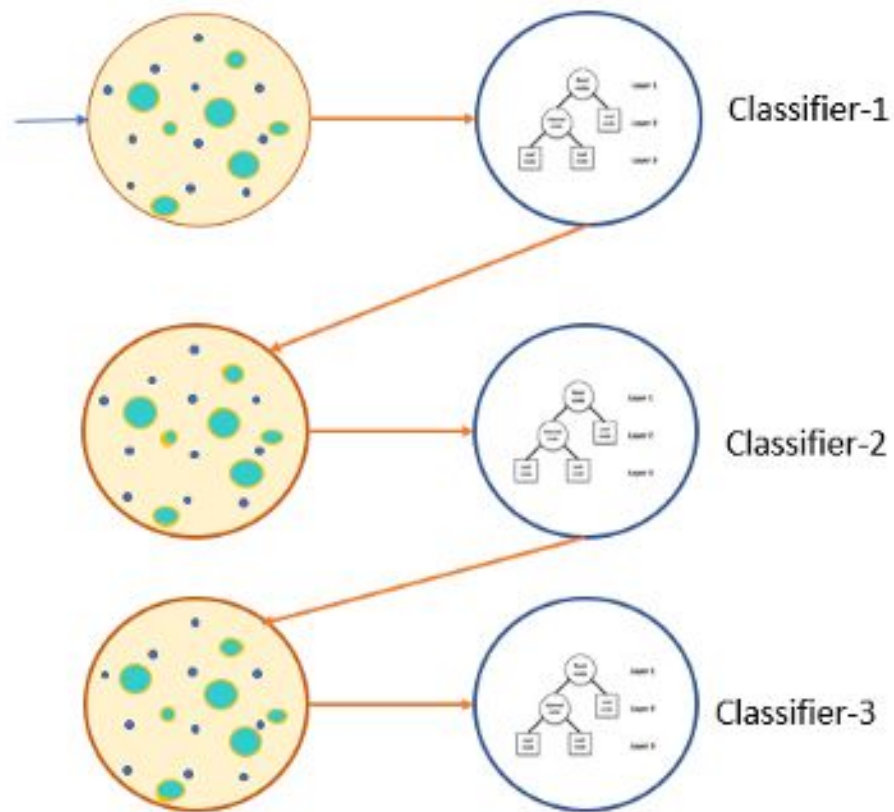


## Bagging



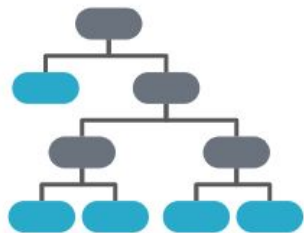
*Parallel*

## Boosting



*Sequential*

First Estimator



Test



Boosting round #1

Second Estimator

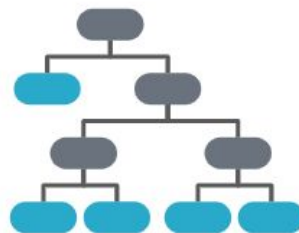


Test



Boosting round #2

Third Estimator



Test



Boosting round #3

Fourth Estimator



Boosting round #4

Train

Train

Train

# Activity Time Part 2:



survived	0,1
pclass	1, 2, 3
sibsp	# of siblings
parch	# of parents
fare	price of ticket
is_female	yes or no

# Sample Workflow

Frame the problem

Install and Load the Libraries

Collect/Load the Data

Conduct Exploratory Data Analysis (EDA)

Prepare the Data

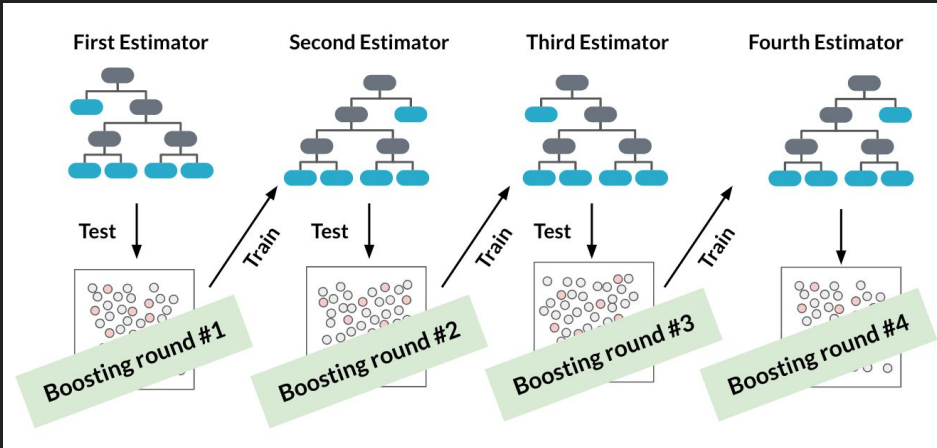
Build and Evaluate the Model

Inspect the Feature Importance

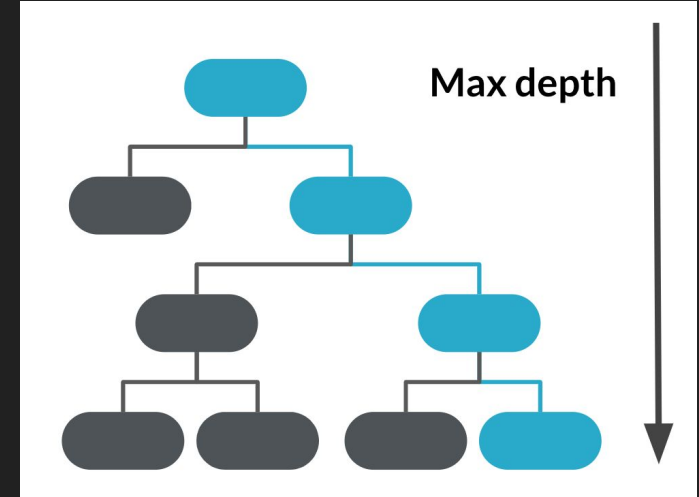




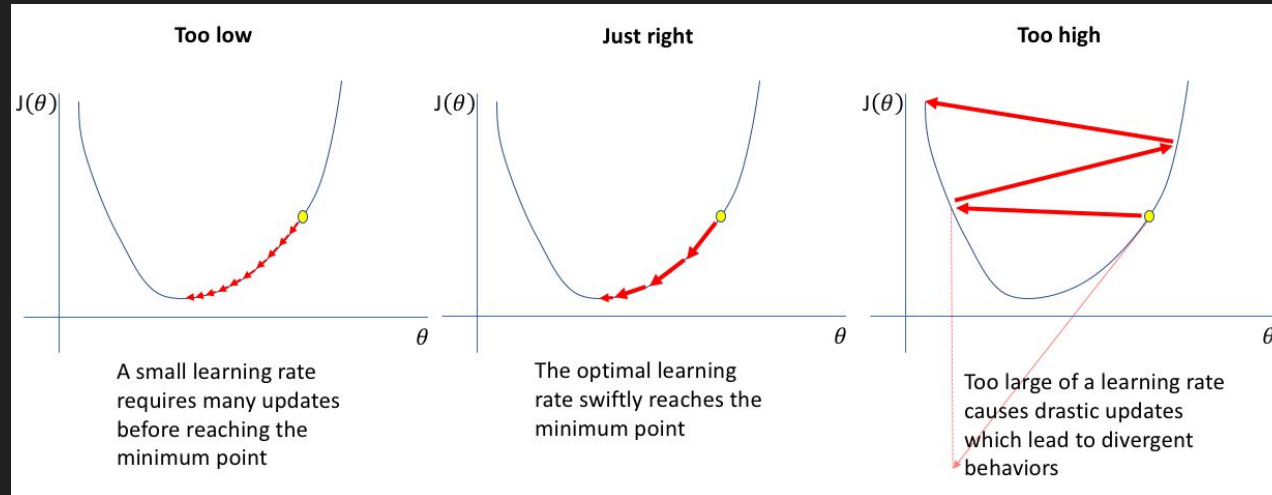
## n\_estimators



## max\_depth



## eta (aka, "learning rate")



## early\_stopping





# Next Steps



O P T U N A



LightGBM



CatBoost

# Resources

[Kaggle Tutorial using Random Forest](#)

[DataCamp Decision Tree Tutorial](#)

[DataCamp XGBoost Tutorial](#)

[Hands On Machine Learning Chapters 6 and 7](#)

[Numsense! Data Science for the Layman](#)

[Medium Reading list based on this topic](#)

[YouTube Playlist](#)

[Repo for this Lesson](#)

