

Taxi Fare Prediction

September 23, 2021



Modeling Taxi Demand

September 23, 2021



PREDICT TAXI DEMAND

What's in it?

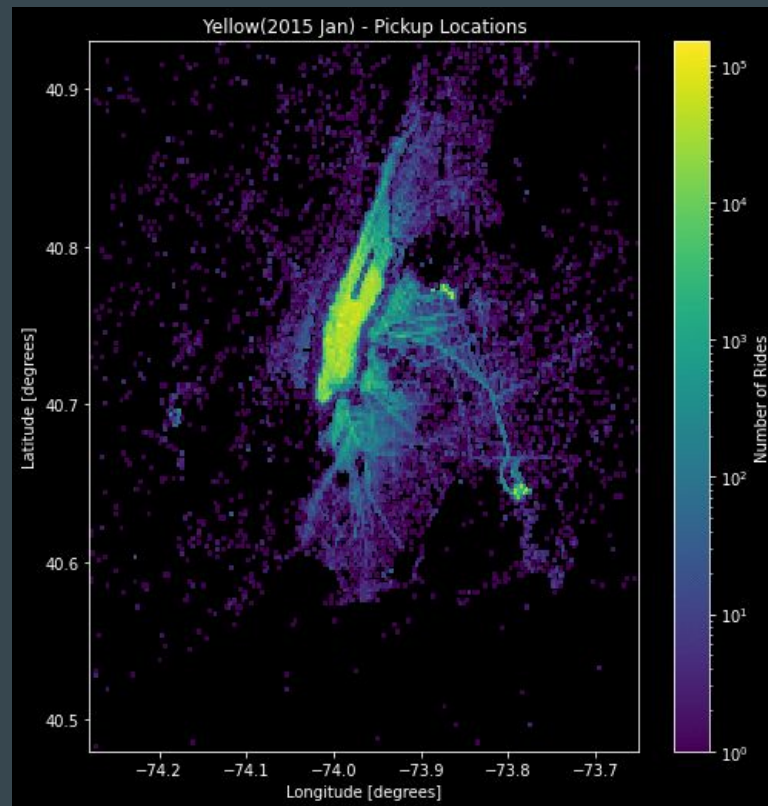
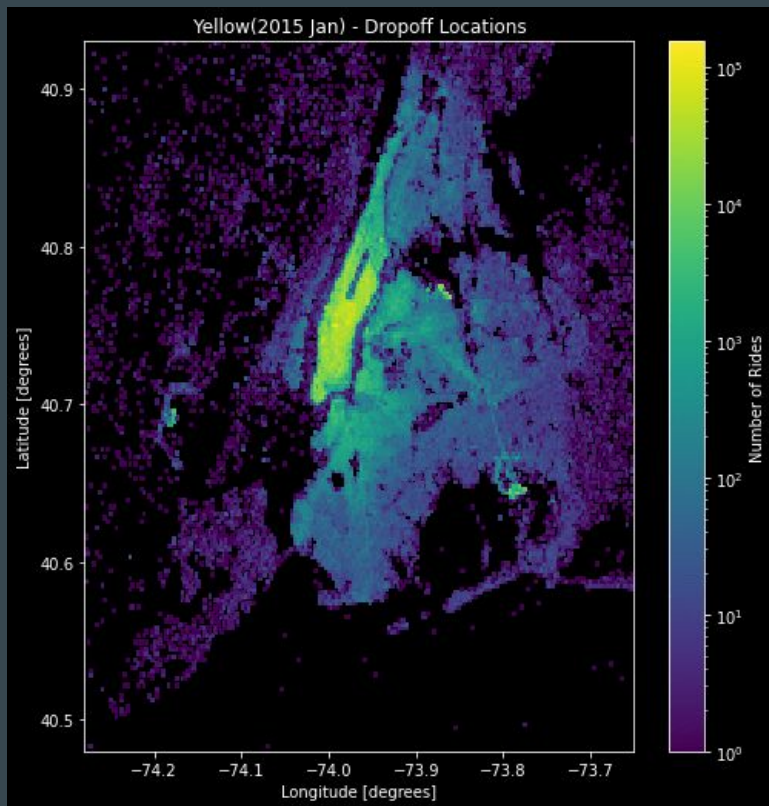
Company:

- Predicting demand and deploy taxis accordingly there by efficiently making use of resources
- Use the demand prediction model to offer discounts or surcharges, thereby increasing the revenue.

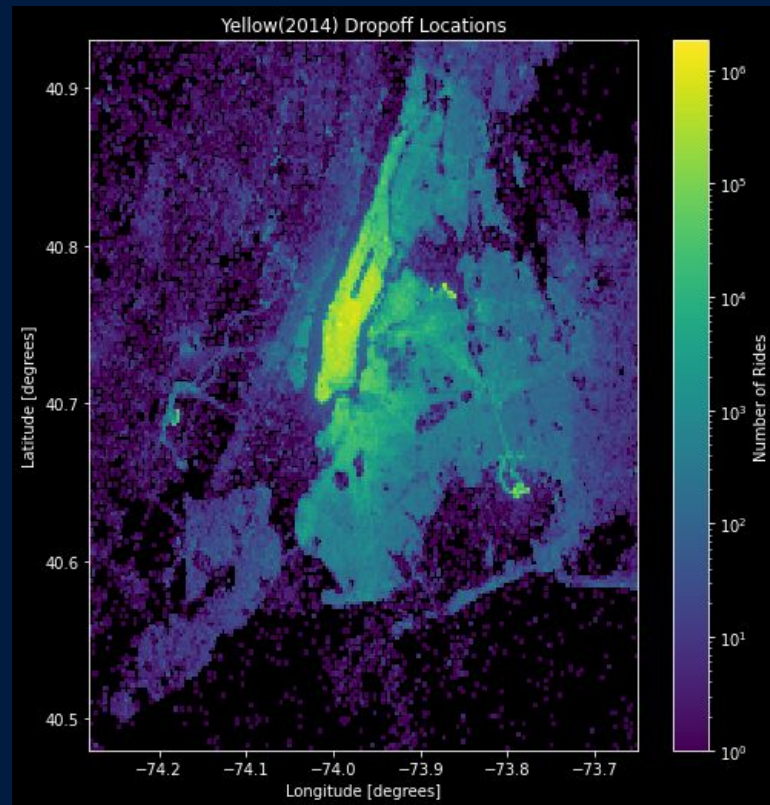
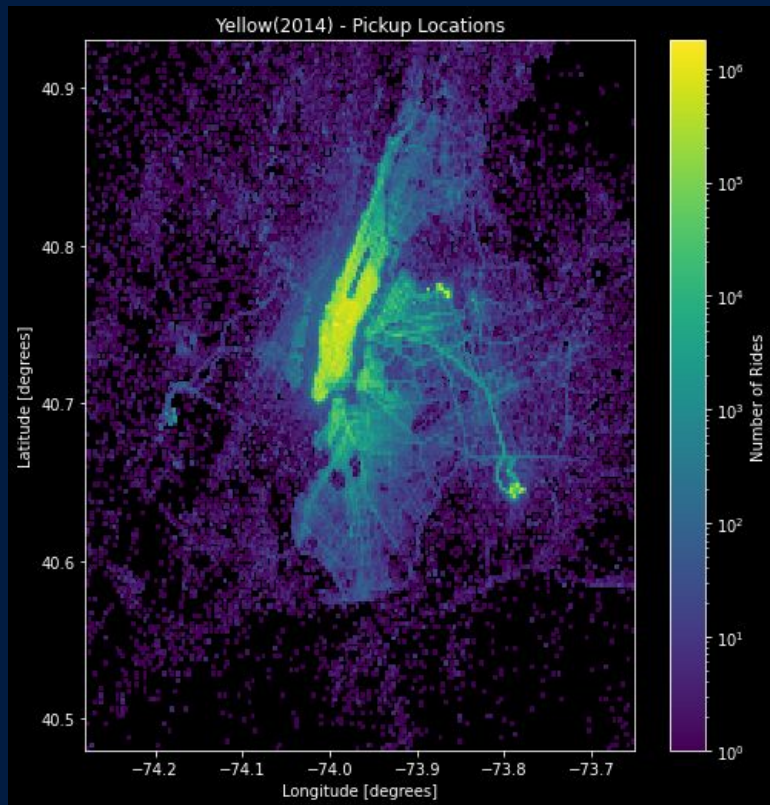


YELLOW TAXI

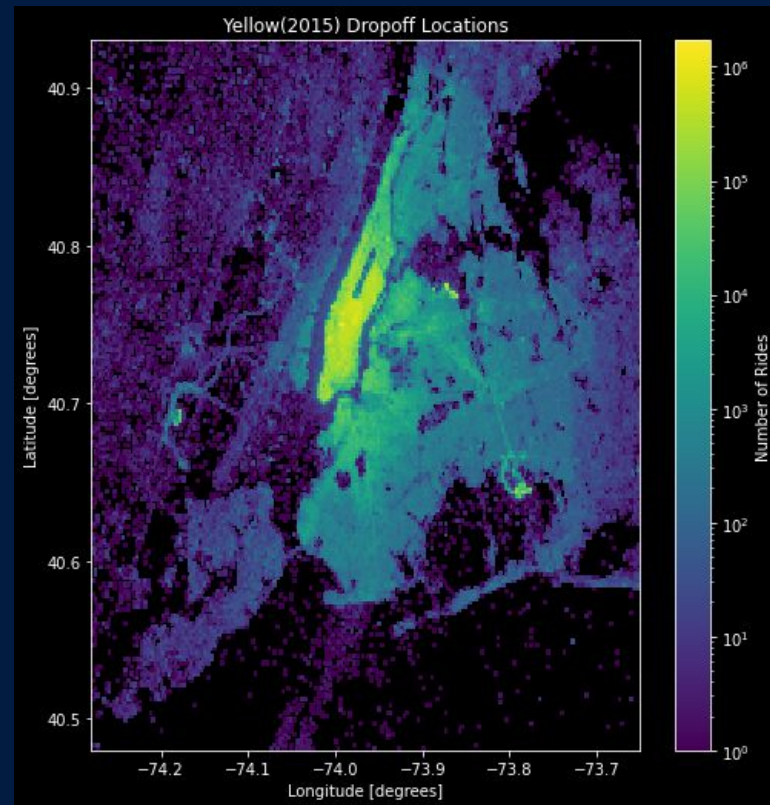
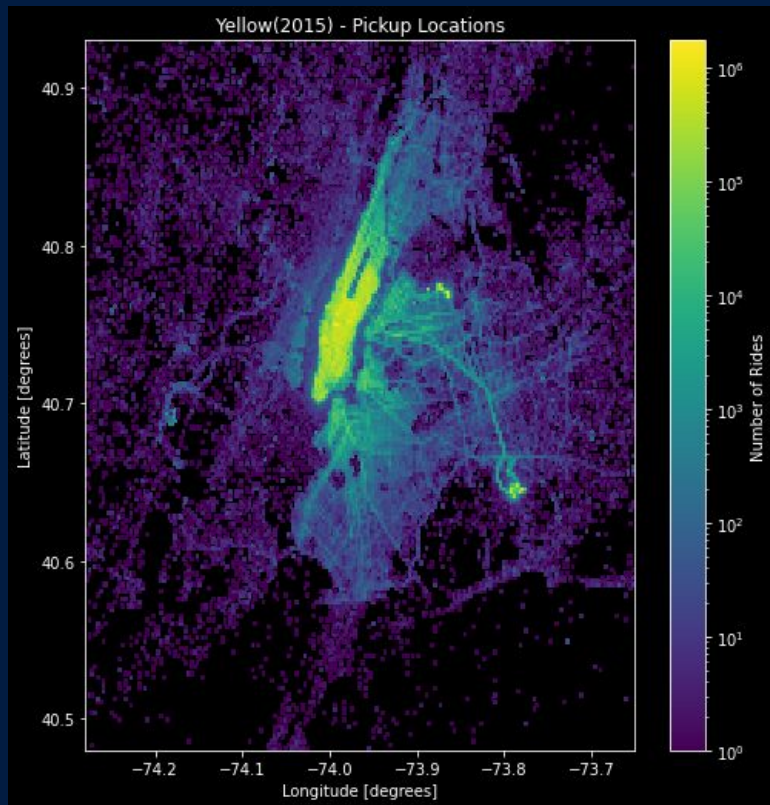
2015 - Pickup/Drops By Months



2014

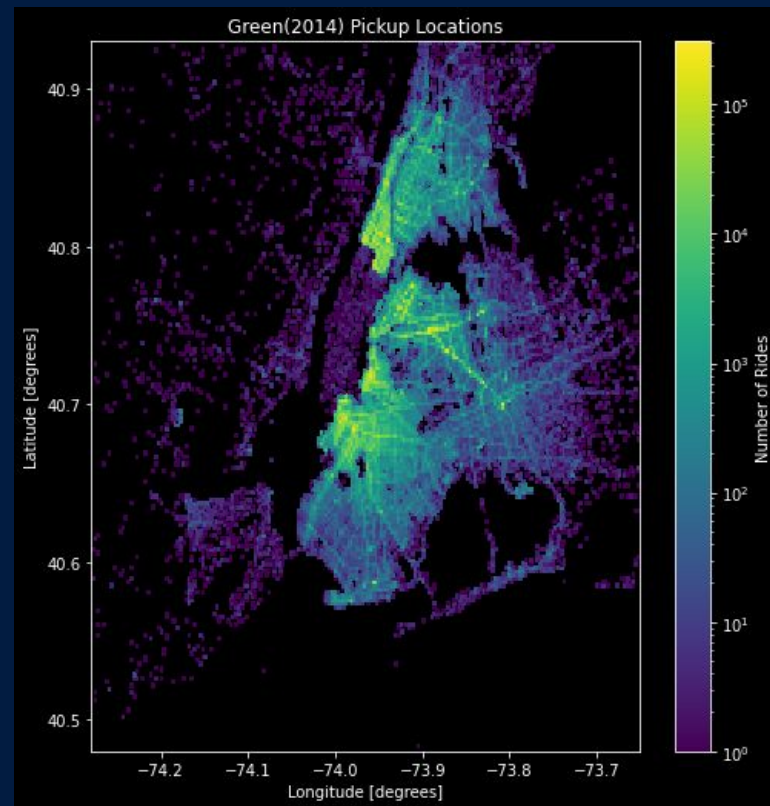
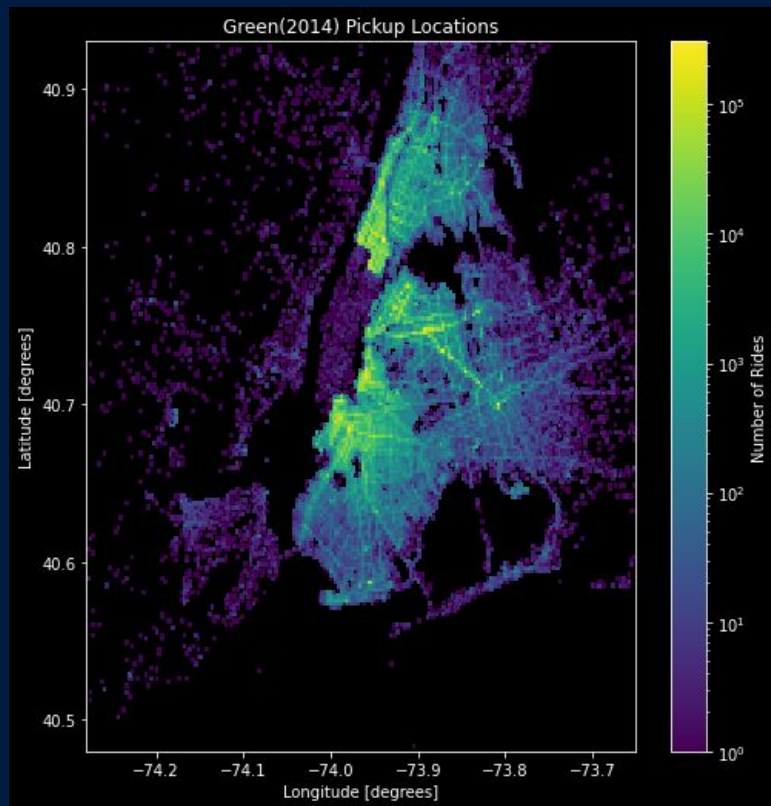


2015

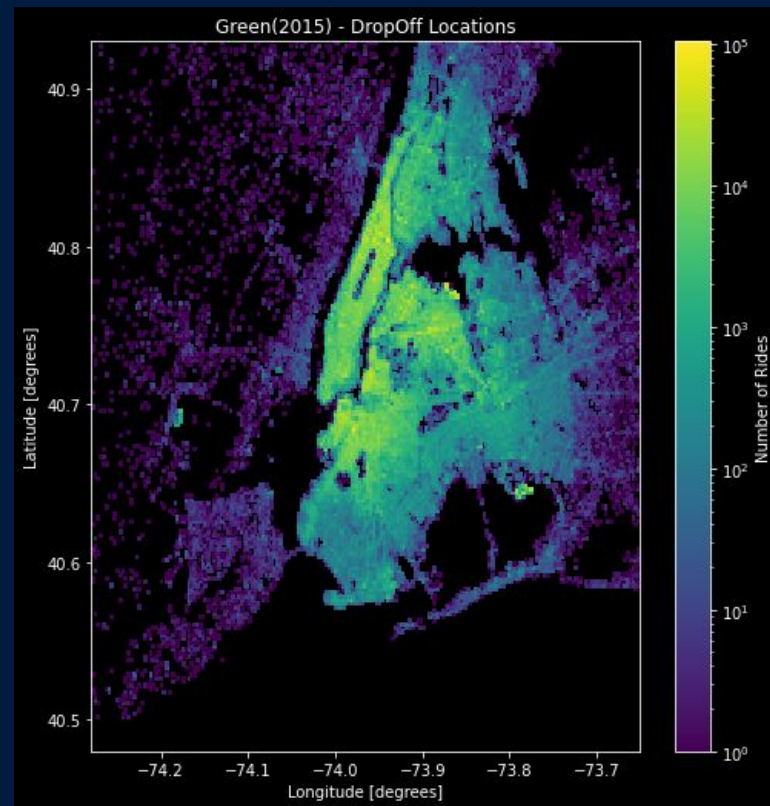
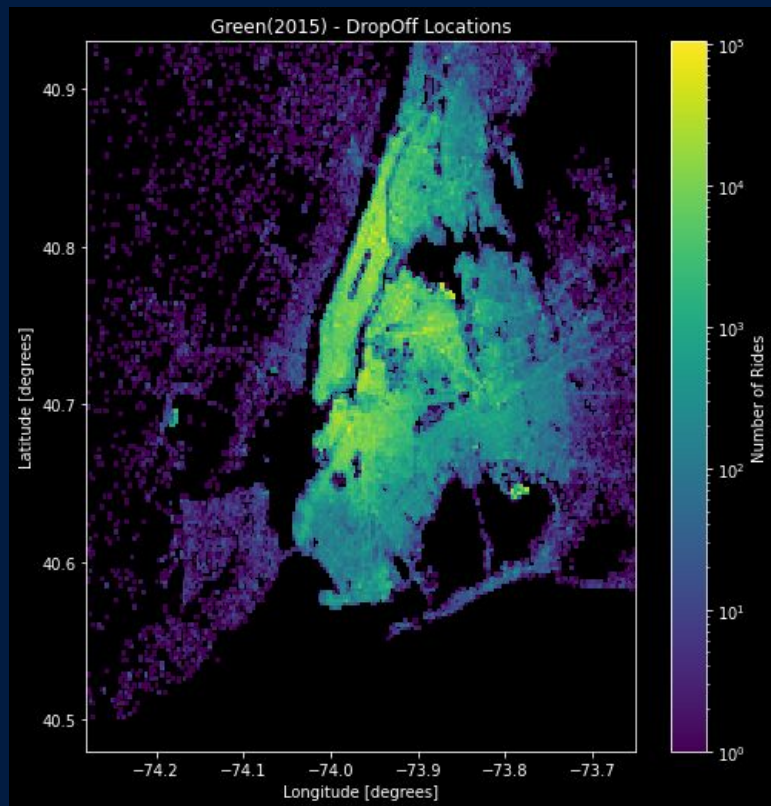


GREEN TAXI

2014



2015



What are the options?

Probabilistic Models

- Linear regression
- Multiple Linear Regression.

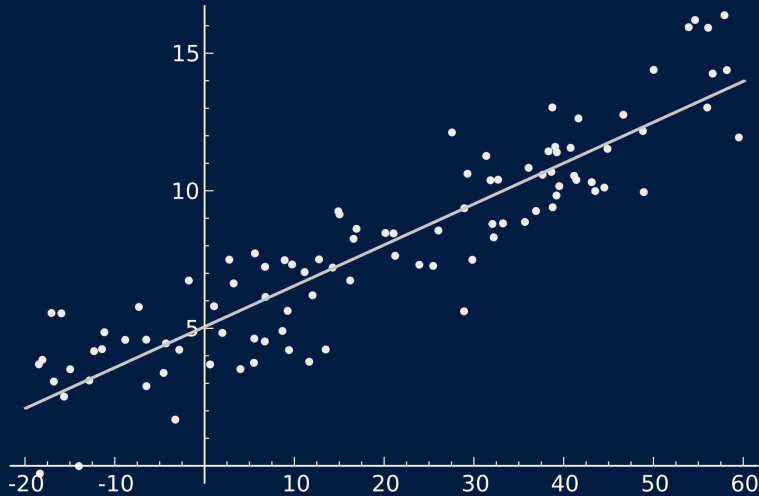
Machine Learning Models

- Artificial Neural Networks
- Decision Tree
- Clustering
- Ensemble Methods

Time Series Models

- Autoregressive(AR)
- Vector AR
- Moving Average(MA)
- ARIMA
- STARIMA

Multiple Regression Model



$$Y = \sum_{i=0}^n \beta_i X_i + \epsilon$$

Use MLE(Least Squares Estimation) to determine the coefficients for the explanatory variables.

Step 1: Check Correlation Coefficients

Strong correlation b/w X_i and Y indicates X_i is important
Strong correlation b/w X_i and X_j - multicollinearity

Step 2: Stepwise Selection

Step 3: Best Subsets Regression



$$y = b + W^T X$$
$$f(x) = \sigma(b + W^T X)$$

$$W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

ARTIFICIAL NEURAL NETWORK

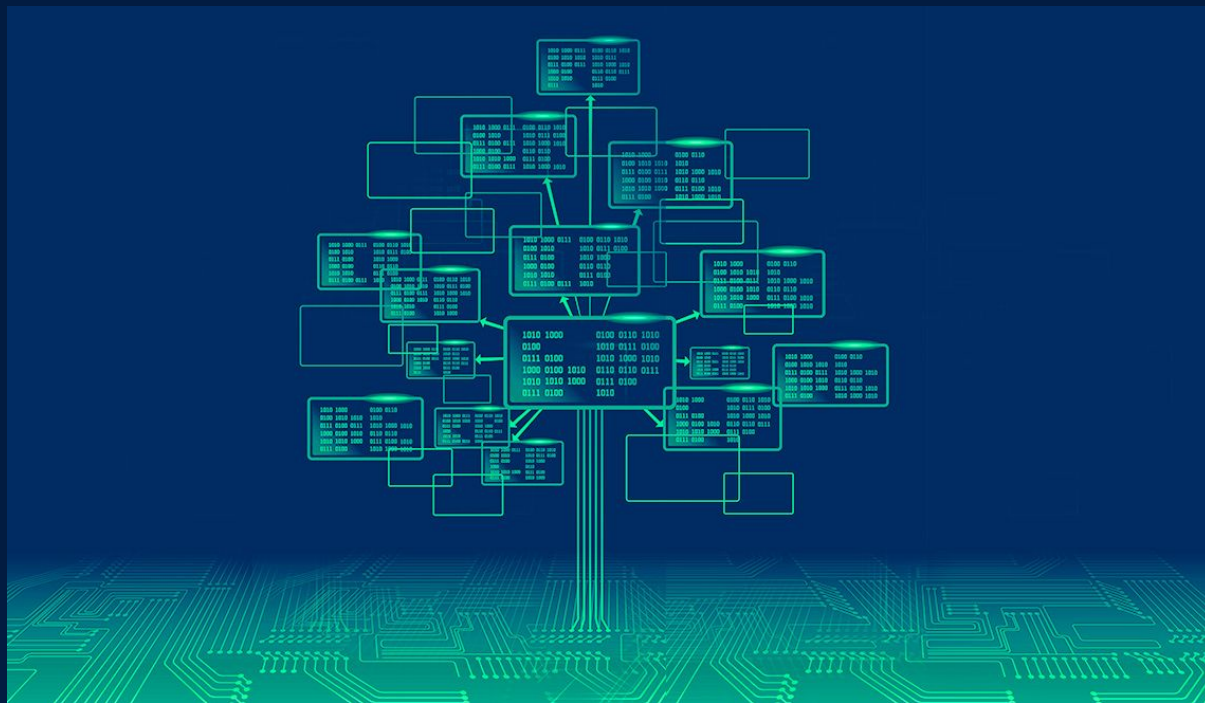
Pros:

- Fast Predictions
- Good with non-linear data
- Can be used for both regression and classification

Cons:

- Computationally intensive
- Needs lot of training data
- Overfitting and generalization

Decision Trees



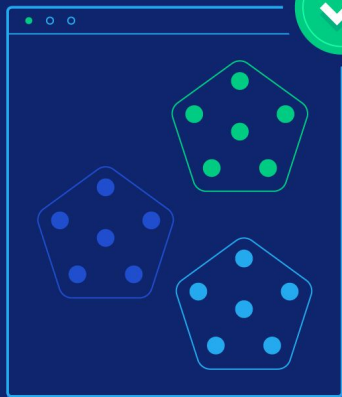
Pros:

- Understandability
- Resistant to Outliers

Cons:

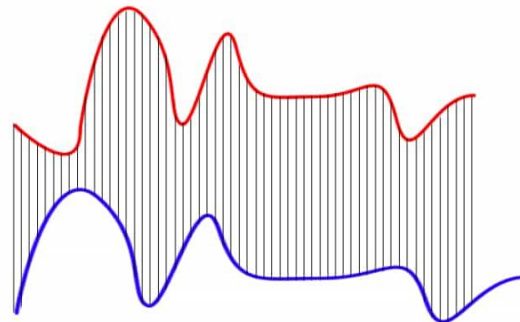
- Prone to overfitting
- Needs careful Parameter tuning
- Biased Tree can be created if there is an imbalance

CLUSTERING WITH DYNAMIC TIME WARPING

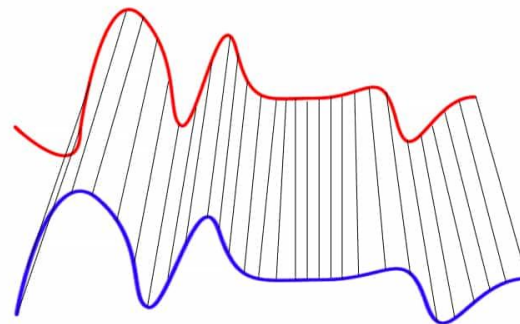


Cons:

- Involves transformation time-series data, which might lead to some loss of information
- Incorrect Parameter setting due to assumptions
- Slow processing time

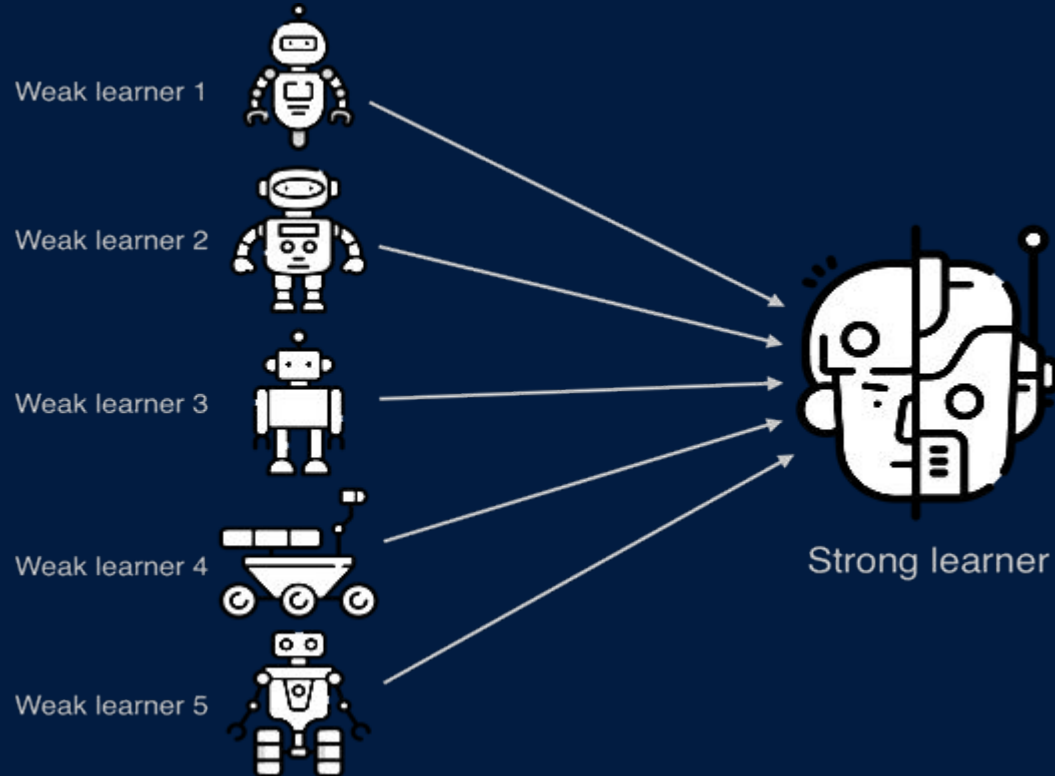


Euclidean Matching



Dynamic Time Warping Matching

ENSEMBLE



TIME SERIES

What kind of Forecasting?

- Point forecast
- Interval forecast
- Density forecast

What kind of data needed?

- Data collected at a single point of time
- Observations of data made over time

Identify the patterns in data:

- Horizontal (Stationary)
- Trend
- Seasonal
- Cyclical

