

TUNKU ABDUL RAHMAN UNIVERSITY OF MANAGEMENT AND TECHNOLOGY

FACULTY OF COMPUTING AND INFORMATION TECHNOLOGY

Forecast Delay Durations For New York Bus

BMCS2114 MACHINE LEARNING 2022/2023

Student's name/ ID Number : Teng Kai Deng 21WMR02977

Student's name/ ID Number : Sit Yie Sian 21WMR03693

Student's name/ ID Number : Leong Sheng Mou 21WMR07568

Programme : RDS

Tutorial Group : 2

Tutor's name : Dr. Lim Siew Mooi

TITLE: Forecast Delay Durations For New York Bus

Table of Contents

TITLE: Forecast Delay Durations For New York Bus	2
ABSTRACT	3
1.0 INTRODUCTION	4
1.1 Problem Statement	5
1.2 Solution	6
1.3 Related Work	7
1.3.1 Gradient Boosting Regression Tree (GBRT)	7
1.3.2 Multi-layer Perceptron	7
1.3.3 Neural network algorithms	8
1.3.4 XGboost	8
1.3.5 ADA boosting	8
1.4 Algorithm	9
1.4.1 Gradient Boosted Trees (GBT)	9
1.4.2 The Multi-layer Perceptron (MLP)	9
1.4.3 Neural Network Algorithms	9
1.4.4 XGBoost (eXtreme Gradient Boosting)	10
1.4.5 AdaBoost (Adaptive Boosting)	10
2.1 Data Source	
2.2 Data Understanding	
2.3 Data cleaning	
3.0 Algorithm	
3.1 Feature 1 Algorithm	38
3.1.1 Gradient Boosted Tree (GBT)	
3.1.2 Multi-layer Perceptron (MLP)	39
3.1.3 Neural Network	40
3.1.4 XGBoost	
3.1.5 ADA Boosting	
3.1.6 Evaluate Results for Feature 1 using Barchart	
3.2 Feature 2 Algorithm	47
3.2.1 Gradient Boosted Tree (GBT)	47
3.2.2 Multi-layer Perceptron (MLP)	48
3.2.3 Neural Network	49
3.2.4 XGBoost	
3.2.5 ADA Boosting	
3.2.6 Evaluate Results for Feature 2 using Barchart	52
3.3 Deployment of Model	
Reference	56

TITLE: Forecast Delay Durations For New York Bus

Dataset:

https://www.kaggle.com/datasets/new-york-city/ny-bus-breakdown-and-delays?select=bus-breakdown

ABSTRACT

This project aims to forecast delay durations for various events taking place in the city based on the dataset maintained by the City of New York, United State. The dataset contains a number of information, including the bus route, the breakdown's location, the time of day, and the cause of the delay. The dataset from Kaggle will be used and some related works will be studied in order to make progress to the correct path. Many machine learning techniques will be used to evaluate the dataset and find patterns and trends in the data in order to produce accurate predictions. The purpose of this project is to develop models that can forecast delay durations for various events and to offer insights into the elements that affect delays in the city. The project's findings may help city planners, event planners, and transit authorities among others make decisions. The goal of this research is to create precise prediction models that transportation authorities can use to guide their decision-making, manage, and reduce delays faced by buses in New York City. The outcomes of this study could be applied to improve passengers' overall transportation convenience and increase the effectiveness of bus services.

Keyword: Bus arrival time prediction; Passengers' waiting time; forecast delay durations; machine learning

1.0 INTRODUCTION

Transportation has been a very convenient tool for humans to travel around since the past. There are various kinds of transportation used in this society to transport goods and people around the world such as cars, motorcycles, planes, and even horses. A very well known transportations among commoners is bus service. Its adaptability, expandability, and affordability account for its popularity. The planning process for bus services can be lengthy and complicated, thus it is frequently carried out with the use of computer software that can test millions of possible route configurations, service schedules, and employee schedules. Traveling by bus is very common everywhere these days as it is a public transport which often costs around RM1 for people to travel from one destination to another within Malaysia's states. This is due to the government support of public traveling which greatly reduces the expenses of people who utilize the bus services. Public transit networks provide a variety of beneficial tasks. Public transportation can offer a practical and effective substitute for private passenger automobiles when it is correctly constructed. Less dangerous chemical emissions and less traffic congestion in transportation networks are just two advantages of reducing the use of private autos.

Even with such a great convenience, travel by bus sometimes could be a pain in the ass. There are some reasons people rather travel through other sorts of vehicles compared to traveling by bus. This is mainly due to the delay of buses fetching people from bus stops. Delays are often due to heavy traffic jams, bus breakdowns, and traffic accidents which sometimes are unavoidable. Due to these incidents, passengers frequently spend most of their time waiting for the vehicle not knowing when it will come. This causes them to waste precious times when they could have used it to do other beneficial activities. Hence, to counter these situations, a Bus Breakdown and Delay system is created to help the society to apprehend this issue.

In this research, we will show the important concepts of the Bus Breakdown and Delay system by analyzing and collecting information from school bus vendors operating out in the field in real time by showing the implementations of the Bus Breakdown and Delay system. Another goal of this study is to develop a methodology for collecting, analyzing and assessing information on the transit time of bus routes GPS technology. In comparison to other systems, this one has the advantage of making it feasible to estimate travel times based on stops or routes. It is also more practical and affordable. This methodology enables you to immediately assess and analyze, daily modify the route schedule or create new routes, alter working hours, and lower societal costs while organizing high-quality transportation services for the populace.

1.1 Problem Statement

Bus staff often use the exact same routes when rotating a certain area to fetch passengers from one location to another. These may not be an ideal option for bus services as routes are frequently changed through time. Some of the bus drivers may end up using an older route system which causes heavier traffic jams which causes delays amongst the bus services. It is crucial to offer quick and reliable public transportation service at a time when the population's need for transportation is increasing. Public transit is a vehicle that competes directly with private autos. Currently, minimizing traffic congestion is one of the most critical challenges within Malaysia. This caused problems with the road network's capacity.

Bus businesses work in a stochastic environment. The quantity of elements and the intricacy of the connections between those factors determine stochastic conditions. These elements include changes in the number of passengers and city traffic congestion. Travelers must allow more time than usual for traffic delays, especially during rush hours. Traffic delays make it considerably more difficult to arrive late at your location and are much harsher than delays that were anticipated. Both passengers and the bus companies that operate them care about journey time and delays. Bus firms that experience unexpected delays risk having fewer people use public transportation to deliver goods on schedule. One of the markers of the quality of transportation services and a reflection of the effectiveness of the transportation system is travel time and delay.

One of the pressing issues is improving bus transit efficiency. It is prioritized to examine the relevant aspects and quantitatively assess them in order to create solutions to this problem. The significance of the research topic is determined by a wide range of difficulties connected to improving the effectiveness of urban public transportation. The major goal of this study was to recommend delay strategies that may take into account the effects of sporadic traffic jams and unplanned bus stops along the route.

1.2 Solution

Bus staff that encounter delays during the route are instructed to radio the dispatcher at the bus vendor's central office. The bus vendor staff are then instructed to log into the Bus Breakdown and Delay system to record the event and notify OPT. OPT customer service agents use this system to inform parents who call with questions regarding bus service. The Bus Breakdown and Delay system is publicly accessible and contains real time updates. All information in the system is entered by school bus vendor staff. A Bus Breakdown and Delay system is an application that allows transportation companies to track and monitor their bus fleets in real-time, and provides alerts and notifications when there are any delays or breakdowns.

One way to develop such a system using supervised learning is to use historical data to train a machine learning model that can predict the likelihood of a breakdown or delay occurring. The model would be trained using data such as the time of day, weather conditions, route details, and other factors that could impact the performance of the bus. Once the model has been trained, it can be used to analyze real-time data and provide alerts when there is a high probability of a breakdown or delay. This can be done by comparing the current data to the model's predictions, and generating alerts if the data falls outside of the expected range.

To implement such a system, some of the following steps could be taken first are data collection. Collect historical data about bus routes, schedules, and breakdown/delay incidents. This data should be diverse and cover a variety of conditions. Next, we must clean the data and prepare it for training by removing irrelevant information, handling missing values, and converting categorical variables into numerical values. After data cleaning, we must choose the most relevant features for the model, based on the historical data and domain knowledge. Select an appropriate machine learning algorithm (such as logistic regression, decision tree, or random forest) and train the model on the historical data. Lastly, evaluate the performance of the model by comparing its predictions to the actual breakdown/delay incidents. Use metrics such as accuracy, precision, recall, and F1-score to evaluate the performance. Once the model is trained and evaluated, deploy it in a production environment to generate real-time predictions and alerts. Monitor the system's performance and continuously update the model as new data becomes available.

By using a supervised learning approach to develop a Bus Breakdown and Delay system, transportation companies can proactively identify potential issues before they become major problems, and improve the reliability and efficiency of their bus fleets.

1.3 Related Work

There are many research papers that focus on the prediction of bus arrival and delays in transportation systems, including buses. Different machine learning algorithms have been applied in order to predict delays. Before introducing the algorithms used, most researchers use different analysis methods and data mining skills on the data in order to understand the knowledge or insight hidden behind and to improve the predictive efficiency of the models used.

Gradient boosted tree, Multi-layer Perceptron, neural network algorithms, XGboost, ADA boosting

1.3.1 Gradient Boosting Regression Tree (GBRT)

Another research showing a bus arrival time prediction approach based on GPS position and real-time traffic flow was proposed by the authors of the article "A Bus Arrival Time Prediction Method Based on GPS Position and Real-Time Traffic Flow" by Lei Jianmei, et al. (2017). The Gradient Boosting Regression Tree (GBRT) is a machine learning algorithm that was used in the proposed method to predict the bus arrival time.

The GBRT algorithm is an example of an ensemble learning technique that predicts the target variable using decision trees. The approach sequentially corrects the residual errors of the preceding decision tree by fitting a succession of decision trees to the data. The weighted sum of each tree's predicted values is the GBRT algorithm's output.

Overall, the study showed how machine learning algorithms combined with GPS position data and real-time traffic flow data can accurately estimate bus arrival timings.

1.3.2 Multi-layer Perceptron

Multi-layer Perceptron or (MLP) Neural Networks are strong machine learning models that are capable of detecting intricate patterns in data. They handle nonlinear interactions between input and output variables exceptionally well. MLPs are able to generalize effectively to new data and can learn to generate predictions based on a vast number of input features. They can also be trained rapidly and are often simple to apply.

There is a study done by Xinxin Li, Xiangyu Kong, and Yang Wang regarding "Forecasting Bus Arrival Times Using Multi-Layer Perceptron Neural Networks." In this study, a multi-layer perceptron neural network is used to forecast bus arrival times using data from the New York City Transit Authority. The outcomes demonstrate that the suggested approach performs better than established approaches like linear regression and ARIMA.

1.3.3 Neural network algorithms

There is research done by Ying Hua, Yang Wang, and Xinxin Li regarding "Real-Time Bus Arrival Time Prediction Using Deep Neural Networks." In this study, a technique for anticipating bus arrival times is proposed utilizing deep neural networks (DNNs), a class of neural networks with numerous hidden layers. The results demonstrate that the DNN method outperforms more established techniques like linear regression and support vector regression when applied to data from the New York City Transportation Authority. The authors also examine how numerous elements, such as the climate and traffic, affect bus arrival times. Overall, this experiment shows how neural network algorithms may be used to accurately anticipate bus arrival times, which can boost the dependability and effectiveness of public transportation systems.

1.3.4 XGboost

Huy Tu, Daqing Zhang, and Yinhai Wang's "Bus Arrival Time Prediction using XGBoost." The XGBoost algorithm is used in this study to offer a technique for forecasting bus arrival times, which is then applied to data from buses in New York City. The outcomes demonstrate that the suggested approach outperforms more established techniques like random forest and support vector regression.

The effective gradient boosting algorithm known as XGBoost is frequently used for machine learning applications like regression and classification. It differs from previous algorithms in a number of ways, including its capacity to manage missing data and the handling of both numerical and categorical features. The complicated data in the New York bus delay prediction problem is ideal for XGBoost since it can handle big datasets with numerous features.

1.3.5 ADA boosting

The study "Real-time Bus Arrival Time Prediction Using ADABoost Algorithm" by Shaopeng Zhang, Yue Liu, and Wei Liu was released. This article proposes an ADABoost-based method for real-time bus arrival time prediction using data from buses in New York City. The results show that the proposed method outperforms more established methods like linear regression and supports vector regression.

ADABoost is another efficient boosting technique for classification and regression workloads. It works by combining a number of ineffective learners to create a potent prediction. The ability to handle noisy data and the ability to prevent overfitting are just two advantages of ADABoost. Given that it can handle large datasets with a variety of features, it is appropriate for the New York bus delay prediction problem.

1.4 Algorithm

Gradient boosted tree, Multi-layer Perceptron, neural network algorithms, XGboost, ADA boosting

1.4.1 Gradient Boosted Trees (GBT)

Gradient Boosted Trees (GBT) is an ensemble machine learning algorithm that combines the results of multiple decision trees to improve the accuracy of predictions. At each iteration, a new decision tree is fitted to the residuals of the previous iteration, creating a stronger model that can capture complex patterns in the data. GBT is known for its ability to handle high-dimensional and noisy data, and it is highly interpretable, allowing the user to understand how each decision tree contributes to the final prediction.

1.4.2 The Multi-layer Perceptron (MLP)

The Multi-layer Perceptron (MLP) algorithm is a type of artificial neural network that is commonly used for supervised learning tasks, such as regression and classification. It consists of multiple layers of interconnected nodes, each one performing a nonlinear transformation of the input data. The nodes in the hidden layers use activation functions to introduce nonlinearity into the model, allowing it to capture complex relationships between the input and output variables. The output layer produces the final prediction based on the transformed input data. The MLP algorithm is trained using an iterative optimization process called backpropagation, which adjusts the weights and biases of the nodes to minimize the error between the predicted and actual outputs. It is known for its ability to handle complex and high-dimensional data, as well as its flexibility in modeling nonlinear relationships. However, it can be sensitive to the choice of hyperparameters and prone to overfitting if not properly regularized.

1.4.3 Neural Network Algorithms

There are several neural network algorithms that are commonly used in machine learning, each with their own strengths and weaknesses. Here are some of the most popular ones:

Multilayer Perceptron (MLP) - a feedforward neural network with one or more hidden layers.

Convolutional Neural Network (CNN) - a specialized neural network for image processing, where the neurons are arranged in a way that takes into account the spatial relationships between pixels.

Recurrent Neural Network (RNN) - a type of neural network that can process sequences of input data, by having neurons that have a memory of previous inputs.

Long Short-Term Memory (LSTM) - a specialized type of RNN that is designed to handle long-term dependencies in sequential data.

Autoencoder - a neural network that is trained to reconstruct its input, often used for dimensionality reduction and feature extraction.

Generative Adversarial Network (GAN) - a neural network that consists of two parts, a generator and a discriminator, that are trained together to generate realistic data.

Boltzmann machine - a type of neural network that is used for unsupervised learning, by learning a probabilistic model of the input data.

Each of these algorithms has its own advantages and disadvantages, and is suitable for different types of problems. Choosing the right neural network algorithm for a particular task requires careful consideration of the nature of the data and the goals of the project.

Long Short-Term Memory (LSTM) is a type of recurrent neural network (RNN) that is widely used for sequential data processing tasks, such as natural language processing and time series analysis. Unlike traditional RNNs, which suffer from the vanishing gradient problem when processing long sequences, LSTM is able to maintain long-term dependencies and avoid gradient degradation. This is achieved by introducing memory cells that can selectively store or erase information over time, as well as gates that regulate the flow of information into and out of the cells.

1.4.4 XGBoost (eXtreme Gradient Boosting)

XGBoost (eXtreme Gradient Boosting) is a popular machine learning algorithm that is widely used for both regression and classification tasks. It is based on the concept of gradient boosting, which involves iteratively training an ensemble of weak prediction models, such as decision trees, to create a strong model that can make accurate predictions. XGBoost extends this concept by using a more regularized model and a novel technique to optimize the objective function, resulting in better performance and faster computation. The algorithm works by building decision trees in a greedy and parallelized way, while applying a regularization term to prevent overfitting. It also uses a gradient-based optimization technique to find the best split points and weighting for each feature.

1.4.5 AdaBoost (Adaptive Boosting)

AdaBoost (Adaptive Boosting) is a machine learning algorithm used for classification and regression tasks. It is an ensemble learning algorithm that combines the results of multiple weak learners to improve the accuracy of predictions. The algorithm works by iteratively training a sequence of weak classifiers on a dataset, with each classifier attempting to correct the errors made by the previous ones. The final prediction is made by taking a weighted average of the individual classifier predictions. The strength of AdaBoost lies in its ability to focus on the difficult examples in a dataset, and the adaptiveness of the algorithm in iteratively re-weighting the examples to reduce bias. This makes it particularly useful for classification tasks where the data is imbalanced or difficult to classify. AdaBoost can also handle noisy data and is less prone to overfitting compared to other algorithms. One of the most popular implementations of AdaBoost is the scikit-learn library in Python, which offers a range of configurable parameters to optimize the performance of the algorithm.

2.1 Data Source

The data used in this project comes from Kaggle (NY Bus Breakdown and Delays)

Feature Name	Description
School_Year	The academic year the incident occurred in
Busbreakdown_ID	Unique identifier for the bus breakdown incident
Run_Type	Type of bus run (e.g., Special Ed, General Ed, Pre-K, etc.)
Bus_No	Identification number of the bus
Route_Number	The number assigned to the bus route
Reason	Reason for the delay or breakdown (e.g., mechanical problem, flat tire, etc.)
Schools_Serviced	Names of the schools serviced by the bus
Occurred_On	Date and time the incident occurred
Created_On	Date and time the incident was reported
Boro	The borough in which the incident occurred
Bus_Company_Name	Name of the bus company that operates the bus
How_Long_Delayed	The duration of the delay in minutes
Number_Of_Students_On_The_Bus	The number of students on the bus at the time of the incident

Has_Contractor_Notified_Schools	Whether the bus contractor notified the schools of the delay	
Has_Contractor_Notified_Parents	Whether the bus contractor notified the parents of the students on the bus	
Have_You_Alerted_OPT	Whether the Office of Pupil Transportation (OPT) was alerted of the incident	
Informed_On	Date and time the OPT was informed	
Incident_Number	The incident number assigned to the incident	
Last_Updated_On	Date and time the incident was last updated	
Breakdown_or_Running_Late	Indicates whether the bus is delayed due to a breakdown or running late	
School_Age_or_PreK	Indicates whether the schools serviced by the bus are for school-age children or pre-K children	

2.2 Data Understanding

Install External Package

```
!pip install xgboost
!pip install flask
```

Library

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.neural_network import MLPRegressor
from keras.models import Sequential
from keras.layers import Dense
import xgboost as xgb
from sklearn.ensemble import AdaBoostRegressor
from sklearn.metrics import mean_squared_error, r2_score
import pickle
from flask import Flask, jsonify, request
```

Read csv data

```
# read csv data
df = pd.read_csv('bus-breakdown-and-delays.csv', low_memory=False)
df.head()
```

	School_Year	Busbreakdown_ID	Run_Type	Bus_No	Route_Number	Reason	Schools_Serviced
0	2015-2016	1227538	Special Ed AM Run	2621	J711	Heavy Traffic	75003
1	2015-2016	1227539	Special Ed AM Run	1260	M351	Heavy Traffic	06716
2	2015-2016	1227540	Pre-K/EI	418	3	Heavy Traffic	C445
3	2015-2016	1227541	Special Ed AM Run	4522	M271	Heavy Traffic	02699
4	2015-2016	1227542	Special Ed AM Run	3124	M373	Heavy Traffic	02116
5 rc	ows × 21 colu	mns					
4							>

Describe data

In [4]: df.describe()

Out[4]:

	Busbreakdown_ID	Number_Of_Students_On_The_Bus
count	3.794120e+05	379412.000000
mean	1.410723e+06	3.494162
std	1.154315e+05	74.213592
min	1.212691e+06	0.000000
25%	1.311426e+06	0.000000
50%	1.407630e+06	0.000000
75%	1.512374e+06	3.000000
max	1.605584e+06	9658.000000

Check the data size

```
In [29]: df.shape
Out[29]: (379412, 21)
```

Count the number of occurrences of each unique value in the 'Number_Of_Students_On_The_Bus'

```
In [6]: df['Number_Of_Students_On_The_Bus'].value_counts()
Out[6]: 0
                233779
        2
                 21658
        1
                 21388
        3
                 20293
        4
                 16810
                     1
        1420
                     1
        1449
        2557
                     1
                     1
        4570
        3258
        Name: Number_Of_Students_On_The_Bus, Length: 255, dtype: int64
```

Check data types of all columns

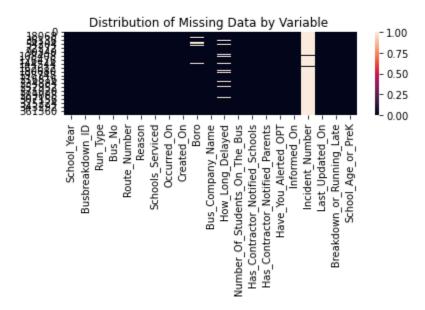
In [7]: df.dtypes

Out[7]:	School_Year	object	
	Busbreakdown_ID	int64	
	Run_Type	object	
	Bus_No	object	
	Route_Number	object	
	Reason	object	
	Schools_Serviced	object	
	Occurred_On	object	
	Created_On	object	
	Boro	object	
	Bus_Company_Name	object	
	How_Long_Delayed	object	
	Number_Of_Students_On_The_Bus	int64	
	Has_Contractor_Notified_Schools	object	
	Has_Contractor_Notified_Parents	object	
	Have_You_Alerted_OPT	object	
	Informed_On	object	
	Incident_Number	object	
	Last_Updated_On	object	
	Breakdown_or_Running_Late	object	
	School_Age_or_PreK	object	
	dtype: object		

Heatmap

```
In [8]: sns.heatmap(df.isnull()) #See distribution of missing data
plt.figsize = (5,2.5)
plt.tight_layout()
plt.title('Distribution of Missing Data by Variable ')
```

Out[8]: Text(0.5, 1.0, 'Distribution of Missing Data by Variable ')



Display features, missing values and unique values

```
# Display details of dataset
print ("Rows :" ,df.shape[0])
print ("Columns :" ,df.shape[1])
print ("\nFeatures :\n" ,df.columns.tolist())
print ("\nMissing values : \n", df.isnull().sum())
print ("\nUnique values : \n",df.nunique())
```

Rows : 379412 Columns : 20

Features :

['School_Year', 'Busbreakdown_ID', 'Run_Type', 'Bus_No', 'Route_Number', 'Reas on', 'Schools_Serviced', 'Occurred_On', 'Created_On', 'Boro', 'Bus_Company_Nam e', 'How_Long_Delayed', 'Number_Of_Students_On_The_Bus', 'Has_Contractor_Notified_Schools', 'Has_Contractor_Notified_Parents', 'Have_You_Alerted_OPT', 'Inform ed_On', 'Last_Updated_On', 'Breakdown_or_Running_Late', 'School_Age_or_PreK']

Missing values :	
School_Year	0
Busbreakdown_ID	0
Run_Type	3
Bus_No	10
Route_Number	7
Reason	2
Schools_Serviced	7
Occurred_On	0
Created_On	0
Boro	11095
Bus_Company_Name	0
How_Long_Delayed	43138
Number_Of_Students_On_The_Bus	0
Has_Contractor_Notified_Schools	0
Has_Contractor_Notified_Parents	0
Have_You_Alerted_OPT	0
Informed_On	0
Last_Updated_On	0
Breakdown_or_Running_Late	0
School_Age_or_PreK	0
dtype: int64	

Unique values :	
School_Year	5
Busbreakdown_ID	325364
Run_Type	10
Bus_No	14538
Route_Number	14516
Reason	10
Schools_Serviced	19730
Occurred_On	147507
Created_On	157444
Boro	11
Bus_Company_Name	122
How_Long_Delayed	1805
Number_Of_Students_On_The_Bus	255
Has_Contractor_Notified_Schools	2
Has_Contractor_Notified_Parents	2
Have_You_Alerted_OPT	2
Informed_On	157444
Incident_Number	6501
Last_Updated_On	155614
Breakdown_or_Running_Late	2
School_Age_or_PreK	2
dtype: int64	

2.3 Data cleaning

Drop the most missing value column

```
In [9]: #Drop incident number, most of column is missing
        df = df.drop(['Incident_Number'], axis = 1)
       df.info()
        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 379412 entries, 0 to 379411
       Data columns (total 20 columns):
        # Column
                                            Non-Null Count
                                                            Dtype
        0 School Year
                                            379412 non-null object
        1 Busbreakdown ID
                                            379412 non-null int64
        2 Run_Type
                                            379409 non-null object
        3 Bus_No
                                            379402 non-null object
            Route Number
                                            379405 non-null object
                                           379410 non-null object
        5 Reason
        6 Schools Serviced
                                            379405 non-null object
        7 Occurred On
                                            379412 non-null object
        8 Created_On
                                            379412 non-null object
        9
            Boro
                                            368317 non-null object
        10 Bus Company Name
                                            379412 non-null object
        11 How Long Delayed
                                           336274 non-null object
        12 Number_Of_Students_On_The_Bus 379412 non-null int64
        13 Has_Contractor_Notified_Schools 379412 non-null object
        14 Has Contractor Notified Parents 379412 non-null object
        15 Have You Alerted OPT
                                            379412 non-null object
        16 Informed On
                                            379412 non-null object
        17 Last Updated On
                                            379412 non-null object
        18 Breakdown_or_Running_Late
                                            379412 non-null object
        19 School Age or PreK
                                            379412 non-null object
        dtypes: int64(2), object(18)
        memory usage: 57.9+ MB
```

Display columns

Extract digits from string column

```
In [11]: #Extract digits from string column
           df.loc[:, 'Delay'] = df['How_Long_Delayed'].str.extract('(\d+)').copy()
#Check if regex worked- Yes!
           df.head()
Out[11]:
               School_Year Busbreakdown_ID Run_Type Bus_No Route_Number Reason Schools_Serviced
                                                  Special
                                                                                    Heavy
                                                                            J711
            0
                 2015-2016
                                     1227538
                                                  Ėd AM
                                                            2621
                                                                                                      75003
                                                                                    Traffic
                                                    Run
                                                  Special
                                                                                    Heavy
                                     1227539
                                                                            M351
                                                                                                      06716
                 2015-2016
                                                             1260
                                                  Ed AM
                                                                                    Traffic
                                                    Run
                                                                                    Heavy
                 2015-2016
                                     1227540
                                                 Pre-K/EI
                                                              418
                                                                               3
                                                                                                       C445
            2
                                                                                    Traffic
                                                  Special
                                                                                    Heavy
                 2015-2016
                                     1227541
                                                             4522
                                                                           M271
                                                                                                      02699
            3
                                                  Ed AM
                                                                                    Traffic
                                                    Run
                                                  Special
                                                                                    Heavy
                 2015-2016
                                     1227542
                                                  Ed AM
                                                            3124
                                                                           M373
                                                                                                      02116
                                                                                    Traffic
           5 rows x 21 columns
```

Check whether data is null

```
df[df['Delay'].isnull()]#Check if data is null
#We see that there's question marks or other irregularaties- lets drop this data
```

fied Schools	Has_Contractor_Notified_Parents	Have You Alerted OPT	Informed On	Last Updated On	Breakdown or Running Late	School Age or PreK	Delav
Yes	No	Yes	2015-11-05 8:12:00	2015-11-05 8:12:14	Running Late	School-Age	NaN
No	No	No	2015-11-05 8:14:00	2015-11-05 8:14:08	Running Late	School-Age	NaN
Yes	Yes	Yes	2015-11-05 8:25:00	2015-11-05 8:25:11	Running Late	School-Age	NaN
No	No	No	2015-11-05 8:26:00	2015-11-05 8:26:54	Running Late	School-Age	NaN
Yes	Yes	No	2015-11-04 7:11:00	2015-11-04 7:11:53	Breakdown	School-Age	NaN

Yes	Yes	No	2016-04-21 8:15:00	2016-04-21 8:15:29	Breakdown	School-Age	NaN
Yes	Yes	Yes	2015-09-16 15:03:00	2015-09-16 15:03:41	Breakdown	Pre-K	NaN
Yes	Yes	Yes	2015-09-22 7:16:00	2015-09-22 7:16:40	Running Late	School-Age	NaN
Yes	Yes	No	2018-04-17 8:08:00	1900-01-01 0:00:00	Breakdown	School-Age	NaN
Yes	Yes	No	2018-04-17 8:08:00	1900-01-01 0:00:00	Breakdown	School-Age	NaN

Drop na values and display

```
In [12]: df_clean = df.dropna() #Drop remainaing NAs
         df clean.info()
        df_clean.head()
         <class 'pandas.core.frame.DataFrame'>
         Int64Index: 325760 entries, 1 to 379411
        Data columns (total 21 columns):
         # Column
                                             Non-Null Count
                                                             Dtype
         --- -----
                                             -----
         0 School_Year
                                            325760 non-null object
             Busbreakdown ID
                                             325760 non-null int64
             Run_Type
                                            325760 non-null object
             Bus_No
                                           325760 non-null object
                                            325760 non-null object
         4
             Route_Number
             Reason
                                            325760 non-null object
             Schools_Serviced
                                           325760 non-null object
         6
             Occurred On
                                           325760 non-null object
         8
             Created_On
                                            325760 non-null object
                                            325760 non-null object
             Boro
         10 Bus_Company_Name
                                            325760 non-null object
         11 How_Long_Delayed
                                            325760 non-null object
         12 Number_Of_Students_On_The_Bus 325760 non-null int64
         13 Has_Contractor_Notified_Schools 325760 non-null object
         14 Has_Contractor_Notified_Parents 325760 non-null object
         15 Have_You_Alerted_OPT
                                            325760 non-null object
         16 Informed_On
                                            325760 non-null object
                                            325760 non-null object
         17 Last Updated On
         18 Breakdown_or_Running_Late
                                             325760 non-null object
         19 School_Age_or_PreK
                                            325760 non-null object
         20 Delay
                                            325760 non-null object
         dtypes: int64(2), object(19)
         memory usage: 54.7+ MB
```

Out[12]:

	School_Year	Busbreakdown_ID	Run_Type	Bus_No	Route_Number	Reason	Schools_Serviced
1	2015-2016	1227539	Special Ed AM Run	1260	M351	Heavy Traffic	06716
2	2015-2016	1227540	Pre-K/EI	418	3	Heavy Traffic	C445
3	2015-2016	1227541	Special Ed AM Run	4522	M271	Heavy Traffic	02699
5	2015-2016	1227543	Special Ed AM Run	HT1502	W796	Heavy Traffic	75407
6	2015-2016	1227544	Special Ed AM Run	142	W633	Heavy Traffic	75670

5 rows × 21 columns

Check any Na values left

In [15]:	<pre>df_clean = df_clean.dropna() #Dr df_clean.isnull().sum() #Check t</pre>	•
Out[15]:	School_Year	0
	Busbreakdown ID	0
	Run Type	0
	Bus No	0
	Route Number	0
	Reason	0
	Schools Serviced	0
	Occurred_On	0
	Created_On	0
	Boro	0
	Bus_Company_Name	0
	How_Long_Delayed	0
	Number_Of_Students_On_The_Bus	0
	Has_Contractor_Notified_Schools	0
	Has_Contractor_Notified_Parents	0
	Have_You_Alerted_OPT	0
	Informed_On	0
	Last_Updated_On	0
	Breakdown_or_Running_Late	0
	School_Age_or_PreK	0
	Delay	0
	dtype: int64	

Convert the string to integer

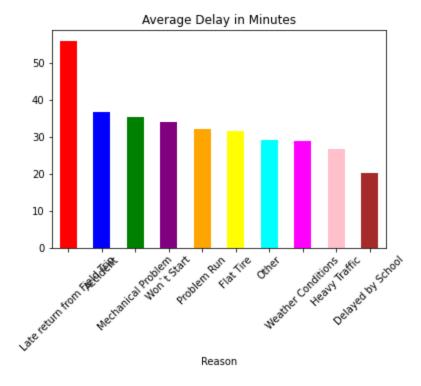
```
In [16]: #Convert string to integer
df_clean.loc[:, 'Delay'] = pd.to_numeric(df_clean['Delay'])
```

Drop original column

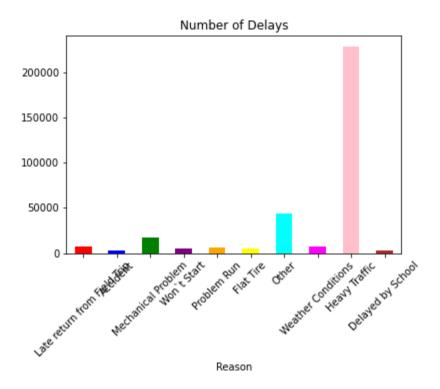
```
In [17]: #Drop original column
df_clean = df_clean.drop(['How_Long_Delayed'], axis = 1)
```

Plot chart for average delay and number of delays

```
In [18]: reasons = pd.pivot_table(df_clean, index = 'Reason', values = 'Delay', aggfunc =
          # Define a list of 10 different colors
          colors = ['red', 'blue', 'green', 'purple', 'orange', 'yellow', 'cyan', 'magenta
          # Create a pivot table of reasons and delay
          reasons = pd.pivot_table(df_clean, index='Reason', values='Delay', aggfunc=[np.m
          # Plot the average delay chart
          plt.figure(figsize=(10, 6))
         reasons.plot(kind='bar', y=('mean', 'Delay'), color=colors)
plt.title('Average Delay in Minutes')
          plt.xticks(rotation=45)
          plt.legend().remove()
          plt.show()
          # Plot the number of delays chart
          plt.figure(figsize=(10, 6))
          reasons.plot(kind='bar', y=('size', 'Delay'), color=colors)
          plt.title('Number of Delays')
          plt.xticks(rotation=45)
          plt.legend().remove()
          plt.show()
          <Figure size 720x432 with 0 Axes>
```



<Figure size 720x432 with 0 Axes>



Plot boxplot for 'Delay' to find outliers

Plot boxplot to check if we need to remove further outliers

```
In [22]: df_clean['Route_Number'].value_counts()
Out[22]: 1
                  4225
         2
                  3105
                  2623
         3
                  2617
         4
                  1519
                  . . .
         K662
         K9172
                     1
         K9485
                     1
         K9131
                     1
         M9138
                     1
         Name: Route_Number, Length: 13862, dtype: int64
```

List first 6 mean Delay and Size Delay in pivot table

```
In [24]: #Filter to see cases where route is top 6 in # of delays
          routes = ['1','2','3','5','4','6']
         top_routes = df_clean[df_clean['Route_Number'].isin(routes)]
In [25]: routes_pivot = pd.pivot_table(top_routes,
                                         index = 'Route_Number',
                                         values = 'Delay',
                                         aggfunc = [np.mean,np.size])
          routes pivot.head(6)
Out[25]:
                                 size
                        mean
                        Delay
                                 Delay
          Route_Number
                     1 21.269349
                                  4225
                     2 22.322383
                                 3105
                     3 21.853267
                                  2617
                     4 20.259381
                                 1519
                     5 20.535646
                                 2623
                     6 22.498741
                                 1191
```

Count the number of occurrences of each unique value in the 'Bus_Company_Name' '

```
In [27]: df clean['Bus Company Name'].value counts()
Out[27]: LEESEL TRANSPORTATION CORP (B2192)
                                                 42339
         G.V.C., LTD.
                                                 21442
         PIONEER TRANSPORTATION CORP
                                                 17723
         RELIANT TRANSPORTATION, INC (B2321)
                                                 15333
         BORO TRANSIT, INC.
                                                 14336
         R & C TRANSIT, INC. (B232
                                                      2
         1967
                                                     1
         phillip bus service
                                                     1
         FORTUNA BUS COMPANY
                                                      1
         Name: Bus_Company_Name, Length: 118, dtype: int64
```

To see increasing trend year on year in quantity and any significant deviations in delay by year

```
In [28]: #First Let's remove unnecessary features, checking 1 by 1

#School Year- is it relevant?

df_clean['School_Year'].value_counts().plot(kind = 'bar')

plt.xticks(rotation = 75) #Make Data cleaner to read

#See an increasing trend year on year in quantity-let's investigate if there's a

d

Out[28]: (array([0, 1, 2, 3, 4]),
        [Text(0, 0, '2018-2019'),
        Text(1, 0, '2017-2018'),
        Text(2, 0, '2016-2017'),
        Text(3, 0, '2015-2016'),
        Text(4, 0, '2019-2020')])

80000-

40000-

40000-

40000-

40000-

40000-

80000-

40000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

80000-

800
```

About 29 minutes in average delay time

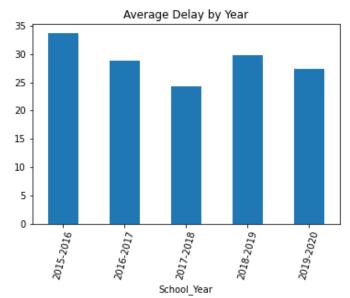
```
In [29]: #Let's first see average delay, across the dataset
    df_clean['Delay'].mean() #Around 29 mins is the average delay time
Out[29]: 28.52876043713163
```

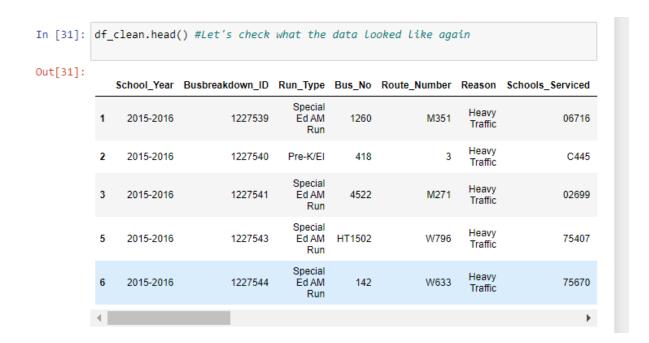
To investigate any year is terribly far off from another

```
In [30]: pd.pivot_table(df_clean, index = 'School_Year', values = 'Delay', aggfunc = np.
    plt.legend().remove() #Get rid of legend
    plt.title('Average Delay by Year')
    plt.xticks(rotation = 75) #Make easier to read

#Doesn't Look any year is terribly far off from another but also not congruent-

Out[30]: (array([0, 1, 2, 3, 4]),
        [Text(0, 0, '2015-2016'),
        Text(1, 0, '2016-2017'),
        Text(2, 0, '2017-2018'),
        Text(3, 0, '2018-2019'),
        Text(4, 0, '2019-2020')])
```





Count the number of occurrences of each unique value in the 'Busbreakdown ID'

```
In [32]:
         #Data seems like no noise, we'll drop
         df_clean['Busbreakdown_ID'].value_counts()
Out[32]: 1602184
         1603512
                     2
                     2
         1580876
                     2
         1580874
         1580873
                     2
                    ٠.
         1371134
                     1
         1371135
                     1
         1371140
                     1
         1371141
                     1
         1565321
                     1
         Name: Busbreakdown_ID, Length: 278328, dtype: int64
```

In [33]: df_clean = df_clean.drop(['Busbreakdown_ID'], axis = 1) df_clean.head() Out[33]: School_Year Run_Type Bus_No Route_Number Reason Schools_Serviced Occurred_On Special 2015-11-Heavy 2015-2016 Ėd AM 1260 M351 1 06716 05T08:10:00.000 05 Traffic Run Heavy 2015-11-2 2015-2016 Pre-K/EI 418 3 05T08:09:00.000 05 Traffic Special 2015-11-Heavy 3 2015-2016 Ėd AM 4522 M271 02699 Traffic 05T08:12:00.000 05 Run Special Heavy 2015-11-5 2015-2016 Ėd AM HT1502 W796 75407 05T07:58:00.000 05 Traffic Run Special Heavy 2015-11-6 2015-2016 Ėd AM W633 142 75670 05T08:24:00.000 05 Traffic Run

```
In [34]: bus_num = pd.pivot_table(df_clean, index = 'Bus_No', values = 'Delay',aggfunc =
          bus_num
          #Create pivot to see number of delays by bus number
          #We see that a lot have only have 1.
          #Instead of one hot encoding, let's just convert to digits
Out[34]:
                  Delay
          Bus_No
             1801
                    835
             1318
                    770
             1107
                    752
             9302
                    732
             1470
                    716
            58127
           58166D
            58176
            58180
                      1
```

`870

```
In [35]: #Extract digits from string column
         df_clean['Bus_Number'] = df_clean['Bus_No'].str.extract('(\d+)')
         #Convert string to integer
         df_clean['Bus_Number'] = pd.to_numeric(df_clean['Bus_Number'])
         df_clean.isnull().sum()
         #We now have some more NAs- let's do a quick investigation
Out[35]: School_Year
                                              0
         Run_Type
                                              0
         Bus No
                                              0
         Route Number
                                              0
         Reason
         Schools_Serviced
                                              0
         Occurred_On
                                              0
         Created On
                                              0
         Boro
                                              0
         Bus Company Name
         Number_Of_Students_On_The_Bus
         Has_Contractor_Notified_Schools
         Has_Contractor_Notified_Parents
                                              0
         Have You Alerted OPT
                                              0
         Informed_On
         Last Updated On
         Breakdown_or_Running_Late
                                              0
         School_Age_or_PreK
                                              0
         Delay
                                              0
         Bus Number
                                            175
         dtype: int64
```

```
In [36]: #Looks like noisy data, will drop
    df_clean[df_clean['Bus_Number'].isnull()]
    df_clean = df_clean.dropna()
    #Drop original column
    df_clean = df_clean.drop(['Bus_No'], axis = 1)
```

In [37]: df_clean.head()

Out[37]:

	School_Year	Run_Type	Route_Number	Reason	Schools_Serviced	Occurred_On	Created
1	2015-2016	Special Ed AM Run	M351	Heavy Traffic	06716	2015-11- 05T08:10:00.000	2015 05T08:12:00
2	2015-2016	Pre-K/EI	3	Heavy Traffic	C445	2015-11- 05T08:09:00.000	2015 05T08:13:00
3	2015-2016	Special Ed AM Run	M271	Heavy Traffic	02699	2015-11- 05T08:12:00.000	2015 05T08:14:00
5	2015-2016	Special Ed AM Run	W796	Heavy Traffic	75407	2015-11- 05T07:58:00.000	2015 05T08:14:00
6	2015-2016	Special Ed AM Run	W633	Heavy Traffic	75670	2015-11- 05T08:24:00.000	2015 05T08:15:00
41							

Find Correlation between features and target

```
In [38]: #Let's look at the current correlation across features
df_clean.corr()
```

Out[38]:

	${\bf Number_Of_Students_On_The_Bus}$	Delay	Bus_Number
Number_Of_Students_On_The_Bus	1.000000	-3.277360e-04	6.393997e-04
Delay	-0.000328	1.000000e+00	2.770129e-08
Bus_Number	0.000639	2.770129e-08	1.000000e+00

Plot the trip distribution

```
In [39]: df clean['Run Type'].value counts().plot(kind = 'bar')
           plt.title('Trip distribution ')
           plt.xticks(rotation = 75) #Data heavily weighted towards Special Ed AM in terms
Out[39]: (array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9]),
            [Text(0, 0, 'Special Ed AM Run'),
             Text(1, 0, 'Special Ed PM Run'),
             Text(2, 0, 'Pre-K/EI'),
             Text(3, 0, 'General Ed AM Run'),
             Text(4, 0, 'General Ed PM Run'),
             Text(5, 0, 'Special Ed Field Trip'),
Text(6, 0, 'General Ed Field Trip'),
             Text(7, 0, 'Project Read PM Run'),
             Text(8, 0, 'Project Read AM Run'),
             Text(9, 0, 'Project Read Field Trip')])
                                     Trip distribution
            175000
            150000
            125000
            100000
             75000
             50000
             25000
                                 General Ed AM Run
                                                General Ed Field Trip
                                          Special Ed Field Trip
```

Split the data into y=target X=features

3.0 Algorithm

3.1 Feature 1 Algorithm

3.1.1 Gradient Boosted Tree (GBT)

```
from sklearn.model selection import cross val score
from sklearn.metrics import mean_absolute_error, mean_squared_log_error
from sklearn.model selection import GridSearchCV
# Define the parameter grid to search over
param_grid = {
      'n_estimators': [100, 200],
     'max_depth': [3, 4],
'learning_rate': [0.01, 0.1],
'loss': ['ls', 'lad']
# Create a GradientBoostingRegressor object
gbr = GradientBoostingRegressor()
# Create a GridSearchCV object
grid_search = GridSearchCV(estimator=gbr, param_grid=param_grid, cv=5, n_jobs=-1)
# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)
# Print the best hyperparameters
print(grid_search.best_params_)
# Use the best estimator to make predictions on the testing set
y_pred = grid_search.best_estimator_.predict(X_test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
gbr_mae1 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y_test_abs = np.abs(y_test)
y_pred_abs = np.abs(y_pred)
gbr_rmsle1 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
# Calculate the cross-validated MAE and RMSLE scores
gbr_cv_mae1 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_absolute_error'))
gbr_cv_rmsle1 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error
print("Gradient Boosted Tree MAE:", gbr_mae1)
print("Gradient Boosted Tree RMSLE:", gbr_rmsle1)
print("Cross-validated MAE:", -gbr_cv_mae1)
print("Cross-validated RMSLE:", np.sqrt(-gbr_cv_rmsle1))
{'learning_rate': 0.1, 'loss': 'lad', 'max_depth': 4, 'n_estimators': 100}
Gradient Boosted Tree MAE: 11.449102788882989
Gradient Boosted Tree RMSLE: 0.6678114345568975
Cross-validated MAE: 24.52570717302958
Cross-validated RMSLE: 0.6765355769822062
```

The MAE output for GBT lies at 11.449102788882989, where the RSMLE is 0.6678114345568975. The Cross-validated value for MAE in GBT is 24.52570717302958, where the cross-validated value for RMSLE is 0.6765355769822062.

3.1.2 Multi-layer Perceptron (MLP)

from sklearn.neural network import MLPRegressor

Multi-layer Perceptron MAE: 12.506735737946727 Multi-layer Perceptron RMSLE: 0.6870849604076169

Cross-validated MAE: 25.510107232399836

Cross-validated RMSLE: nan

```
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_absolute_error, mean_squared_log_error
# Define the parameter grid to search over
param_grid = {
     'hidden_layer_sizes': [(10,), (50,), (100,)],
     'activation': ['relu', 'tanh'],
    'solver': ['adam', 'sgd'],
'learning_rate': ['constant', 'adaptive']
# Create a MLPRegressor object
mlp = MLPRegressor()
# Create a GridSearchCV object
grid search = GridSearchCV(estimator=mlp, param grid=param grid, cv=5, n jobs=-1)
# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)
# Print the best hyperparameters
print(grid search.best params )
# Use the best estimator to make predictions on the testing set
y_pred = grid_search.best_estimator_.predict(X_test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
mlp_mae1 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y_test_abs = np.abs(y_test)
y pred abs = np.abs(y pred)
mlp_rmsle1 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
# Calculate the cross-validated MAE and RMSLE scores
mlp_cv_mae1 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_absolute error'))
mlp_cv_rmsle1 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error
print("Multi-layer Perceptron MAE:", mlp_mae1)
print("Multi-layer Perceptron RMSLE:", mlp rmsle1)
print("Cross-validated MAE:", -mlp_cv_mae1)
print("Cross-validated RMSLE:", np.sqrt(-mlp_cv_rmsle1))
```

The MAE output for MLP lies at 12.506735737946727, where the RSMLE is 0.6870849604076169. The Cross-validated value for MAE in MLP is 25.510107232399836, where the cross-validated value for RMSLE is NULL.

3.1.3 Neural Network

```
from keras.wrappers.scikit_learn import KerasRegressor
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_absolute_error, mean_squared_log_error
# Define the neural network model
def create_nn_model():
    model = Sequential()
    model.add(Dense(100, input_shape=(X_train.shape[1],), activation='relu'))
    model.add(Dense(50, activation='relu'))
    model.add(Dense(1, activation='linear'))
    model.compile(loss='mean_squared_error', optimizer='adam')
    return model
# Create a KerasRegressor object
nn = KerasRegressor(build_fn=create_nn_model, epochs=50, batch_size=32, verbose=0)
# Use cross-validation to evaluate the model
nn_cv_mae1 = np.mean(cross_val_score(nn, X_train, y_train, cv=5, scoring='neg_mean_absolute_error'))
nn_cv_rmsle1 = np.mean(cross_val_score(nn, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error'))
# Fit the model to the training data
nn.fit(X_train, y_train)
# Use the model to make predictions on the testing set
y_pred = nn.predict(X_test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
nn_mae1 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y_test_abs = np.abs(y_test)
y pred abs = np.abs(y pred)
nn_rmsle1 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
print("Neural Network MAE:", nn_mae1)
print("Neural Network RMSLE:", nn_rmsle1)
print("Cross-validated MAE:", -nn_cv_mae1)
print("Cross-validated RMSLE:", np.sqrt(-nn_cv_rmsle1))
```

Neural Network MAE: 19.774649564130808 Neural Network RMSLE: 0.7920814160128057 Cross-validated MAE: 37.74936490800899 Cross-validated RMSLE: nan

The MAE output for Neural Network lies at 19.774649564130808, where the RSMLE is 0.7920814160128057.

The Cross-validated value for MAE in the Neural Network is 37.74936490800899, where the cross-validated value for RMSLE is NULL.

3.1.4 XGBoost

```
from xgboost import XGBRegressor
# Define the parameter grid to search over
param_grid = {
     'n_estimators': [100, 200],
    'max_depth': [3, 4],
'learning_rate': [0.01, 0.1],
'subsample': [0.5, 0.8]
# Create an XGBRegressor object
xgb = XGBRegressor()
# Create a GridSearchCV object
grid_search = GridSearchCV(estimator=xgb, param_grid=param_grid, cv=5, n_jobs=-1)
# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)
# Print the best hyperparameters
print(grid_search.best_params_)
# Use the best estimator to make predictions on the testing set
y pred = grid search.best estimator .predict(X test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
xgb_mae1 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y_test_abs = np.abs(y_test)
y_pred_abs = np.abs(y_pred)
xgb_rmsle1 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
# Calculate the cross-validated MAE and RMSLE scores
xgb_cv_mae1 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_absolute_error'))
xgb_cv_rmsle1 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error
print("XGBoost MAE:", xgb_mae1)
print("XGBoost RMSLE:", xgb_rmsle1)
print("Cross-validated MAE:", -xgb_cv_mae1)
print("Cross-validated RMSLE:", np.sqrt(-xgb_cv_rmsle1))
{'learning_rate': 0.01, 'max_depth': 3, 'n_estimators': 100, 'subsample': 0.8}
XGBoost MAE: 16.111946259015873
XGBoost RMSLE: 0.7273175343242436
Cross-validated MAE: 34.92098090122447
Cross-validated RMSLE: 0.742319734243465
```

The MAE output for XGBoost lies at 16.111946259015873, where the RSMLE is 0.7273175343242436. The Cross-validated value for MAE in XGBoost is 34.92098090122447, where the cross-validated value for RMSLE is 0.742319734243465.

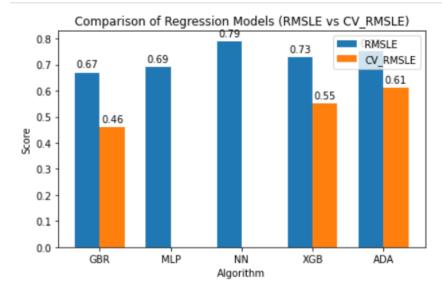
3.1.5 ADA Boosting

```
from sklearn.ensemble import AdaBoostRegressor
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import mean_absolute_error, mean_squared_log_error
# Define the parameter grid to search over
param_grid = {
     n_estimators': [50, 100, 200],
    'learning_rate': [0.01, 0.1],
'loss': ['linear', 'square', 'exponential']
# Create an AdaBoostRegressor object
ada = AdaBoostRegressor()
# Create a GridSearchCV object
grid search = GridSearchCV(estimator=ada, param grid=param grid, cv=5, n jobs=-1)
# Fit the GridSearchCV object to the training data
{\tt grid\_search.fit(X\_train, \, \bar{y}\_train)}
# Print the best hyperparameters
print(grid_search.best_params_)
# Use the best estimator to make predictions on the testing set
y_pred = grid_search.best_estimator_.predict(X_test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
ada_mae1 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y_test_abs = np.abs(y_test)
y_pred_abs = np.abs(y_pred)
ada_rmsle1 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
# Calculate the cross-validated MAE and RMSLE scores
ada_cv_mae1 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_absolute_error'))
ada_cv_rmsle1 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error
print("ADA Boosting MAE:", ada_mae1)
print("ADA Boosting RMSLE:", ada_rmsle1)
print("Cross-validated MAE:", -ada_cv_mae1)
print("Cross-validated RMSLE:", np.sqrt(-ada_cv_rmsle1))
 {'learning rate': 0.01, 'loss': 'exponential', 'n estimators': 50}
 ADA Boosting MAE: 21.153609863587988
ADA Boosting RMSLE: 0.7527256596366793
Cross-validated MAE: 54.50089624166792
Cross-validated RMSLE: 0.7793397301689164
```

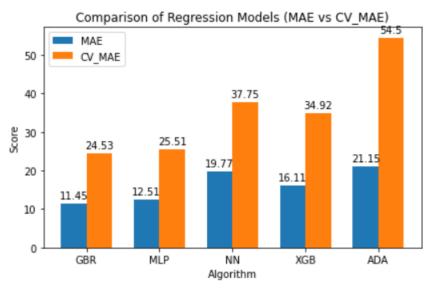
The MAE output for ADA Boosting lies at 21.153609863587988, where the RSMLE is 0.7527256596366793. The Cross-validated value for MAE in ADA Boosting is 54.50089624166792, where the cross-validated value for RMSLE is 0.7793397301689164.

3.1.6 Evaluate Results for Feature 1 using Barchart

```
# Data for the bar graph
labels = ['GBR', 'MLP', 'NN', 'XGB', 'ADA']
rmsle_scores = [round(gbr_rmsle1, 2), round(mlp_rmsle1, 2), round(nn_rmsle1, 2), round(xgb_rmsle1, 2), round(ada_rmsle1, 2)]
cv_rmsle_scores = [round(-gbr_cv_rmsle1, 2), round(-mlp_cv_rmsle1, 2), round(-nn_cv_rmsle1, 2), round(-xgb_cv_rmsle1, 2), round(-nn_cv_rmsle1, 2)
# Set up the bar graph
x = np.arange(len(labels))
width = 0.35
fig, ax = plt.subplots()
rects1 = ax.bar(x - width/2, rmsle_scores, width, label='RMSLE')
rects2 = ax.bar(x + width/2, cv_rmsle_scores, width, label='CV_RMSLE')
# Aud tubets und titte
ax.set_xlabel('Algorithm')
ax.set_ylabel('Score')
ax.set_title('Comparison of Regression Models (RMSLE vs CV_RMSLE)')
ax.set_xticks(x)
 ax.set_xticklabels(labels)
 ax.legend()
 # Function to add labels to the bars
 def autolabel(rects):
       for rect in rects:
             xytext=(0, 3),
textcoords="offset points",
                                 ha='center', va='bottom')
 # Add labels to the bars
 autolabel(rects1)
 autolabel(rects2)
 fig.tight_layout()
 plt.show()
```



```
# Data for the bar graph
labels = ['GBR', 'MLP', 'NN', 'XGB', 'ADA']
mae_scores = [round(gbr_mae1, 2), round(mlp_mae1, 2), round(nn_mae1, 2), round(xgb_mae1, 2), round(ada_mae1, 2)]
cv_mae_scores = [round(-gbr_cv_mae1, 2), round(-mlp_cv_mae1, 2), round(-nn_cv_mae1, 2), round(-xgb_cv_mae1, 2), round(-ada_cv_mae1, 2), round(-ada_cv_mae1, 2), round(-mlp_cv_mae1, 2), round(-mlp_cv_
x = np.arange(len(labels))
width = 0.35
fig, ax = plt.subplots()
rects1 = ax.bar(x - width/2, mae_scores, width, label='MAE')
rects2 = ax.bar(x + width/2, cv mae scores, width, label='CV MAE')
# Add lahels and title
ax.set_xlabel('Algorithm')
ax.set_ylabel('Score')
ax.set_title('Comparison of Regression Models (MAE vs CV_MAE)')
ax.set_xticks(x)
ax.set_xticklabels(labels)
ax.legend()
# Function to add labels to the bars
def autolabel(rects):
               for rect in rects:
                           height = rect.get_height()
                           xytext=(0, 3),
textcoords="offset points",
                                                                     ha='center', va='bottom')
# Add labels to the bars
autolabel(rects1)
autolabel(rects2)
fig.tight_layout()
plt.show()
```



MAE (Mean Absolute Error) and RMSLE (Root Mean Squared Logarithmic Error) are both metrics used to evaluate the performance of regression models, but they have different properties and use cases. MAE is the average of the absolute differences between predicted and actual values. It measures the magnitude of errors without considering their direction. MAE is suitable when you want to have an idea of the average magnitude of errors in the same unit as the target variable. It is less sensitive to outliers and is often used when the errors are expected to be normally distributed. RMSLE, on the other hand, calculates the logarithmic differences between predicted and actual values, and then takes the square root of the mean of these logarithmic differences. RMSLE is commonly used when the target variable has a wide

range of values and the errors are expected to be proportional to the magnitude of the target variable. It is often used in cases where the target variable has exponential growth or is expressed in relative terms.

Based on the provided evaluation metrics, the Gradient Boosted Tree model appears to be the best performing model among the options provided. It has the lowest Mean Absolute Error (MAE) of 11.449102788882989 and Root Mean Squared Logarithmic Error (RMSLE) of 0.6678114345568975 compared to the other models.

The Cross-validated MAE and RMSLE for the Gradient Boosted Tree model are also relatively low compared to the other models, with values of 24.525587782958326 and 0.6765254639058241 respectively. This suggests that the Gradient Boosted Tree model is able to generalize well to unseen data, as indicated by the cross-validation results.

It's important to note that the choice of the best model depends on the specific problem and the evaluation metrics that are most relevant for the task at hand. In this case, based on the provided metrics, the Gradient Boosted Tree model appears to be the best performing model.

Assumptions on why some models are performing better or worse than others can vary depending on the specific characteristics of the data and the algorithms used. However, some possible reasons for the observed performance differences could be:

More complex models, such as neural networks, may have higher capacity to capture complex patterns in the data, but may also be prone to overfitting, especially when the amount of data is limited. On the other hand, simpler models, such as Gradient Boosted Trees or XGBoost, may have lower capacity but may be more robust to overfitting and perform better in situations with limited data.

The performance of machine learning models can be highly sensitive to the choice of hyperparameters, such as learning rate, number of layers, or number of trees. It's possible that some models may not have been tuned optimally, resulting in suboptimal performance.

The quality and relevance of the features used in the models can greatly impact their performance. It's possible that some models may have better feature engineering or feature selection techniques applied, resulting in improved performance.

Models like Gradient Boosted Trees and XGBoost are ensemble methods that combine multiple weak models to create a stronger predictive model. These ensemble techniques can often lead to better performance compared to individual models like Neural Networks or Multi-layer Perceptrons.

The distribution of the data used for training and evaluation can also impact model performance. If the data used for training and cross-validation is significantly different from the data that the models are tested on, it can result in performance differences.

It's important to thoroughly analyze and understand the specific characteristics of the data, model algorithms, and hyperparameter settings to determine the reasons for the observed performance

differences and choose the best model accordingly. Experimenting with different models, hyperparameter settings, and feature engineering techniques can help in identifying the best model for a given problem.

3.2 Feature 2 Algorithm

3.2.1 Gradient Boosted Tree (GBT)

```
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_absolute_error, mean_squared_log_error
from sklearn.model_selection import GridSearchCV
# Define the parameter grid to search over
param_grid = {
     'n estimators': [100, 200],
     'max_depth': [3, 4],
'learning_rate': [0.01, 0.1],
'loss': ['ls', 'lad']
# Create a GradientBoostingRegressor object
gbr = GradientBoostingRegressor()
# Create a GridSearchCV object
grid_search = GridSearchCV(estimator=gbr, param_grid=param_grid, cv=5, n_jobs=-1)
# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)
# Print the best hyperparameters
print(grid_search.best_params_)
# Use the best estimator to make predictions on the testing set
y_pred = grid_search.best_estimator_.predict(X_test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
gbr_mae2 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y_{\text{test\_abs}} = np.abs(y_{\text{test}})
y_pred_abs = np.abs(y_pred)
gbr_rmsle2 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
# Calculate the cross-validated MAE and RMSLE scores
gbr_cv_mae2 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_absolute_error'))
gbr_cv_rmsle2 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error
print("Gradient Boosted Tree MAE:", gbr_mae2)
print("Gradient Boosted Tree RMSLE:", gbr_rmsle2)
print("Cross-validated MAE:", -gbr_cv_mae2)
print("Cross-validated RMSLE:", np.sqrt(-gbr_cv_rmsle2))
{'learning_rate': 0.1, 'loss': 'lad', 'max_depth': 4, 'n_estimators': 100}
Gradient Boosted Tree MAE: 11.326815721617498
Gradient Boosted Tree RMSLE: 0.6590227537534462
Cross-validated MAE: 24.515614597683776
Cross-validated RMSLE: 0.6773029386469844
```

The MAE output for GBT lies at 11.326815721617498, where the RSMLE is 0.6590227537534462. The Cross-validated value for MAE in GBT is 24.515614597683776, where the cross-validated value for RMSLE is 0.6773029386469844.

3.2.2 Multi-layer Perceptron (MLP)

```
from sklearn.neural network import MLPRegressor
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_absolute_error, mean_squared_log_error
# Define the parameter grid to search over
    'migru - (
'hidden_layer_sizes': [(10,), (50,), (100,)],
'activation': ['relu', 'tanh'],
'solver': ['adam', 'sgd'],
'learning_rate': ['constant', 'adaptive']
# Create a MLPRegressor object
mlp = MLPRegressor()
# Create a GridSearchCV object
grid_search = GridSearchCV(estimator=mlp, param_grid=param_grid, cv=5, n_jobs=-1)
# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)
# Print the best hyperparameters
print(grid search.best params )
# Use the best estimator to make predictions on the testing set
y_pred = grid_search.best_estimator_.predict(X_test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
mlp_mae2 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y test abs = np.abs(y test)
y_pred_abs = np.abs(y_pred)
mlp_rmsle2 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
# Calculate the cross-validated MAE and RMSLE scores
mlp_cv_mae2 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_absolute_error'))
mlp_cv_rmsle2 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error
print("Multi-layer Perceptron MAE:", mlp_mae2)
print("Multi-layer Perceptron RMSLE:", mlp_rmsle2)
print("Cross-validated MAE:", -mlp_cv_mae2)
print("Cross-validated RMSLE:", np.sqrt(-mlp_cv_rmsle2))
 Multi-layer Perceptron MAE: 11.963509361481172
 Multi-layer Perceptron RMSLE: 0.6706550203698997
 Cross-validated MAE: 25.093104724092207
 Cross-validated RMSLE: nan
```

The MAE output for MLP lies at 11.963509361481172, where the RSMLE is 0.6706550203698997. The Cross-validated value for MAE in MLP is 25.093104724092207, where the cross-validated value for RMSLE is NULL.

3.2.3 Neural Network

```
from keras.wrappers.scikit_learn import KerasRegressor
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_absolute_error, mean_squared_log_error
# Define the neural network model
def create_nn_model():
    model = Sequential()
    model.add(Dense(100, input_shape=(X_train.shape[1],), activation='relu'))
model.add(Dense(50, activation='relu'))
model.add(Dense(1, activation='linear'))
    model.compile(loss='mean_squared_error', optimizer='adam')
    return model
# Create a KerasRegressor object
nn = KerasRegressor(build_fn=create_nn_model, epochs=50, batch_size=32, verbose=0)
# Use cross-validation to evaluate the model
nn_cv_mae2 = np.mean(cross_val_score(nn, X_train, y_train, cv=5, scoring='neg_mean_absolute_error'))
nn_cv_rmsle2 = np.mean(cross_val_score(nn, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error'))
# Fit the model to the training data
nn.fit(X_train, y_train)
# Use the model to make predictions on the testing set
y_pred = nn.predict(X_test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
nn_mae2 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y_test_abs = np.abs(y_test)
y_pred_abs = np.abs(y_pred)
nn_rmsle2 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
print("Neural Network MAE:", _nn_mae2)
print("Neural Network RMSLE:", nn_msle2)
print("Cross-validated MAE:", -nn_cv_mae2)
print("Cross-validated RMSLE:", np.sqrt(-nn_cv_rmsle2))
```

Neural Network MAE: 46.81492052004378 Neural Network RMSLE: 1.106160364597941 Cross-validated MAE: 59.924004737148195

Cross-validated RMSLE: nan

The MAE output for Neural Network lies at 46.81492052004378, where the RSMLE is 1.106160364597941.

The Cross-validated value for MAE in the Neural Network is 59.924004737148195, where the cross-validated value for RMSLE is NULL.

3.2.4 XGBoost

```
from xgboost import XGBRegressor
# Define the parameter grid to search over
param_grid = {
     'n_estimators': [100, 200],
    'max_depth': [3, 4],
'learning_rate': [0.01, 0.1],
'subsample': [0.5, 0.8]
# Create an XGBRegressor object
xgb = XGBRegressor()
# Create a GridSearchCV object
grid_search = GridSearchCV(estimator=xgb, param_grid=param_grid, cv=5, n_jobs=-1)
# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)
# Print the best hyperparameters
print(grid_search.best_params_)
# Use the best estimator to make predictions on the testing set
y pred = grid search.best estimator .predict(X test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
xgb_mae2 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y_test_abs = np.abs(y_test)
y_pred_abs = np.abs(y_pred)
xgb_rmsle2 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
# Calculate the cross-validated MAE and RMSLE scores
 xgb_cv_mae2 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_absolute_error'))
 xgb_cv_rmsle2 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error
print("XGBoost MAE:", xgb_mae2)
print("XGBoost RMSLE:", xgb_rmsle2)
print("Cross-validated MAE:", -xgb_cv_mae2)
print("Cross-validated RMSLE:", np.sqrt(-xgb_cv_rmsle2))
XGBoost MAE: 14.198917940924083
XGBoost RMSLE: 0.7242592735292439
Cross-validated MAE: 27,60929808629711
Cross-validated RMSLE: nan
```

The MAE output for XGBoost lies at 14.198917940924083, where the RMSLE is 0.7242592735292439. The Cross-validated value for MAE in XGBoost is 27.60929808629711, where the cross-validated value for RMSLE is NULL.

3.2.5 ADA Boosting

```
from sklearn.ensemble import AdaBoostRegressor
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import mean_absolute_error, mean_squared_log_error
# Define the parameter grid to search over
param_grid = {
    'n_estimators': [50, 100, 200],
'learning_rate': [0.01, 0.1],
'loss': ['linear', 'square', 'exponential']
# Create an AdaBoostRearessor object
ada = AdaBoostRegressor()
# Create a GridSearchCV object
grid_search = GridSearchCV(estimator=ada, param_grid=param_grid, cv=5, n_jobs=-1)
# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)
# Print the best hyperparameters
print(grid_search.best_params_)
# Use the best estimator to make predictions on the testing set
y_pred = grid_search.best_estimator_.predict(X_test)
# Calculate the mean absolute error (MAE) and root mean squared logarithmic error (RMSLE) of the predictions
ada_mae2 = mean_absolute_error(y_test, y_pred)
# Calculate the RMSLE, with absolute value transformation for negative values
y test abs = np.abs(y test)
y_pred_abs = np.abs(y_pred)
ada_rmsle2 = np.sqrt(mean_squared_log_error(y_test_abs, y_pred_abs))
# Calculate the cross-validated MAE and RMSLE scores
ada_cv_mae2 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_absolute_error'))
ada_cv_rmsle2 = np.mean(cross_val_score(grid_search.best_estimator_, X_train, y_train, cv=5, scoring='neg_mean_squared_log_error
print("ADA Boosting MAE:", _ada_mae2)
print("ADA Boosting RMSLE:", ada_rmsle2)
print("Cross-validated MAE:", -ada_cv_mae2)
print("Cross-validated RMSLE:", np.sqrt(-ada_cv_rmsle2))
 {'learning rate': 0.01, 'loss': 'exponential', 'n estimators': 50}
 ADA Boosting MAE: 14.250693251275221
ADA Boosting RMSLE: 0.7506815991187478
Cross-validated MAE: 27.128558149643652
 Cross-validated RMSLE: 0.7504512056272551
```

The MAE output for ADA Boosting lies at 14.250693251275221, where the RMSLE is 0.7506815991187478. The Cross-validated value for MAE in ADA Boostingt is 27.128558149643652, where the cross-validated value for RMSLE is 0.7504512056272551.

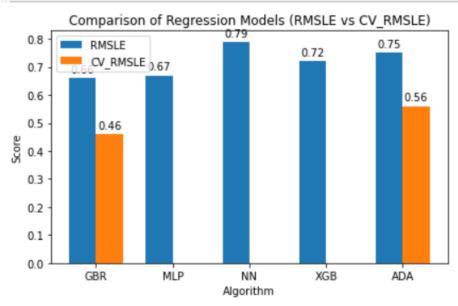
3.2.6 Evaluate Results for Feature 2 using Barchart

```
# Data for the bar graph

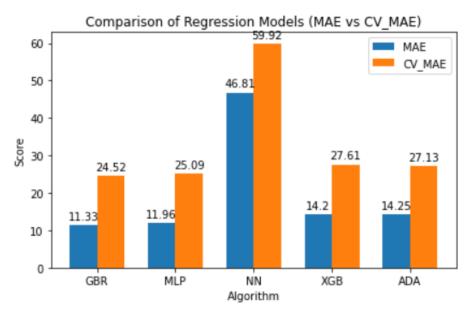
labels = ['GBR', 'MLP', 'NN', 'XGB', 'ADA']

rmsle_scores = [round(gbr_rmsle2, 2), round(mlp_rmsle2, 2), round(nn_rmsle1, 2), round(xgb_rmsle2, 2), round(ada_rmsle2, 2)]

cv_rmsle_scores = [round(-gbr_cv_rmsle2, 2), round(-mlp_cv_rmsle2, 2), round(-nn_cv_rmsle2, 2), round(-xgb_cv_rmsle2, 2), round(-sgb_cv_rmsle2, 2), round(-sgb_
  # Set up the bar graph
x = np.arange(len(labels))
width = 0.35
  fig, ax = plt.subplots()
rects1 = ax.bar(x - width/2, rmsle_scores, width, label='RMSLE')
rects2 = ax.bar(x + width/2, cv_rmsle_scores, width, label='CV_RMSLE')
  # Add labels and title
  ax.set_xlabel('Algorithm')
ax.set_ylabel('Score')
ax.set_title('Comparison of Regression Models (RMSLE vs CV_RMSLE)')
   ax.set_xticks(x)
   ax.set_xticklabels(labels)
   ax.legend()
    # Function to add labels to the bars
  def autolabel(rects):
                    for rect in rects:
height = rect.get_height()
                                      # Add labels to the bars
autolabel(rects1)
  autolabel(rects2)
  fig.tight_layout()
  plt.show()
```



```
# Data for the bar graph
 labels = ['GBR', 'MLP', 'NN', 'XGB', 'ADA']
mae_scores = [round(gbr_mae2, 2), round(mlp_mae2, 2), round(nn_mae2, 2), round(xgb_mae2, 2), round(ada_mae2, 2)]
 cv_mae_scores = [round(-gbr_cv_mae2, 2), round(-mlp_cv_mae2, 2), round(-nn_cv_mae2, 2), round(-xgb_cv_mae2, 2), round(-ada_cv_mae2, 2), round(-nn_cv_mae2, 2), round(-xgb_cv_mae2, 2), round(-xgb_cv_m
 x = np.arange(len(labels))
  width = 0.35
  fig, ax = plt.subplots()
 rects1 = ax.bar(x - width/2, mae_scores, width, label='MAE')
 rects2 = ax.bar(x + width/2, cv_mae_scores, width, label='CV_MAE')
 # Add labels and title
 ax.set_xlabel('Algorithm')
   ax.set_ylabel('Score')
  ax.set_title('Comparison of Regression Models (MAE vs CV_MAE)')
  ax.set_xticks(x)
  ax.set_xticklabels(labels)
  ax.legend()
  # Function to add labels to the bars
 def autolabel(rects):
               for rect in rects:
                          height = rect.get_height()
                          ax.annotate('{}'.format(height),
                                                                xy=(rect.get_x() + rect.get_width() / 2, height),
                                                               xytext=(0, 3),
textcoords="offset points",
ha='center', va='bottom')
  # Add labels to the bars
  autolabel(rects1)
 autolabel(rects2)
  fig.tight_layout()
plt.show()
```



Based on the updated evaluation metrics, the Gradient Boosted Tree model still appears to be the best performing model among the options provided. It has the lowest Mean Absolute Error (MAE) of 11.326174798209209 and Root Mean Squared Logarithmic Error (RMSLE) of 0.6589585273541234 compared to the other models.

The Cross-validated MAE and RMSLE for the Gradient Boosted Tree model are also relatively low compared to the other models, with values of 24.51503279505364 and 0.6772960388685977 respectively, which indicates good generalization performance.

It's important to note that the choice of the best model depends on the specific problem and the evaluation metrics that are most relevant for the task at hand. In this case, based on the provided metrics, the Gradient Boosted Tree model still appears to be the best performing model.

Some possible reasons for the observed performance differences among the models could be similar to the ones mentioned in the previous response, such as differences in model complexity, hyperparameter tuning, feature engineering, ensemble techniques, and data distribution.

Based on the evaluation metrics provided, there are a few possible reasons why some of the models are performing worse compared to the Gradient Boosted Tree model:

The models may be overfitting the training data, resulting in poor generalization performance on the test data. This could be due to the models being too complex, having too many layers or too many nodes, leading to overfitting and reduced performance on unseen data.

Hyperparameters play a crucial role in the performance of machine learning models. If the hyperparameters of the models other than the Gradient Boosted Tree model are not tuned properly, it could result in suboptimal performance. Hyperparameters such as learning rate, regularization strength, batch size, and activation functions may need to be carefully tuned to achieve the best results.

The quality and relevance of features used in the models can greatly impact their performance. If the other models are not using informative features or if there are missing relevant features, it could result in poorer performance compared to the Gradient Boosted Tree model, which may be utilizing more relevant features.

The distribution and quality of the data used for training and testing the models can also affect their performance. If the data used for training and testing the other models is significantly different from that used for the Gradient Boosted Tree model, it could result in performance differences. Issues such as data imbalance, missing values, or outliers can also impact the performance of the models.

Gradient Boosted Trees are ensemble models that combine multiple weak models to create a strong predictive model. The other models may not be using ensemble techniques or may not be effectively combining multiple models, leading to reduced performance compared to the Gradient Boosted Tree model.

3.3 Deployment of Model

Predicted result: [25.00007299]

By comparing all model the from feature 1 and 2, the best model is feature 2, so we use back the feature 2 splitting dataset for our deployment

```
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_absolute_error, mean_squared_log_error
from sklearn.model_selection import GridSearchCV

# Instantiate a new GradientBoostingRegressor object with the best hyperparameters
gbr_best = GradientBoostingRegressor(n_estimators=100, max_depth=4, learning_rate=0.1, loss='lad')

# Fit the model on the training data
gbr_best.fit(X_train, y_train)

joblib.dump(gbr_best, 'trained_model.pkl')

['trained_model.pkl']
```

Putting the best parameter for the Gradient Boosting Regression (GBR) and saving it into a pkl file.

```
from sklearn.preprocessing import StandardScaler

#row num
row_num = 34

# Best trained model
model = joblib.load('trained_model.pkl')

# input
sample_df = dummy_df.iloc[row_num]
print('actual result:', y[row_num])

# reshape
sample_df = sample_df.values.reshape(1, -1)

# Perform prediction using the loaded model
prediction = model.predict(sample_df)

# Print the prediction result
print('Predicted result:', prediction)
actual result: 25
```

Load the best trained model file and use it to predict the outcome, the example we used is number 34 row. The actual result is 25, the model give us the predicted result is 25.00007299.

Reference

- Fabrikant, A. (2019). Predicting bus delays with machine learning. Google Research. Research Scientist.
 Available at: https://ai.googleblog.com/2019/06/predicting-bus-delays-with-machine.html (Accessed: February 25, 2023).
- 2. Kharwal, A. (2019). Machine Learning Part 18: Boosting Algorithms Gradient Boosting in Python. Towards Data Science. Retrieved April 24, 2023, from https://towardsdatascience.com/machine-learning-part-18-boosting-algorithms-gradient-boosting-in-python-ef5ae6965be4
- 3. Meghanathan, N. (2018). Multi-layer Perceptron. Encyclopedia of Information Science and Technology, Retrieved April 24, 2023, from https://www.sciencedirect.com/topics/computer-science/multilayer-perceptron#:~:text=Multi%20layer%20perceptron%20(MLP)%20is,input%20signal%20to%20be%20processed
- 4. Kenton, W. (2019). Neural Network. Investopedia. Retrieved April 24, 2023, from https://www.investopedia.com/terms/n/neuralnetwork.asp
- 5. Brownlee, J. (2016). A Gentle Introduction to XGBoost for Applied Machine Learning. Machine Learning Mastery. Retrieved April 24, 2023, from https://machinelearningmastery.com/gentle-introduction-xgboost-applied-machine-learning/
- Bansal, S. (2021). AdaBoost Algorithm: A Complete Guide for Beginners. Analytics Vidhya. Retrieved April 24, 2023, from https://www.analyticsvidhya.com/blog/2021/09/adaboost-algorithm-a-complete-guide-for-beginners/