# Technical Report: Final Project DS 5110: Introduction to Data Management and Processing

Team Members:
Rishi Kamtam and Sabari Mathavan
Khoury College of Computer Sciences
Data Science Program
kamtam.r@northeastern.edu, ramasamybalasubram.s@northeastern.edu

December 11, 2024

## Contents

# 1   Introduction

**Background**: Living in a large city like Boston where common use of a car is not very popular, many of us rely on modes of public transportation like buses and trains. However, with the growing popularity of Ride share apps like Uber and Lyft, the market dynamic has greatly shifted with over 40 million ride share fares in Boston this past year. A large issue with ride shares however is pricing and variability across different apps like Uber and Lyft, and how prices are actually determined. Many times, we see drastically different prices across different locations like the North End, South End, and Back Bay which is the biggest motivation behind our project, to determine what really impacts the variability in ride share prices.

**Objectives**: For our project, we had 2 main objectives that would help us address both the instructions and requirements of the project as well as exploring the issues in price variability as mentioned in the background.

**1. Develop a Prediction Model for Ride share Fares in Boston** This involves utilizing Machine Learning Models, specifically supervised models that take in previous features that impact the target variable (price). This was a large objective for us because we believe that it would be helpful for students like us, who are on tight budgets, to have an application that can predict fare prices to see the most economical methods of transportation to wherever we are planning on traveling to.

**2. Enhance Fare Estimation Accuracy**  This second objective involves doing exploratory data analysis to explore and identify common trends across our data set to give the users a comprehensive overview of our machine learning model. By looking into the correlation between features, differences in average and median pricing between both Uber and Lyft, it allows us and the user to better understand how the ML model actually predicts fare prices.

**Goals:**

- Deliver a solution that benefits the consumer

- Use Advanced ML Data Science Tools and Visualizations

**Scope:** The scope of our project first starts with utilizing a dataset from both Uber and Lyft that includes the source, destination, price, distance, car type, and time of hundreds of thousands of rides. Next, our project involves matching weather data to the rides based on the timestamps. This further involves cleaning both datasets then aggregating them together to match weather data with specific rides. The next step is visualizing trends and insights. This step is vital because it gives us and the user an understanding of general correlation and distribution metrics across out dataset, enhancing their understanding of the ML model that will be implemented. The final two steps relate to building multiple ML models and evaluating them through various common metrics and visualizations to see how the model actually performed in predicting ride share fare prices.

# 2    Literature Review

Existing work regarding predicting ride share fare prices includes utilizing Deep Learning, specifically Neural Networks, to predict prices. Important features for this research were distance traveled, time elapsed, and number of passengers. Based on the study's evaluation, they found that distance traveled, time elapsed, and number of passengers were the most important and had by far the largest impact compared to time of day, cab type, ect.. Feature Scaling and Selection were also large parts of existing work already done as they are important for distinguishing what metrics actually influence the pricing of ride shares.

A gap in this literature and how our project fits into the existing work is with the weather data and aggregating weather data to already existing ride share data. Weather has a large influence on demand of Rideshares and specific metrics like Rain and others could play a role in the variability of pricing. For example, someone is more likely to use ride share apps like Uber and Lyft when it is raining rather than walking or waiting outside for a train or bus. By matching rides to weather patterns, the ML model and exploratory analysis can provide insights into if Weather has a role in how Uber and Lyft determine prices.

# 3    Methodology

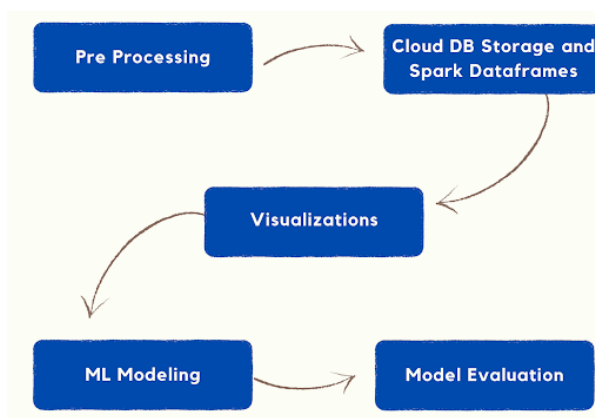The workflow for the project can be presented as below.



Figure 1: Working mechanism

## 3.1    Data Collection

The dataset utilized for this project is a static resource available on Kaggle (`https://www.kaggle.com/datasets/ravi72munde/uber-lyft-cab-prices`). It comprises two CSV files:

- **cab_rides.csv:** This file includes data on ride pricing, locations, and timestamps for both Uber and Lyft during the period from November 2018 to December 2018. The records were collected at intervals of 5 minutes, resulting in over 690,000 data points.

- **weather.csv:** This file contains weather metrics corresponding to the same time-frame as the `cab_rides.csv` file. The data were recorded at hourly intervals, amounting to more than 6,000 rows.

## 3.2   Data Preprocessing

The preprocessing phase ensured the dataset was clean and standardized for analysis. The steps included:

- **Rideshare Data (rides.csv):**

  - Renamed to "Rideshare Data" for clarity.
  - Removed null values to maintain data integrity.
  - Converted the `Time` column from Unix format to human-readable `Datetime`, truncating to the nearest 30 minutes.

- **Weather Data (weather.csv):**

  - Renamed to "Weather Data" for clarity.
  - Converted the `Time` column from Unix format to `Datetime`, truncated to the nearest 30 minutes.
  - Replaced null values in the `Rain` column with 0 to indicate no rain.

- Additional preprocessing tasks included removing duplicate entries, eliminating unnecessary columns, and splitting datasets as needed for analysis.

The datasets were transformed into Pandas DataFrames to facilitate the preprocessing tasks. The rides dataset was further divided into two subsets: *Uber Data* and *Lyft Data*, based on the ride-hailing service type. These subsets, along with the weather DataFrame, were uploaded to Google Cloud BigQuery using the BigQuery API.

The schema for the ride-based data table is as follows:

| Column Name | Data Type |
|---|---|
| distance | float |
| cab_type | string |
| time_stamp | datetime |
| destination | string |
| source | string |
| price | float |
| surge_multiplier | float |
| id | integer |
| name | string |

Table 1: Schema for the ride-based data table

The schema for the weather data table is as follows:

| Column Name | Data Type |
|---|---|
| temp | float |
| location | string |
| clouds | float |
| pressure | float |
| rain | float |
| time_stamp | datetime |
| humidity | float |
| wind | float |

Table 2: Schema for the weather data table

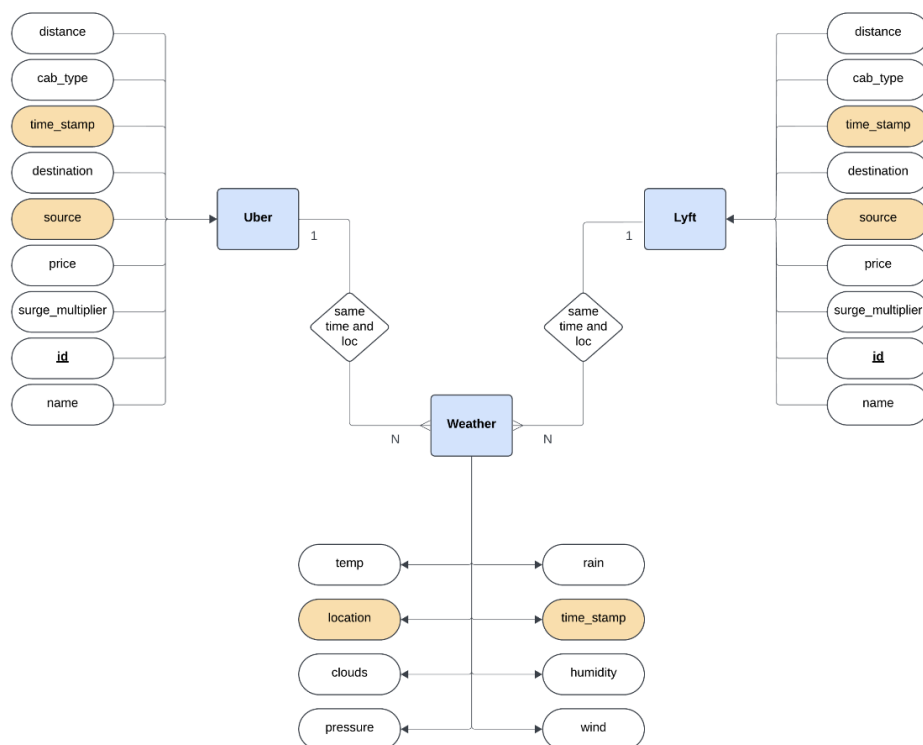Below is the entity-relationship diagram representing the tables:



Figure 2: ER Diagram

The preprocessed datasets were stored in Google Cloud BigQuery using the BigQuery API with the following table names:

- Uber ride data: `uber_data`

- Lyft ride data: `lyft_data`

- Weather data: `weather_data`

## 3.3   Analysis Techniques

The data stored in Google Cloud's BigQuery was accessed using PySpark to enable efficient and faster computations. For example, retrieving data from the `uber_data` table

using PySpark took 12 seconds, while using the BigQuery API for the same data retrieval took over 4 minutes.

Several analyses and visualizations were performed on the datasets, including:

- Statistical comparisons for the ride-hailing services (mean, median, standard deviations, totals, range)

- Box plots (Trip frequency vs Fare)

- Bar charts (Trip frequency vs Fare)

- Scatter plot (Distance vs Fare)

- Heat map (Weather metrics vs Price)

- Line-Bar chart (Rain's influence on Price)

- Line chart (Hourly pricing)

- Line heat map (Hourly x Distance x Pricing)

The visualizations were presented on webpages using the Flask framework. Python Flask is a lightweight, yet powerful web framework suitable for developing web applications. In addition to the aforementioned static visualizations, a dynamic visualization dashboard was included. The dynamic dashboard allows user interaction, where the user inputs source and destination points. The dashboard then calculates and provides the average pricing for each ride-hailing service across the given locations at an hourly level. Upon entering the source and destination, these values are processed and passed as parameters to the Google BigQuery API. The corresponding data (hourly pricing for each cab service at the given locations) is then returned and displayed in the dashboard as a line graph and a table.

As for the machine learning portion of our project, the methodology we took first involved defining the features and target variable. Features are variables that impact the target variable and in the context of our project, features include weather data such as rain, clouds, ride distance, etc.. and the target variable is the price because it is what we are trying to predict. After defining these, we had to implement feature engineering for the categorical data. This is an important step because a supervised ML model cannot read text, only interpret numbers so we had to change our categorical data to numerical data. Upon doing this, we were ready to implement our first ML model (Random Forest) which combines independent decision trees. Then, we cross validated the model which splits our dataset and trains it on unseen data to avoid over fitting within our model. Over fitting in machine learning refers to "creating a model that matches (memorizes) the training set so closely that the model fails to make correct predictions on new data." Then we built our second ML model (XGBoost) which works with decision trees but the main difference being it develops one tree at a time, correcting faults caused by previous tree. Finally, we evaluated our model using metrics such as variations of Mean Absolute Error scores and R-Squared scores as well as visualizations regarding how well our model predicted fare prices.

# 4   Results

**Data Analysis and Visualization Results:** Through the analysis, the following key insights were obtained:

- **Uber's Pricing:** It is evident that Uber offers lower average pricing and has a lower standard deviation (SD) compared to Lyft. This indicates that Uber's fare prices are generally more consistent and predictable.

- **Weather Metrics Influence:** The analysis indicates that weather metrics do not significantly affect ride pricing. Despite initial expectations, factors such as temperature, rain, and humidity showed little correlation with fare fluctuations.

- **Influence of Hour and Location:** The time of day and specific location have a significant impact on pricing. Hourly variations and location-based trends were found to play a crucial role in determining fare prices for ride-hailing services.

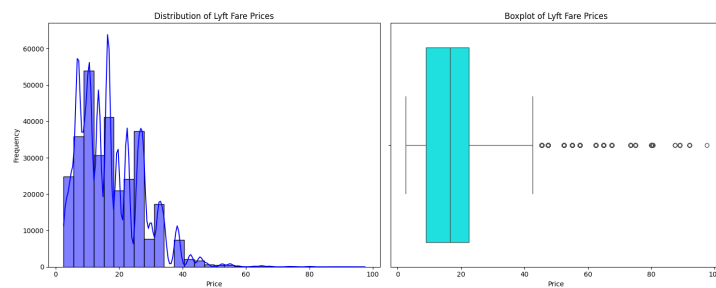Below are some visualization charts that demonstrate these findings:



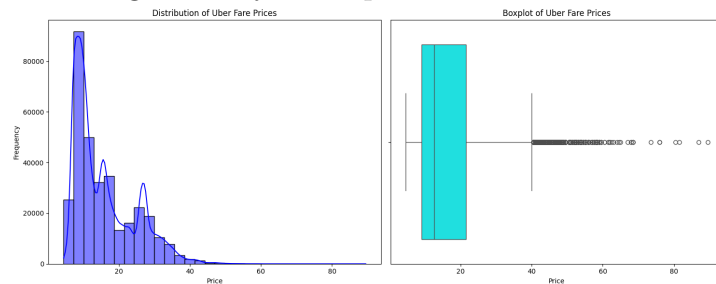Figure 3: Lyft's Trip-Price Distribution



Figure 4: Uber's Trip-Price Distribution

The analysis using Python revealed that the number of outliers in the Lyft data set represented 1.33% of the total data for Lyft, while for Uber, the outliers accounted only 0.91%.
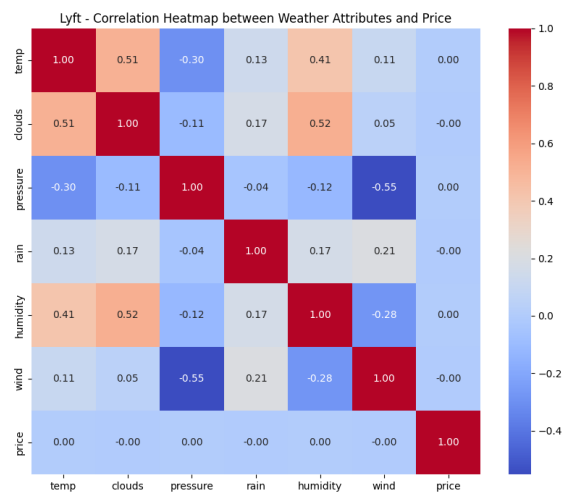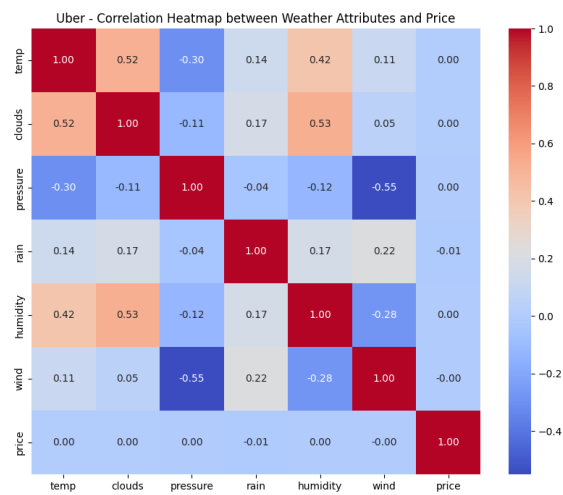
Figure 5: Lyft's Heatmap



Figure 6: Uber's Heatmap

As we could see, the weather metrics have a very minimal effect on pricing based on the obtained data set.
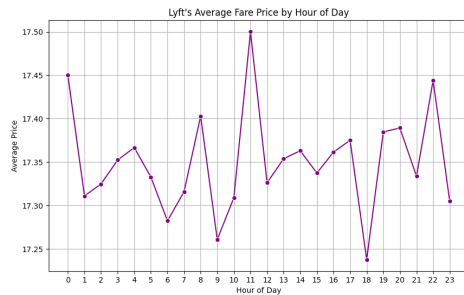


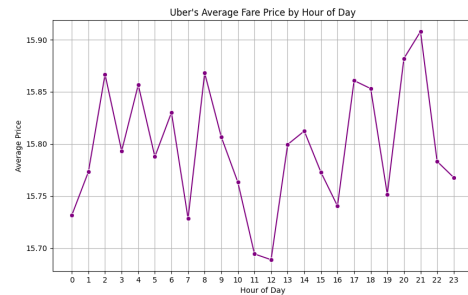Figure 7: Lyft's Hourly Pricing



Figure 8: Uber's Hourly Pricing

From the above graph, we can find that the average pricing of Uber is lower compared to Lyft's pricing across all hours of the day.
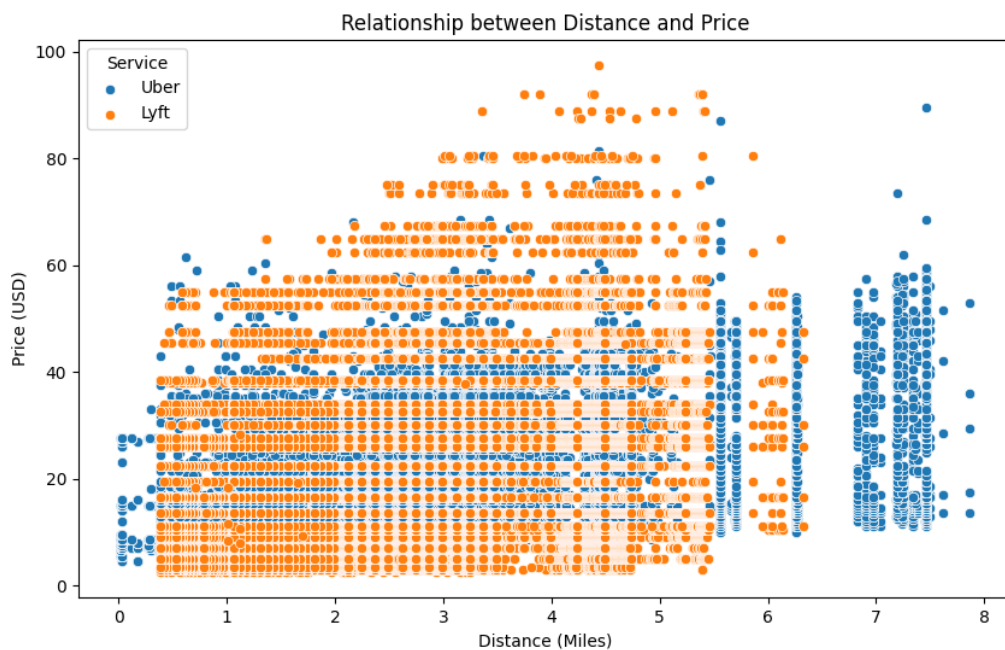


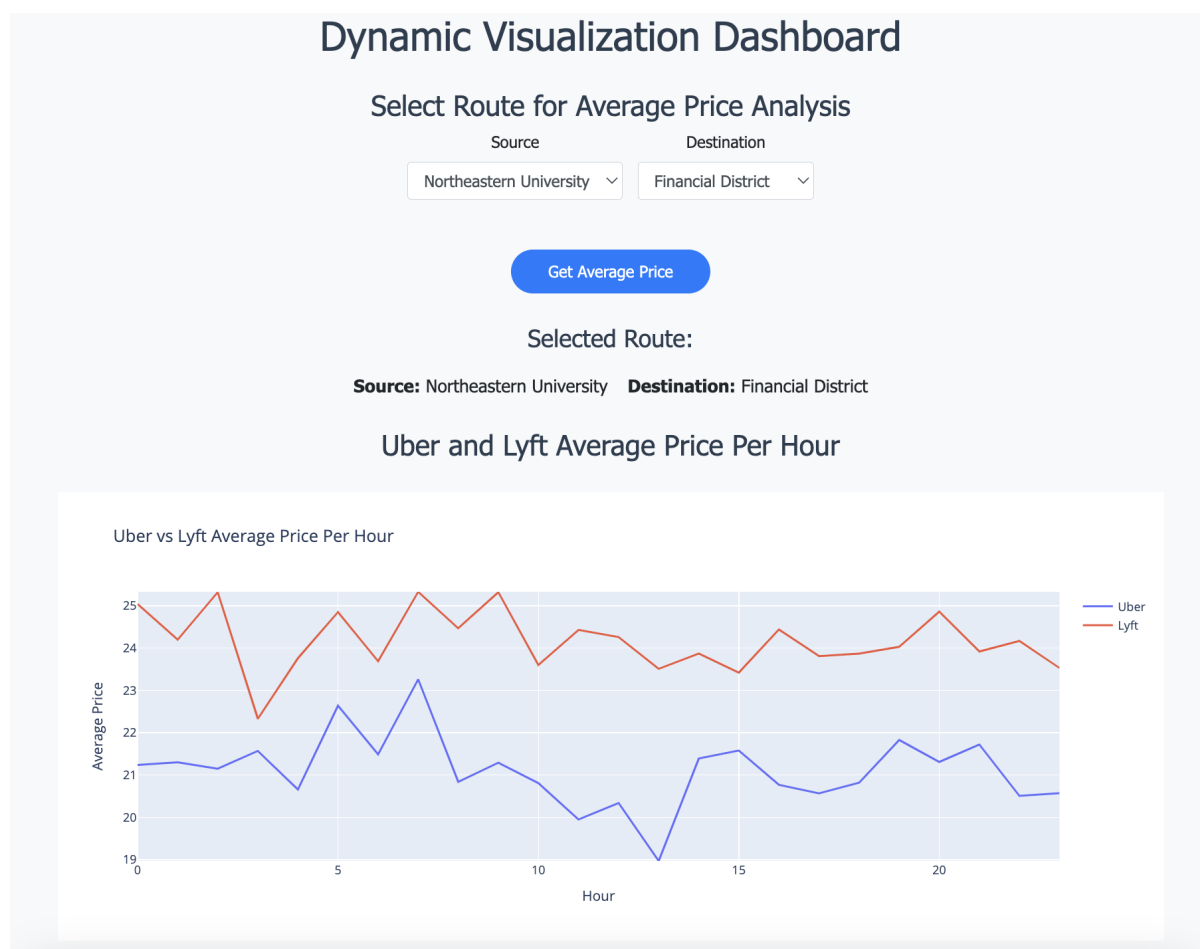Figure 9: Scatter Plot for Price vs Distance

Figure 10: Dynamic Visualization

The image above displays a snippet of the dynamic visualization feature. Based on the user-provided source and destination inputs, the graph dynamically adjusts to show the hourly pricing for each ride-hailing service between the two locations.

**Machine Learning Results:**

Considering the table of the results, the XGBoost Model had a better performance across all the metrics. Specifically, considering the lower Mean Absolute Error, we can conclude that the variability across pricing predictions was lower for the XGBoost model by about 0.4 dollars than the Random Forest Model. However, in a grander context, both models had a decent performance as there was still variability within the models' pricing predictions. Also in terms of the features, the R-Squared scores indicate that the features do not explain the pricing very well.

Note: Residual = Actual Y Value - Predicted Y Value Looking at the comparison of Residual Visualizations, it is evident that the XGBoost model performed better because there are less data points spread outside the red line (indicating perfect model performance). Also, it can be determined through the visualization that the model performed better for mid tier pricing (15-30) dollars as there is a higher concentration of data points in that range closer to the red line.

|                          | Random Forest Model | XGBoost Model |
|--------------------------|---------------------|---------------|
| Mean Absolute Error      | 7.11                | 6.69          |
| Mean Squared Error       | 74.94               | 64.75         |
| Root Mean Squared Error  | 8.65                | 8.04          |
| R^2 Score                | 0.029               | 0.11          |

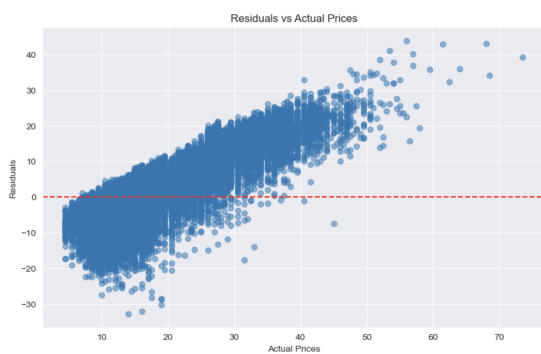Figure 11: Table of ML Model MAE and R-Squared Result Metrics



Figure 12: Residual Visualization for RF Model



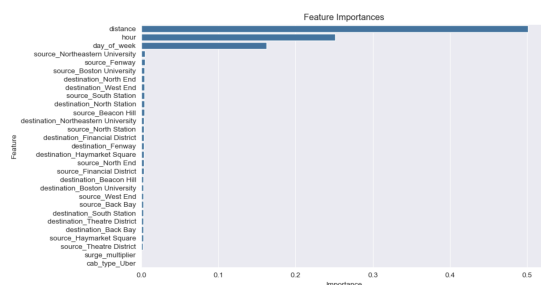Figure 13: Residual Visualization for XGBoost Model
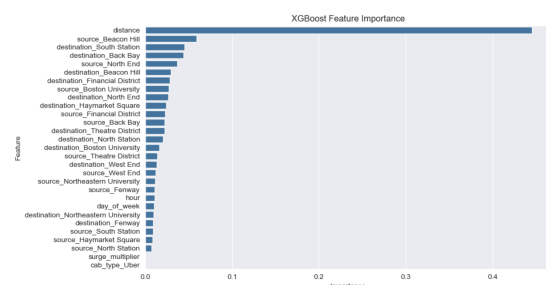


Figure 14: Feature Importance for RF Model



Figure 15: Feature Importance for XG-Boost Model

Upon building and analyzing the visualization of feature importance, one can conclude that distance is the most important in determining the price of a ride share for both models. Something that we found interesting was that the second most important feature was different for both models. For the RF model it was the hour (of the day) while for the XGBoost model it was the source (specifically Beacon Hill). The source of Beacon Hill being important for the XGBoost model we think can be explained by the high income of the Beacon Hill neighborhood and ride share prices most likely being high for that area. Overall, for both ML models, the results indicate that there was a decent performance especially for middle tier pricing and distance was the most important feature, impacting the models the most.

# 5   Discussion

The ride share prices predicted depend on different metrics, making it challenging to predict fares using machine learning models. Some factors include driver behavior, passenger preferences, location, time, weather conditions, etc.
Based on the visualization performed on the dataset accumulated, an inference can be made that there is roughly no correlation between price and specific weather metrics. This suggests that weather conditions (e.g., temperature, rain, wind, pressure, humidity, clouds) do not significantly influence fare. On the other hand, factors such as distance, location, and time were found to have a much larger impact on price fluctuations. For the better working of this project, the previously prescribed metrics should be prioritized for more accurate price predictions.

The project's results were restricted due to a few limitations, such as:

- **Inflation Impact:** The dataset is 6 years old (Nov 2018-Dec 2018) and does not account for inflation in rideshare prices over the years.

- **Changing Location Demographics:** Over the past 6 years, the demographic and economic makeup of the locations may have changed, causing a change in current pricing.
  The machine learning model developed for price prediction performs decently on the training data but struggles with unseen data. The model could not work properly when presented with geographical locations apart from the dataset, making the model look general and specific case-oriented to locations and time.

  Our findings align closely with those reported in the literature. While our Mean Absolute Error (MAE), a key metric for evaluating model performance, was 6.69, the literature reported a slightly better MAE of 5.96. This difference can be attributed to the use of Deep Learning Models in the literature, as opposed to our Supervised Machine Learning Model. Deep Learning Models are more effective at capturing and predicting non-linear relationships, suggesting that the features in our dataset may exhibit significant non-linearity with respect to the target variable (price).

  Notably, our analysis highlights the potential impact of weather on predictions. Although our results showed no correlation across weather metrics, this may indicate a complex, non-linear relationship that our model struggled to capture. Integrating

weather data into a Deep Learning Model could provide valuable insights, as these models are better equipped to handle non-linear interactions.

# 6    Conclusion

The results for the project were well executed, as mentioned in the scope of the project. Compiling them together:

- **Data Pre-processing:** Static data regarding rideshare pricing and weather conditions were accumulated from Kaggle. The data were preprocessed and cleaned to ensure that they could be stored and used for analysis.
- **Data Storage:** The data were stored in Google Cloud BigQuery as it involved cloud storage facilities that met the project's needs.
- **Data Analysis:** The data were later analyzed through PySpark to facilitate quicker and more reliable analysis on large datasets. The important findings from the analysis include:
  * It was found that weather metrics do not hold much significance in the pricing model. Despite initial expectations, factors such as temperature, rain, and humidity did not have a strong correlation with price fluctuations.
  * The scores of mean, median, and standard deviation favored Uber more than Lyft.
- **Data Visualization:** Both static and dynamic methods of data visualization were performed and showcased using a Flask webpage to ensure that no data leaks or unnecessary information were passed to the end user.
- **Predictive Modeling:** A predictive model was designed to predict the pricing of the ride-hailing services between two locations at a given time and weather conditions. The machine learning model performed well for mid-tier pricing but struggled with predicting low and high fare values accurately.

The final analysis suggests that Uber's pricing model tends to be more accurate and competitive compared to Lyft's pricing.

Building on the findings of this project, several areas for future improvement and expansion are proposed:

- **Incorporating Real-Time Data:** To enhance the accuracy and relevance of fare predictions, integrating live traffic, weather, and surge pricing data would be crucial. Real-time inputs would allow the model to adjust predictions dynamically, improving its utility in everyday rideshare operations.
- **Deep Learning Models:** The use of deep learning techniques could be explored to better capture non-linear relationships between pricing factors. This approach could improve model performance.
- **Geographical Expansion:** The model should be expanded to cover a wider geographical scope beyond Boston. Exploring areas such as suburban Boston and other major cities like New York and Chicago will provide a more comprehensive view of rideshare pricing trends across diverse regions.

Incorporating these changes would improve both the accuracy and scalability of the pricing model, making it applicable to a broader range of situations and locations.

# 7   References

## References

1. Connecting Apache Spark to BigQuery. *Google Cloud Documentation*, Google, https://cloud.google.com/bigquery/docs/connect-to-spark.  Accessed 3 Dec. 2024.

2. Random Forest Algorithm in Machine Learning. *GeeksforGeeks*, https://www.geeksforgeeks.o forest-algorithm-in-machine-learning/. Accessed 3 Dec. 2024.

3. Thomas, Moses Moncy. Predictive Analysis: Estimating UBER Fare Prices Using ML. *Medium*, 26 Apr. 2023, https://medium.com/@mosesmoncy1626/predictive-analysis-estimating-uber-fare-prices-using-ml-7bb8c54507e9.

4. XGBoost. *NVIDIA Glossary*, NVIDIA, https://www.nvidia.com/en-us/glossary/xgboost/. Accessed 3 Dec. 2024.

# A   Appendix A: Code

```python
# Converting Unix timestamp to Datetime, truncated to nearest hour
rides_df = rides_df.copy()
rides_df['time_stamp'] = pd.to_datetime(rides_df['time_stamp'],
    unit='ms').dt.floor('30min')
weather_df = weather_df.copy()
weather_df['time_stamp'] = pd.to_datetime(rides_df['time_stamp'],
    unit='ms').dt.round('30min')
```
Listing 1: Sample Data Preprocessing

```python
from google.cloud import bigquery
from google.api_core.exceptions import NotFound

# Initialize BigQuery client
client = bigquery.Client.from_service_account_json("secrets/
    serviceKey.json")

# Define the dataset reference
project_id = "xxx"
dataset_id = "yyy"
dataset_ref = f"{project_id}.{dataset_id}"

# List all tables in the dataset
tablesList = []
try:
    tables = client.list_tables(dataset_ref)
    print(f"Tables in dataset {dataset_id}:")
    for table in tables:
        tablesList.append(table.table_id)
```

```
19        print(table.table_id)
20 except NotFound:
21     print(f"Dataset {dataset_id} does not exist.")
```

Listing 2: Google Cloud BigQuery: Table creation

```
1 def overWriteDataToTable(dataFrameName, table_ref):
2     job_config = bigquery.LoadJobConfig(
3         write_disposition="WRITE_TRUNCATE",  # Overwrite the data
    to the table
4     )
5     job = client.load_table_from_dataframe(dataFrameName, table_ref
    , job_config=job_config)
6
7     # Wait for the job to complete
8     job.result()
9
10    print(f"Over written {len(dataFrameName)} rows to the table {
    table_id} in dataset {dataset_id}.")
```

Listing 3: Google Cloud BigQuery: Data Push

```
1 from pyspark.sql import SparkSession
2 import os
3 os.environ['GOOGLE_APPLICATION_CREDENTIALS'] = 'secrets/serviceKey.
    json'
4
5
6 # JAR paths for BigQuery and GCS connectors
7 bigquery_connector_jar = "spark-bigquery-connector.jar"
8 gcs_connector_jar = "gcs-connector.jar"
9
10
11 # Create SparkSession with both connectors
12 spark = SparkSession.builder \
13     .appName("PySpark with BigQuery and GCS") \
14     .config("spark.jars", f"{bigquery_connector_jar},{
    gcs_connector_jar}") \
15     .config("spark.sql.catalog.spark_bigquery", "com.google.cloud.
    spark.bigquery.BigQueryCatalog") \
16     .config("spark.hadoop.google.cloud.auth.service.account.json.
    keyfile", "secrets/serviceKey.json") \
17     .config("spark.bigquery.projectId", "idmpproject-441123") \
18     .getOrCreate()
19
20 spark
```

Listing 4: Google Cloud BigQuery: Data Pull using PySpark

```
1 # Visualization 1 Fare price distribution (histogram and bar plot)
    - Uber
2
3 from pyspark.sql.functions import col, lit
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 # Specify the output file path
8 output_image_path = os.path.join(output_folder, "
    uber_fare_price_distribution.png")
```

```python
 9
10  # Convert the 'price' column to Pandas for visualization
11  price_data_pd = uber_df.select("price").filter(col("price").
       isNotNull()).toPandas()
12
13  # Create a figure with two subplots: one for the histogram and one
        for the boxplot
14  fig, axes = plt.subplots(1, 2, figsize=(15, 6))
15
16  # Histogram Plot
17  sns.histplot(price_data_pd['price'], bins=30, kde=True, color='blue
       ', ax=axes[0])
18  axes[0].set_title('Distribution of Uber Fare Prices')
19  axes[0].set_xlabel('Price')
20  axes[0].set_ylabel('Frequency')
21
22  # Boxplot Plot
23  sns.boxplot(x=price_data_pd['price'], color='cyan', ax=axes[1])
24  axes[1].set_title('Boxplot of Uber Fare Prices')
25  axes[1].set_xlabel('Price')
26
27  # Display the plots
28  plt.tight_layout()
29
30  # Save the figure as a PNG image
31  plt.savefig(output_image_path)
32
33  # Inform the user
34  print(f"Visualization saved to '{output_image_path}'")
35
36  plt.show()
```

Listing 5: Sample Visualization using PySpark

```python
 1  import os
 2  import pandas as pd
 3  from flask import Flask, render_template, request
 4  from google.cloud import bigquery
 5  from datetime import datetime
 6
 7
 8  # Dynamically determine the absolute path to the 'static' folder
 9  BASE_DIR = os.path.abspath(os.path.dirname(__file__))
10  STATIC_FOLDER = os.path.join(BASE_DIR, 'static')
11  os.environ['GOOGLE_APPLICATION_CREDENTIALS'] = 'secrets/serviceKey.
       json'
12
13  client = bigquery.Client()
14  project_id = "xxx"
15  dataset_id = "yyy"
16
17
18  app = Flask(__name__)
19
20  @app.route('/')
21  def home():
22      return render_template('index.html')
```

Listing 6: Flask initialization

```python
import plotly.express as px
import pandas as pd
import plotly.graph_objects as go

@app.route('/dynamic-visualizations', methods=['GET', 'POST'])
def dynamic_visualization():
    # Define project_id and dataset_id variables
    project_id = "idmpproject-441123"
    dataset_id = "uberFareEstimation"

    # Define the BigQuery queries to get unique sources and
    destinations from both tables
    source_query = f"""
    SELECT DISTINCT source FROM `{project_id}.{dataset_id}.
    uber_data`
    UNION DISTINCT
    SELECT DISTINCT source FROM `{project_id}.{dataset_id}.
    lyft_data`
    """

    destination_query = f"""
    SELECT DISTINCT destination FROM `{project_id}.{dataset_id}.
    uber_data`
    UNION DISTINCT
    SELECT DISTINCT destination FROM `{project_id}.{dataset_id}.
    lyft_data`
    """

    # Execute the queries
    sources = client.query(source_query).result()
    destinations = client.query(destination_query).result()

    source_list = [row.source for row in sources]
    destination_list = [row.destination for row in destinations]

    # Create a dictionary mapping sources to possible destinations
    source_to_dest = {}
    for source in source_list:
        query = f"""
        SELECT DISTINCT destination FROM `{project_id}.{dataset_id
}.uber_data`
        WHERE source = '{source}'
        UNION DISTINCT
        SELECT DISTINCT destination FROM `{project_id}.{dataset_id
}.lyft_data`
        WHERE source = '{source}'
        """
        dest_result = client.query(query).result()
        source_to_dest[source] = [row.destination for row in
dest_result]

    # Initialize variables for source and destination
    source = None
    destination = None
    uber_avg_price_per_hour = None
    lyft_avg_price_per_hour = None
    df_html = None  # To store the HTML table
```

```python
51      if request.method == 'POST':
52          source = request.form.get('source')
53          destination = request.form.get('destination')
54
55          # BigQuery query to get the data based on selected source
    and destination from both tables
56          query = f"""
57          WITH uber_data AS (
58              SELECT
59                  EXTRACT(HOUR FROM time_stamp) AS hour,
60                  ROUND(AVG(price), 2) AS avg_price
61              FROM `{project_id}.{dataset_id}.uber_data`
62              WHERE source = '{source}' AND destination = '{
    destination}'
63              GROUP BY hour
64          ),
65          lyft_data AS (
66              SELECT
67                  EXTRACT(HOUR FROM time_stamp) AS hour,
68                  ROUND(AVG(price), 2) AS avg_price
69              FROM `{project_id}.{dataset_id}.lyft_data`
70              WHERE source = '{source}' AND destination = '{
    destination}'
71              GROUP BY hour
72          )
73          SELECT
74              COALESCE(u.hour, l.hour) AS hour,
75              u.avg_price AS uber_avg_price,
76              l.avg_price AS lyft_avg_price
77          FROM uber_data u
78          FULL OUTER JOIN lyft_data l ON u.hour = l.hour
79          ORDER BY hour
80          """
81
82          # Execute the query and fetch the results
83          result = client.query(query).result()
84
85          # Convert the result to a DataFrame
86          df = pd.DataFrame([dict(row) for row in result])
87
88          if df.empty:
89              uber_avg_price_per_hour = "No data available for the
    selected source and destination."
90              lyft_avg_price_per_hour = "No data available for the
    selected source and destination."
91          else:
92              # Plot the line graph
93              fig = go.Figure()
94
95              fig.add_trace(go.Scatter(x=df['hour'], y=df['
    uber_avg_price'], mode='lines', name='Uber'))
96              fig.add_trace(go.Scatter(x=df['hour'], y=df['
    lyft_avg_price'], mode='lines', name='Lyft'))
97
98              # Set the range for the x-axis and y-axis
99              fig.update_layout(
100                 title="Uber vs Lyft Average Price Per Hour",
101                 xaxis_title="Hour",
```

```
102                    yaxis_title="Average Price",
103                    xaxis=dict(
104                        range=[df['hour'].min(), df['hour'].max()]  #
       Limit x-axis range to the hours in the data
105                    ),
106                    yaxis=dict(
107                        range=[df[['uber_avg_price', 'lyft_avg_price'
       ]].min().min(), df[['uber_avg_price', 'lyft_avg_price']].max().
       max()]  # Adjust y-axis range dynamically
108                    )
109                )
110
111            # Render the plot in the HTML template
112            graph_html = fig.to_html(full_html=False)
113
114            # Convert the DataFrame to an HTML table
115            df_html = df.to_html(classes='data', header=True, index
       =False)
116            uber_avg_price_per_hour = graph_html  # Ensure this
       contains the full HTML of the graph
117
118         return render_template('dynamic_visualizations.html',
119                            sources=source_list,
120                            destinations=destination_list,
121                            source_to_dest=source_to_dest,
122                            uber_avg_price_per_hour=
       uber_avg_price_per_hour,
123                            df_html=df_html,
124                            selected_source=source,
125                            selected_destination=destination)
```

Listing 7: Flask dynamic visualization and data retrieval using BigQuery and SQL

```
1  # Define categorical and numerical features
2  categorical_features = ['cab_type', 'source', 'destination']
3  numerical_features = ['distance', 'surge_multiplier', 'hour', '
      day_of_week']
4
5  # Set up the ColumnTransformer
6  preprocessor = ColumnTransformer(
7      transformers=[
8          ('cat', OneHotEncoder(handle_unknown='ignore'),
      categorical_features),
9          ('num', 'passthrough', numerical_features)
10     ]
11 )
12
13 # Display the preprocessor to confirm setup
14 print(preprocessor)
```

Listing 8: ML Modeling: Configures a ColumnTransformer to preprocess categorical features using OneHotEncoder

```
1  # Cross-validation and CV metrics for the RF model using MAE
2
3  # Sample 25% of the data for cross-validation
4  X_sample = X.sample(frac=0.25, random_state=42)
5  y_sample = y.loc[X_sample.index]
6
```

```
 7  kf = KFold(n_splits=5, shuffle=True, random_state=42)
 8
 9  # Define the scoring metric (negative MAE)
10  scorer = make_scorer(mean_absolute_error, greater_is_better=False)
11
12  # Perform cross-validation using the pipeline
13  cv_scores = cross_val_score(pipeline, X_sample, y_sample, scoring=
        scorer, cv=kf, n_jobs=-1)
14
15  # Convert negative MAE to positive MAE
16  cv_mae_scores = -cv_scores
17
18  # Print the results
19  print(f"Cross-Validation MAE Scores: {cv_mae_scores}")
20  print(f"Mean MAE: {np.mean(cv_mae_scores)}")
21  print(f"Standard Deviation of MAE: {np.std(cv_mae_scores)}")
```

Listing 9: Cross Validating RF Model and Getting Metrics for RF Model

```
 1  # Import the XGBoost library
 2  from xgboost import XGBRegressor
 3
 4  # Define the XGBoost model
 5  xgb_model = XGBRegressor(n_estimators=150, learning_rate=0.1,
        max_depth=6, random_state=42, enable_categorical=True)
 6
 7  # Create a pipeline with preprocessing and the XGBoost model
 8  pipeline_xgb = Pipeline(steps=[
 9      ('preprocessor', preprocessor),  # Use the same preprocessor as
        Random Forest
10      ('regressor', xgb_model)
11  ])
12
13  # Train the pipeline on the training data
14  pipeline_xgb.fit(X_train, y_train)
15
16  print("XGBoost pipeline training complete.")
17
18  # Predict on the test set
19  y_pred_xgb = pipeline_xgb.predict(X_test)
20
21
22  # Evaluate performance
23  mae_xgb = mean_absolute_error(y_test, y_pred_xgb)
24  mse_xgb = mean_squared_error(y_test, y_pred_xgb)
25  rmse_xgb = np.sqrt(mse_xgb)
26  r2_xgb = r2_score(y_test, y_pred_xgb)
27
28  # Print the metrics
29  print(f"XGBoost Evaluation Metrics:")
30  print(f" - Mean Absolute Error (MAE): {mae_xgb}")
31  print(f" - Mean Squared Error (MSE): {mse_xgb}")
32  print(f" - Root Mean Squared Error (RMSE): {rmse_xgb}")
33  print(f" - R  Score: {r2_xgb}"
```

Listing 10: XGBoost Model Code

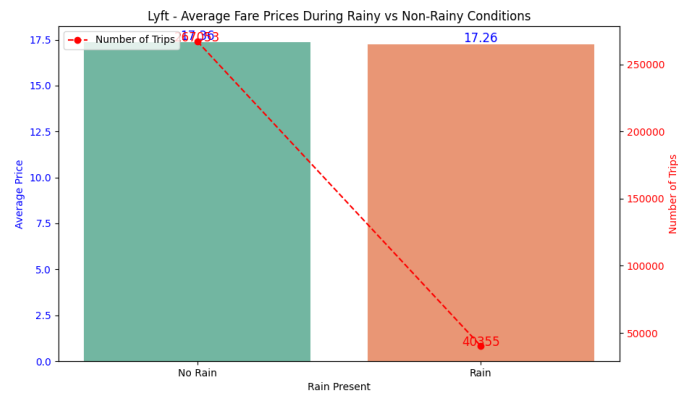# B    Appendix B: Additional Figures



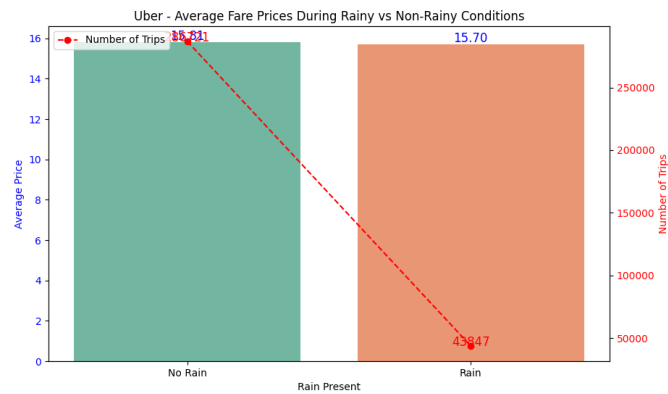Figure 16: Line Graph on Lyft's pricing on rainy and non-rainy days



Figure 17: Line Graph on Uber's pricing on rainy and non-rainy days

Figure 18: Lyft's Place x Hourly Pricing (Line heat map)



Figure 19: Uber's Place x Hourly Pricing (Line heat map)