# UBER VS LYFT

Comparison Between Rivals in Boston, MA

# WHERE DO I FIND THE DATA SOURCE ?

Find the data source from Kaggle, and they had a great source that predicts the cab prices for Uber vs. Lyft in Boston, MA.

```python
In [3]: # Clean up the dataframe
        # Dropping all NaN values
        # Renaming the columns
        # Converting time_stamp to date format

        clean_prices = cab_prices_df.dropna(how='any')
        clean_prices_df = pd.DataFrame(clean_prices)
        clean_prices_df.head()

        clean_prices_rename = clean_prices_df.rename(columns={"distance": "Distance", "cab_type": "Type of Cab",
                                                              "time_stamp":"Date", "source":"Source",
                                                              "price":"Price", "surge_multiplier":"Surge Multiplier",
                                                              "id":"User ID", "product_id":"Service Type", "name":"Vehicle Type"})

        clean_prices_rename_df = pd.DataFrame(clean_prices_rename)

        clean_prices_rename_df['Date'] = pd.to_datetime(clean_prices_rename_df['Date']/1000, unit='s')
        clean_prices_rename_df.head()
```

Out[3]:

| | Distance | Type of Cab | Date | destination | Source | Price | Surge Multiplier | User ID | Service Type | Vehicle Type |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.44 | Lyft | 2018-12-16 09:30:07.890000105 | North Station | Haymarket Square | 5.0 | 1.0 | 424553bb-7174-41ea-aeb4-fe06d4f4b9d7 | lyft_line | Shared |
| 1 | 0.44 | Lyft | 2018-11-27 02:00:23.677000046 | North Station | Haymarket Square | 11.0 | 1.0 | 4bd23055-6827-41c6-b23b-3c491f24e74d | lyft_premier | Lux |
| 2 | 0.44 | Lyft | 2018-11-28 01:00:22.197999954 | North Station | Haymarket Square | 7.0 | 1.0 | 981a3613-77af-4620-a42a-0c0866077d1e | lyft | Lyft |
| 3 | 0.44 | Lyft | 2018-11-30 04:53:02.749000072 | North Station | Haymarket Square | 26.0 | 1.0 | c2d88af2-d278-4bfd-a8d0-29ca77cc5512 | lyft_luxsuv | Lux Black XL |
| 4 | 0.44 | Lyft | 2018-11-29 03:49:20.223000050 | North Station | Haymarket Square | 9.0 | 1.0 | e0126e1f-8ca9-4f2e-82b3-50505a09db9a | lyft_plus | Lyft XL |

```python
In [4]: # Clean up the dataframe
        # Dropping all NaN values
        # Renaming the columns
        # Converting time_stamp to date format

        clean_weather = weather_df.dropna(how='any')
        clean_weather_df = pd.DataFrame(clean_weather)

        clean_weather_rename = clean_weather_df.rename(columns={'temp':'Temperature in Fahrenheit', 'location':'Location',
                                                                'clouds':'Cloudiness', 'pressure':'Pressure', 'rain':'Rainfall in inches',
                                                                'time_stamp':'Date', 'humidity':'Humidity', 'wind':'Wind Speed'})

        clean_weather_rename_df = pd.DataFrame(clean_weather_rename)

        clean_weather_rename_df['Date'] = pd.to_datetime(clean_weather_rename_df['Date'], unit='s')
        clean_weather_rename_df.head()
```

Out[4]:

| | Temperature in Fahrenheit | Location | Cloudiness | Pressure | Rainfall in inches | Date | Humidity | Wind Speed |
|---|---|---|---|---|---|---|---|---|
| 0 | 42.42 | Back Bay | 1.0 | 1012.14 | 0.1228 | 2018-12-16 23:45:01 | 0.77 | 11.25 |
| 1 | 42.43 | Beacon Hill | 1.0 | 1012.15 | 0.1846 | 2018-12-16 23:45:01 | 0.76 | 11.32 |
| 2 | 42.50 | Boston University | 1.0 | 1012.15 | 0.1089 | 2018-12-16 23:45:01 | 0.76 | 11.07 |
| 3 | 42.11 | Fenway | 1.0 | 1012.13 | 0.0969 | 2018-12-16 23:45:01 | 0.77 | 11.09 |
| 4 | 43.13 | Financial District | 1.0 | 1012.14 | 0.1786 | 2018-12-16 23:45:01 | 0.75 | 11.49 |

# CLEANING UP THE DATA

MERGING DATA TOGETHER

# Importing and Finding Objects…

- Importing the dependencies and files

- Finding necessary services for Uber and Lyft

- Finding the total values, after merging the data together

```
In [6]:  # bar chart for comparing uber and lyft rides
         labels_x = ['Uber','Lyft']
         counts_y = [clean_prices_df.cab_type[(clean_prices_df.cab_type) == 'Ube
                     clean_prices_df.cab_type[(clean_prices_df.cab_type)=='Lyft

         plt.bar(labels_x, counts_y, color="lightgreen", align="center", width

         plt.xlabel("Uber vs. Lyft")
         plt.ylabel("Number Of Rides")
         plt.savefig("uber_vs_lyft_number_of_rides.png")
         plt.show()
```



```
In [5]:  # Uber vs Lyft Usage Comparison in pie chart
         cab_counts = [clean_prices_df.cab_type[(clean_prices_df.cab_type) == 'Lyft']
                       clean_prices_df.cab_type[(clean_prices_df.cab_type) == 'Uber']

         explode = (0.08, 0)

         cab_types = ['Uber','Lyft']

         colors = ["lightcoral", "lightskyblue"]

         plt.pie(cab_counts, explode=explode, labels=cab_types, colors=colors,
                 autopct="%1.1f%%", shadow=False, startangle=100)

         plt.axis("equal")
         plt.savefig("uber_vs_lyft_usage_comparison.png")
         plt.show()
```



# Questions: Is there more Uber users than Lyft users?

- Yes, there are more rides with Uber than Lyft rides.

- However, I find out that Lyft has silently more users than Uber users in percentages in the pie chart.

```
In [8]: # Uber vs Lyft Price Comparison bar chart

average_price = clean_prices_df.groupby('name')['price'].mean()

bar_chart_price = average_price.plot.bar(x="name", y="price",
                                          color="pink", figsize=(10,8), fontsize= 12)

bar_chart_price

plt.title("Uber & Lyft Services", fontsize = 18)
plt.xlabel("Type Of Services", fontsize = 14)
plt.ylabel("Average Prices in Dollars", fontsize = 14)
plt.tight_layout()
plt.savefig("uber_vs_lyft_avg_price_comparison_by_service_type.png")
plt.show()
```

**Lyft Rides**

Uber & Lyft Services

Questions: How many average prices are there in each Uber and Lyft services?

■ Lyft Shared has approx. 5-6 dollars in average prices.

■ Lyft Lux Black XL has slightly more than 30 dollars in average prices.

■ Lyft Shared has a lowest average prices in dollars, while Lyft Lux Black XL has the highest average prices in dollars.

■ Lyft Rides are in the red box.

■ Uber Rides are outside of the box.

```
In [71]:  # Average Price Vs Distance traveled

          # Collect Uber and Lyft in the data
          uber_df = merged_df[merged_df['Type of Cab'] == 'Uber']
          lyft_df = merged_df[merged_df['Type of Cab'] == 'Lyft']

          # Find the averages in Uber and Lyft prices
          uber_avgprice = uber_df.groupby('Distance')['Price'].mean()
          lyft_avgprice = lyft_df.groupby('Distance')['Price'].mean()

          # Plot the charts and apply some styling
          fig1, ax1 = plt.subplots(figsize=(10,8))

          plt.plot(uber_avgprice, label='Uber')
          plt.plot(lyft_avgprice, label='Lyft')

          plt.title('Average Price in Dollars VS Distance Traveled', fontsize=16)
          plt.xlabel('Distance Traveled in Miles', fontsize=16)
          plt.ylabel('Average Price in Dollars', fontsize=16)
          plt.legend()
          plt.savefig('Average_Price_vs_Distance_Traveled')
          plt.show()
```

- Lyft has the highest amount of average price than Uber as the distance increased.

- Both average prices are increasing as the distance traveled increases.
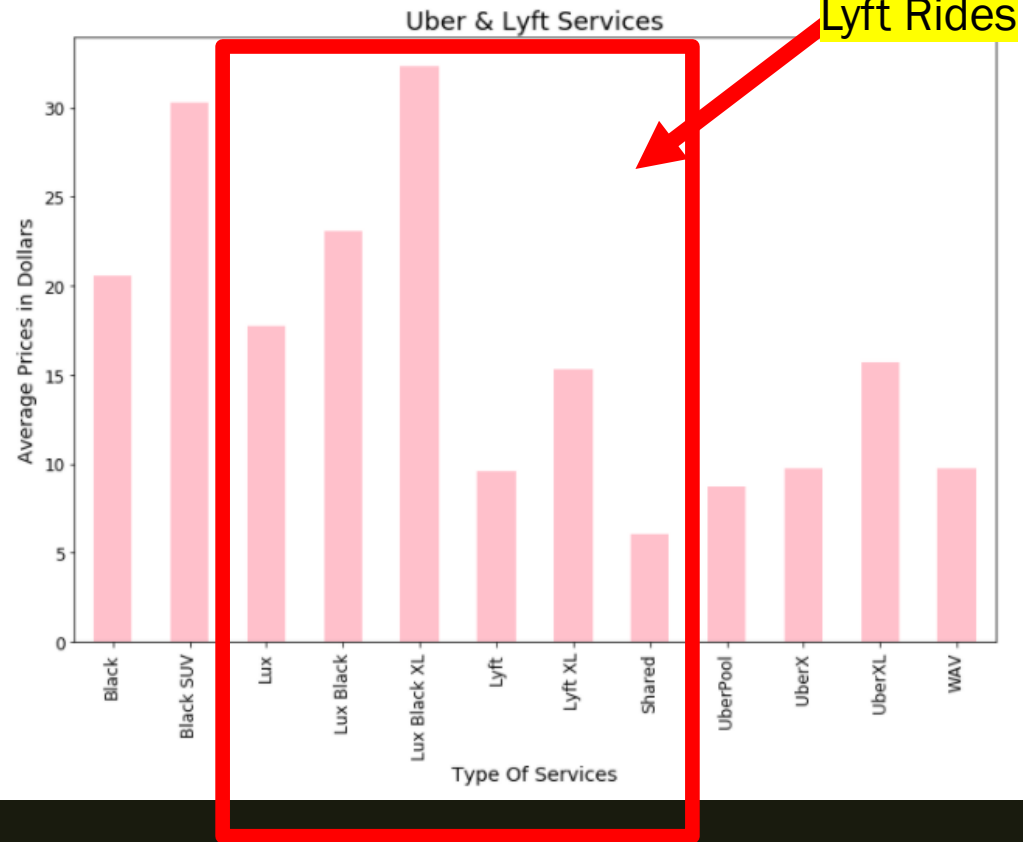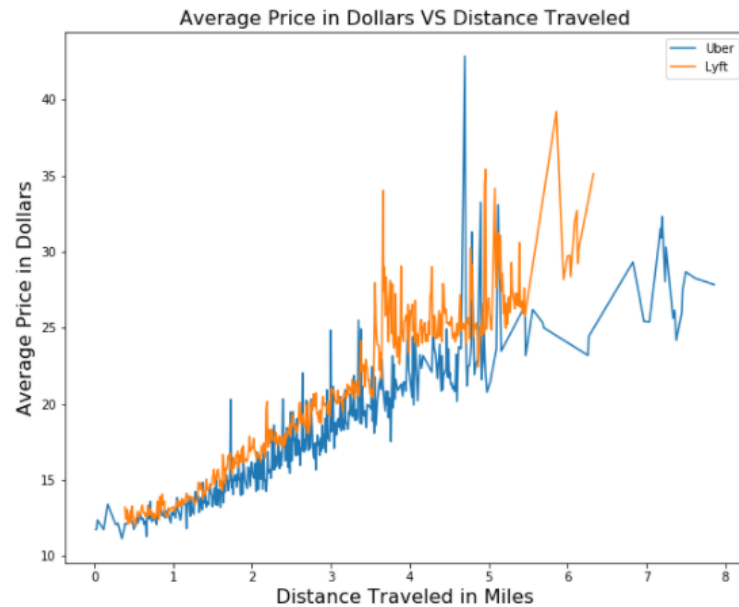
# Summary Analysis

- **Is there more Uber users than Lyft users?**
  - *Yes, there are more rides with Uber than Lyft rides.*
  - *However, I find out that Lyft has silently more users than Uber users in percentages in the pie chart.*
- **How many average prices are there in each Uber and Lyft services?**
  - *Lyft Shared has approx. 5-6 dollars in average prices.*
  - *Lyft Lux Black XL has slightly more than 30 dollars in average prices.*
  - *Lyft Shared has a lowest average prices in dollars, while Lyft Lux Black XL has the highest average prices in dollars.*
  - *Lyft Rides are in the red box.*
  - *Uber Rides are outside of the box.*
- **Lyft has the highest amount of average price than Uber as the distance increased.**
- **Both average prices are increasing as the distance traveled increases.**