

# Problem Statement

As some point or the other, almost each one of us has used an Ola or Uber for taking a ride.

Ride-hailing services are services that use online-enabled platforms to connect between passengers and local drivers using their personal vehicles. In most cases they are a comfortable method for door-to-door transport. However they are cheaper than using licensed taxicabs. Examples of ride-hailing services include Uber and Lyft.

To improve the efficiency of taxi dispatching systems for such services, it is important to be able to predict how long a driver will have his taxi occupied. If a dispatcher knew approximately when a taxi driver would be ending their current ride, they would be better able to identify which driver to assign to each pickup request.

In this competition, we are challenged to build a model that predicts the total ride duration of taxi trips in New York City.

## 1. Exploratory Data Analysis

Let's check the data file! According to the data description we should find the following columns:

- id - a unique identifier for each trip
- vendor\_id - a code indicating the provider associated with the trip record
- pickup\_datetime - date and time when the meter was engaged
- dropoff\_datetime - date and time when the meter was disengaged
- passenger\_count - the number of passengers in the vehicle (driver entered value)
- pickup\_longitude - the longitude where the meter was engaged
- pickup\_latitude - the latitude where the meter was engaged
- dropoff\_longitude - the longitude where the meter was disengaged
- dropoff\_latitude - the latitude where the meter was disengaged
- store\_and\_fwd\_flag - This flag indicates whether the trip record was held in vehicle memory before sending to the vendor because the vehicle did not have a connection to the server (Y=store and forward; N=not a store and forward trip)
- trip\_duration - (target) duration of the trip in seconds

Here, we have 2 variables dropoff\_datetime and store\_and\_fwd\_flag which are not available before the trip starts and hence will not be used as features to the model.

### 1.1 Load Libraries

```
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from datetime import datetime
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
import warnings
warnings.filterwarnings('ignore')
```

### Load Data

```
df = pd.read_csv('nyc_taxi_trip_duration.csv')
```

### File structure and content

```
print('We have {} rows.'.format(df.shape[0]))
print('We have {} columns.'.format(df.shape[1]))
df.info()
```

We have 729322 rows.  
We have 11 columns.

```
Out[52]:
```

	id	vendor_id	pickup_datetime	dropoff_datetime	passenger_count	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	store_and_fwd_flag	trip_duration
0	16988885	1	2016-03-11 23:35:37	2016-03-11 23:53:57	2	-73.9883	40.7357	-73.9848	40.6469	N	1199

At first glance, we can see the types of each variable and what they look like.

### Missing Values

Knowing about missing values is important because they indicate how much we don't know about our data. Making inferences based on just a few cases is often unwise. In addition, many modelling procedures break down when missing values are involved and the corresponding rows will either have to be removed completely or the values need to be estimated somehow.

```
Out[53]:
```

np.sum(pd.isnull(df))

```
Out[52]:
```

	id	vendor_id	pickup_datetime	dropoff_datetime	passenger_count	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	store_and_fwd_flag	trip_duration
0	16988885	1	2016-03-11 23:35:37	2016-03-11 23:53:57	2	-73.9883	40.7357	-73.9848	40.6469	N	1199

Fortunately, in this dataset we do not have any missing values which is great.

### Reformatting features & Checking consistency

There are a variety of features within the dataset and it is important to convert them into the right format such that we can analyse them easily. This would include converting datetime features and string features.

Also, one important thing is never to take assumptions without backing it with data. Here, as you can see the trip duration can also be calculated pick up and drop off datetime. We will check whether the given duration is consistent with the calculated trip duration

```
Out[54]:
```

```
df['pickup_datetime'] = pd.to_datetime(df['pickup_datetime'])
df['dropoff_datetime'] = pd.to_datetime(df['dropoff_datetime'])

# Converting yes/no flag to 1 and 0
df['store_and_fwd_flag'] = 1 * (df['store_and_fwd_flag'].values == 'Y')

df['check_trip_duration'] = (df['dropoff_datetime'] - df['pickup_datetime']).map(lambda x: x.total_seconds())

duration_difference = df[np.abs(df['check_trip_duration'].values - df['trip_duration'].values) > 1]
duration_difference.shape
```

(6, 12)

This implies that there is no inconsistency in data w.r.t the drop location and trip duration

### Target Exploration

In this section we will take a look at the trip duration which is the target variable. It is crucial to understand it in detail as this is what we are trying to predict accurately.

```
Out[55]:
```

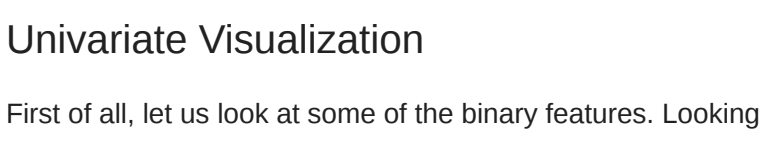
```
df['trip_duration'].describe()/3600 # Trip duration in hours
```

	count	292	589444
mean	8.54640	1.97597	1.0
std	8.110278	0.110278	0
25%	0.110278	0.110278	0
50%	0.110278	0.110278	0
75%	0.110278	0.110278	0
max	0.110278	0.110278	0

Wow! There is a trip with duration of 599 hours. This is a huge outlier and might create problems in the prediction stage. One idea is to log transform the trip duration before prediction to visualise it better.

```
Out[56]:
```

```
df['log_trip_duration'] = np.log(df['trip_duration'] + 1)
sns.distplot(df['log_trip_duration'], kde=False, bins=200)
plt.show()
```



We find:

- The majority of rides follow a rather smooth distribution that looks almost log-normal with a peak just around exp(6.5) i.e. about 17 minutes.
- There are several suspiciously short rides with less than 10 seconds duration.
- As discussed earlier, there are a few huge outliers near 12.

### Univariate Visualization

First of all, let us look at some of the binary features. Looking at each feature might uncover some insight that might be useful at later modelling stages

```
Out[57]:
```

```
# Binary Features
plt.figure(figsize=(22, 6))
fig, axs = plt.subplots(ncols=2)

# Passenger Count
sns.countplot(df['passenger_count'])
plt.xlabel('Passenger Count')
plt.ylabel('Frequency')
```

```
Out[58]:
```

```
# vendor_id
sns.countplot(df['vendor_id'])
plt.xlabel('Vendor ID')
plt.ylabel('Frequency')
```

```
Out[59]:
```

```
# store_and_fwd_flag
sns.countplot(df['store_and_fwd_flag'])
plt.xlabel('store_and_fwd_flag')
plt.ylabel('Frequency')
```

```
Out[60]:
```

```
Test(0.5, 'Frequency')
```



Observations:

- Most of the trips involve only 1 passenger. There are trips with 7-9 passengers but they are very low in number.
- Vendor 2 has more number of trips as compared to vendor 1.
- The store\_and\_fwd\_flag values, indicating whether the trip data was sent immediately to the server (Y) or held in the memory of the taxi because there was no connection to the server (N), show that there was almost no storing taking place.

Now, we will delve into the datetime features to understand the trend of number of hourly/monthly/daily taxi trips

```
Out[61]:
```

```
df['pickup_datetime'].min(), df['pickup_datetime'].max()
```

(Timestamp('2016-01-01 00:01:14'), Timestamp('2016-06-30 23:59:37'))

Clearly, These trips are for first 6 months of 2016. To look at trends, we first need to extract week days and hour of day from the pickup date.

```
Out[62]:
```

```
df['day_of_week'] = df['pickup_datetime'].dt.weekday
df['hour_of_day'] = df['pickup_datetime'].dt.hour
```

```
Out[63]:
```

```
# Datetime Features
plt.figure(figsize=(22, 6))

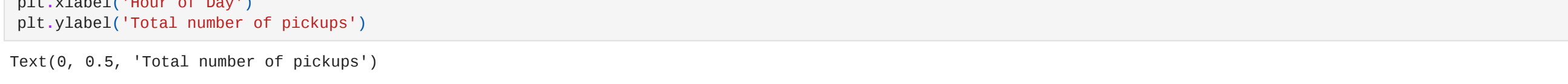
# Passenger Count
sns.countplot(df['day_of_week'])
plt.xlabel('Day of Week')
plt.ylabel('Total Number of pickups')
```

```
Out[64]:
```

```
# vendor_id
sns.countplot(df['hour_of_day'])
plt.xlabel('Hour of Day')
plt.ylabel('Total number of pickups')
```

```
Out[65]:
```

```
Test(0.5, 'Frequency')
```



Observations:

- Number of pickups for weekends is much lower than week days with a peak on Thursday (4). Note that here weekday is a decimal number, where 0 is Sunday and 6 is Saturday.
- Number of pickups as expected is highest in late evenings. However, it is much lower during the morning peak hours.

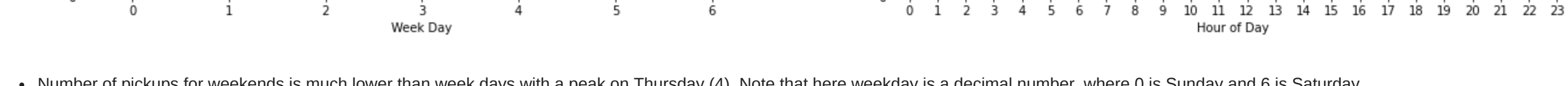
### Latitude & Longitude

Let's look at the geospatial or location features to check consistency. They should not vary much as we are only considering trips within New York city.

```
Out[66]:
```

```
sns.set(style='white', palette='muted', color_codes=True)
f, axes = plt.subplots(2,2,figsize=(30, 30), sharex=False, sharey=False)

# pickup_latitude
sns.distplot(df['pickup_latitude'].values, label='pickup_latitude', color='b', bins=300, ax=axes[0,0])
sns.distplot(df['dropoff_latitude'].values, label='dropoff_latitude', color='r', bins=300, ax=axes[0,1])
sns.distplot(df['pickup_longitude'].values, label='pickup_longitude', color='b', bins=300, ax=axes[1,0])
sns.distplot(df['dropoff_longitude'].values, label='dropoff_longitude', color='r', bins=300, ax=axes[1,1])
plt.tight_layout()
plt.show()
```



Findings - (Here, red represents pickup and dropoff Longitudes & blue represents pickup & dropoff latitudes)

- From the plot above it is clear that pick and drop latitudes are centered around 40 to 41, and longitude are situated around -74 to -73.
- Some extreme co-ordinates has squeezed the plot such that we see a spike here
- A good idea is to remove these outliers and look at the distribution more closely

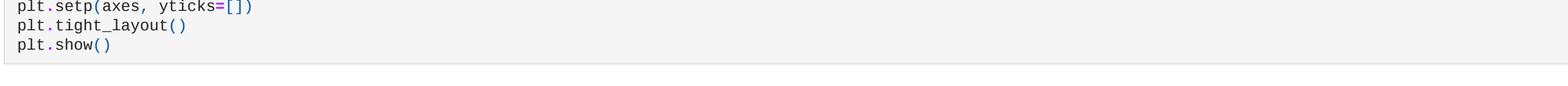
```
Out[67]:
```

```
df = df[df['pickup_latitude'] > 40.6 & (df['pickup_latitude'] < 40.9)]
df = df[df['dropoff_latitude'] > 40.6 & (df['dropoff_latitude'] < 40.9)]
df = df[df['pickup_longitude'] > -74.05 & (df['dropoff_longitude'] < -73.7)]
df = df[df['dropoff_longitude'] > -74.05 & (df['dropoff_longitude'] < -73.7)]
df_data_new = df.copy()
```

```
Out[68]:
```

```
sns.set(style='white', palette='muted', color_codes=True)
f, axes = plt.subplots(2,2,figsize=(10, 10), sharex=False, sharey=False)

# pickup_latitude
sns.distplot(df_data_new['pickup_latitude'].values, label='pickup_latitude', color='b', bins=300, ax=axes[0,0])
sns.distplot(df_data_new['dropoff_latitude'].values, label='dropoff_latitude', color='r', bins=300, ax=axes[0,1])
sns.distplot(df_data_new['pickup_longitude'].values, label='pickup_longitude', color='b', bins=180, ax=axes[1,0])
sns.distplot(df_data_new['dropoff_longitude'].values, label='dropoff_longitude', color='r', bins=180, ax=axes[1,1])
plt.tight_layout()
plt.show()
```



- We have a much better view of the distribution of coordinates instead of spikes. And we see that most trips are concentrated between these lat long only with a few significant clusters.
- These clusters are represented by the numerous peaks in the latitude and longitude histograms

### Bivariate Relations with Target

Now that we have gone through all the basic features one by one. Let us start looking at their relation with the target. This will help us in selecting and extracting features at the modelling stage.

Out[69]:

```
df[['id', 'vendor_id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'pickup_longitude', 'dropoff_longitude', 'store_and_fwd_flag',
    'trip_duration', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day']]
```

### vendor\_id vs Trip Duration

Let's check how the trip duration varies for different vendors.

```
Out[70]:
```

```
plt.figure(figsize=(22, 6))
sns.boxplot(x='vendor_id', y='trip_duration', data=df)
plt.show()
```



Wow! This did not came out as expected. The only thing I can see from this boxplot is that for vendor 2, there are a number of outliers exceeding 24 hours while vendor 1 does not have such long trips.

There could be 2 solutions to this:

- Remove the huge outliers and plot again
- Look at median trip duration for both vendors on hourly basis

Let's try the first technique now and check trips below 50000 seconds only

```
Out[71]:
```

```
df_sub = df[df['trip_duration'] < 50000]
sns.boxplot(x='vendor_id', y='trip_duration', data=df_sub)
plt.show()
```



As you can see, we were in a false perception earlier that vendor 1 had more outliers. Since the median is just around 600 seconds, we observe that vendor 2 has many more outliers as compared to vendor 1. Next, to confirm this, we will quickly look at the mean w.r.t day of week for both vendors using boxplot (time series plot) from seaborn.

### Trip Duration vs Passenger Count

Again as we are aware, there are a large number of outliers for trip duration and we will not be able to observe the differences. For this, we have taken a cutoff of 12000 seconds and used a boxplot.

```
Out[72]:
```

```
df['passenger_count'].value_counts()
```

```
Out[73]:
```

```
1    542476
2    164021
3    140211
4    23952
5    14021
6    13972
7    13972
8    13972
9    13972
Name: passenger_count, dtype: int64
```

```
Out[74]:
```

```
df['passenger_count'].value_counts()
sns.boxplot(x='passenger_count', y='trip_duration', data=df_sub)
plt.show()
```



- The boxplot clearly shows that there is not much of a difference in distribution for the most frequently occurring passenger count values - 1, 2, 3.
- Another key observation is that the number of outliers are reduced for higher passenger counts but that only comes down to the individual frequencies of each passenger count.

### Visualise most frequently occurring Pickup points on the latitude-longitude Map

Here, we try to visualise the most frequently occurring pickup points on the map and check how it is distributed spatially.

```
Out[75]:
```

```
np.zeros((3888, 3500, 3)), dtype=np.uint8)
rgb1, rgb2, rgb3 = 0, 0, 0
df_data_new['pick_lat_new'] = list(map(int, (df['pickup_latitude'] - (40.6888)) * 100000))
df_data_new['pick_lon_new'] = list(map(int, (df['pickup_longitude'] - (-74.859)) * 100000))
df_data_new['pick_lat_new'] = list(map(int, (df['pickup_latitude'] - (-74.859)) * 100000))
df_data_new['pick_lon_new'] = list(map(int, (df['pickup_longitude'] - (-74.859)) * 100000))
summary_plot = pd.DataFrame(df_data_new.groupby(['pick_lat_new', 'pick_lon_new'])['id'].count())
```

```
Out[76]:
```

```
summary_plot.reset_index(inplace=True)
summary_plot.head(100)
fig, lon_lat_list = plt.subplots(2,2,figsize=(14, 14))

lon_lat_list[0,0] = summary_plot[summary_plot['pick_lat_new'] > 0 & summary_plot['pick_lon_new'] > 0]
lon_lat_list[0,1] = summary_plot[summary_plot['pick_lat_new'] > 0 & summary_plot['pick_lon_new'] < 0]
lon_lat_list[1,0] = summary_plot[summary_plot['pick_lat_new'] < 0 & summary_plot['pick_lon_new'] > 0]
lon_lat_list[1,1] = summary_plot[summary_plot['pick_lat_new'] < 0 & summary_plot['pick_lon_new'] < 0]

for j in lon_lat_list:
    for i in lon_lat_list:
        if (i[0,0] > 0 & i[1,0] > 0):
            rgb1[i,0] = 255
            rgb1[i,1] = 0
            rgb1[i,2] = 0
        elif (i[0,0] > 0 & i[1,0] < 0):
            rgb1[i,0] = 0
            rgb1[i,1] = 255
            rgb1[i,2] = 0
        elif (i[0,0] < 0 & i[1,0] > 0):
            rgb1[i,0] = 0
            rgb1[i,1] = 0
            rgb1[i,2] = 255
        else:
            rgb1[i,0] = 255
            rgb1[i,1] = 255
            rgb1[i,2] = 255
fig, ax = plt.subplots(figsize=(14, 14))
ax.imshow(rgb1, cmap='hot')
ax.set_axis_off()
```



Findings - From the heatmap kind of image above -

- White points - 1-30 trips have white as pickup point
- Green points - 30-25 trips have green as pickup point
- Red points - More than 25 trips have red as pickup point

As expected there are a few small clusters for hot pickup points as displayed by red in the above plot. Most pickup points have less than 10 trips and distributed all over the city.

If you go and have a look at an actual map of New York City, red and green points are mostly concentrated around the Manhattan Area

### Correlation Heatmap

Let us quickly look at the correlation heatmap to check the correlations amongst all features.

```
Out[77]:
```

```
plt.figure(figsize=(12, 6))
df = df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day']]
corr = df.apply(lambda x: pd.factorize(x)[0]).corr()
ax = sns.heatmap(corr, yticklabels=corr.columns, xticklabels=corr.columns,
    linewidths=2, cmap='magma')
```



### Conclusions

- The majority of rides follow a rather smooth distribution that looks almost log-normal with a peak just around exp(6.5) i.e. about 17 minutes.
- There are several suspiciously short rides with less than 10 seconds duration.
- As discussed earlier, there are a few huge outliers near 12.
- Most of the trips involve only 1 passenger. There are trips with 7-9 passengers but they are very low in number.
- Vendor 2 has more number of trips as compared to vendor 1.
- Number of pickups for weekends is much lower than week days with a peak on Thursday (4). Note that here weekday is a decimal number, where 0 is Sunday and 6 is Saturday.
- Number of pickups as expected is highest in late evenings. However, it is much lower during the morning peak hours.
- We see that most trips are concentrated between these lat long only with a few significant clusters. These clusters are represented by the numerous peaks in the latitude and longitude histograms
- Vendor 2 has more number of trips as compared to vendor 1.
- These clusters are represented by the numerous peaks in the latitude and longitude histograms
- Median trip duration does not vary much as can be seen from the above plot for different vendors.
- The boxplot clearly shows that there is not much of a difference in distribution for the most frequently occurring passenger count values - 1, 2, 3.
- Another key observation is that the number of outliers are reduced for higher passenger counts but that only comes down to the individual frequencies of each passenger count.
- From the correlation heatmap we see that the latitude and longitude features have higher correlation with the target as compared to the other features.

```
Out[78]:
```

```
df = pd.read_csv('updated_data.csv')
```

### Calculating distance

```
Out[79]:
```

```
def find_dist(x1, y1, x2, y2):
    return ((x2-x1)**2 + (y2-y1)**2)**0.5
```

```
Out[80]:
```

```
df['distance_miles'] = df.apply(lambda x: find_dist(x['pickup_latitude'], x['dropoff_latitude'], x['pickup_longitude'], x['dropoff_longitude']), axis=1)
```

```
Out[81]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[82]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[83]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[84]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[85]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[86]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[87]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[88]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[89]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[90]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[91]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[92]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[93]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[94]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[95]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[96]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[97]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[98]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[99]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[100]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[101]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[102]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[103]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[104]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[105]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[106]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[107]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[108]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[109]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[110]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[111]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[112]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[113]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[114]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[115]:
```

```
df[['id', 'pickup_datetime', 'dropoff_datetime',
    'passenger_count', 'check_trip_duration', 'log_trip_duration',
    'day_of_week', 'hour_of_day', 'distance_miles']]
```

```
Out[116]:
```

```
df[['id', 'pickup_datetime',
```