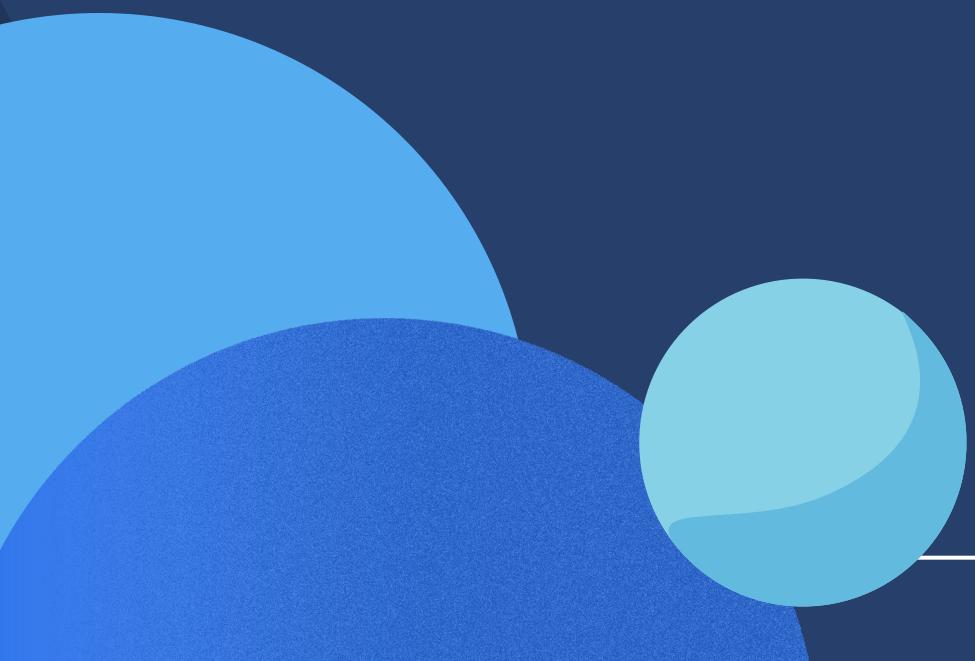


# Taxi Trip Duration prediction



# TABLE OF CONTENTS

- Objective
  - Data Gathering
  - Data Pre-processing
- Data Cleaning
- Missing value analysis
- Exploratory Data analysis
- Feature selection
- Train Model
  - Test model
  - Conclusion



# OBJECTIVE

**To build a Machine Learning model that predicts the total ride duration of taxi trips in New York City with minimum error.**

**To Business Goal : Improve the efficiency of taxi systems to be able to predict how long a driver will have his taxi occupied**

# Data Gathering

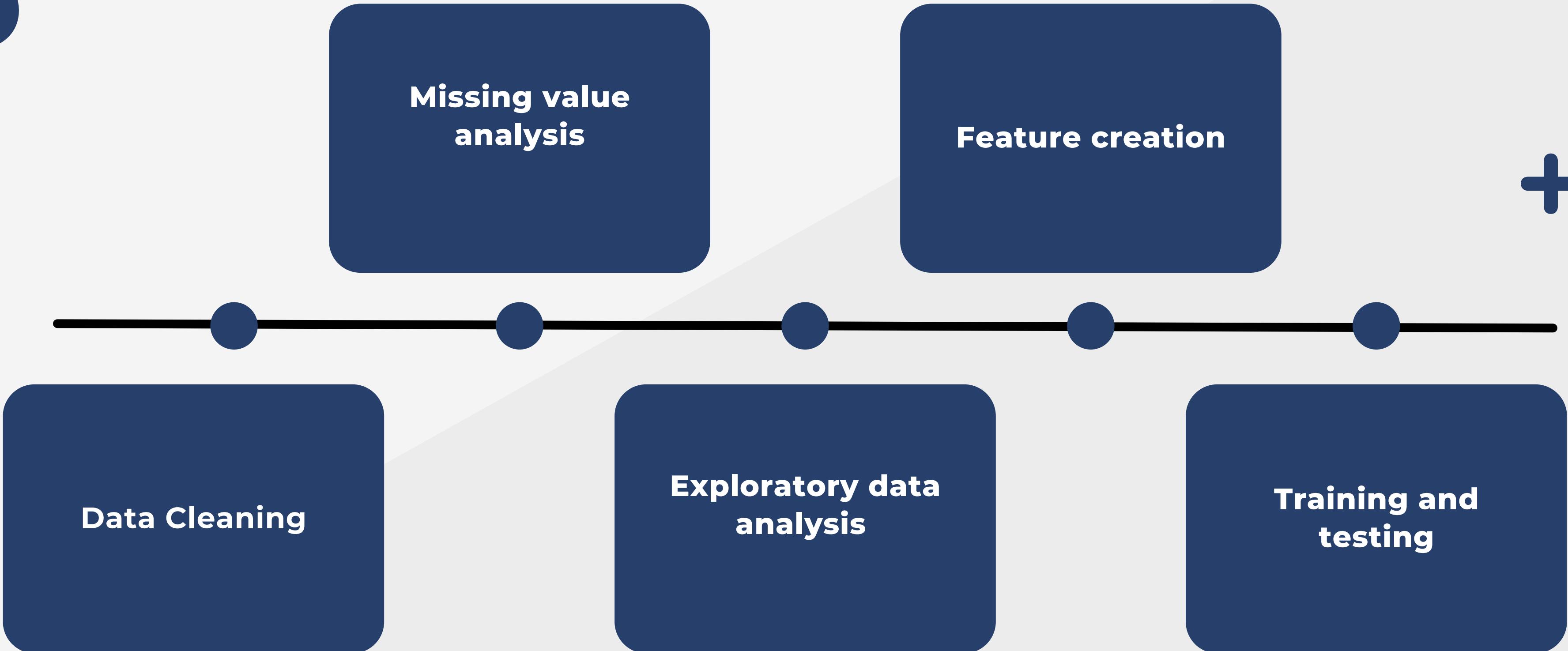
## Dataset: NYC\_taxi

The NYC taxi trip prediction dataset is a collection of data that captures various attributes and information about taxi trips in New York City. It provides insights into the duration and distance of the trips

The dataset includes

<b>id</b>	a unique identifier for each trip
<b>vendor_id</b>	a code indicating the provider associated with the trip record
<b>pickup_datetime</b>	date and time when the meter was engaged
<b>dropoff_datetime</b>	date and time when the meter was disengaged
<b>passenger_count</b>	the number of passengers in the vehicle (driver entered value)
<b>pickup_longitude</b>	the longitude where the meter was engaged
<b>pickup_latitude</b>	the latitude where the meter was engaged
<b>dropoff_longitude</b>	the longitude where the meter was disengaged
<b>dropoff_latitude</b>	the latitude where the meter was disengaged
<b>store_and_fwd_flag</b>	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
<b>trip_duration</b>	duration of the trip in seconds

# PROJECT TIMELINE



# Data preprocessing

## Data Cleaning

id, vendor\_id, pickup time and drop off time doesn't contribute to our target variable

They could be dropped from the training set

```
[3]: (1458644, 11)

▶ #Khushi 21BCE1282
df=df.drop(['id','vendor_id','pickup_datetime','dropoff_datetime'],axis=1)
df.dtypes

[5]: passenger_count      int64
pickup_longitude     float64
pickup_latitude      float64
dropoff_longitude    float64
dropoff_latitude     float64
store_and_fwd_flag   object
trip_duration        int64
dtype: object

[4]: #Khushi 21BCE1282
df.isnull().sum()

+ Code + Markdown
```

# Missing values

- Missing values can be identified by

```
df.isnull().sum()
```

The screenshot shows a Jupyter Notebook cell with the following code and output:

```
#Khushi 21BCE1282
df.isnull().sum()
```

[6]:

	passenger_count	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	store_and_fwd_flag	trip_duration
dtype:	int64						

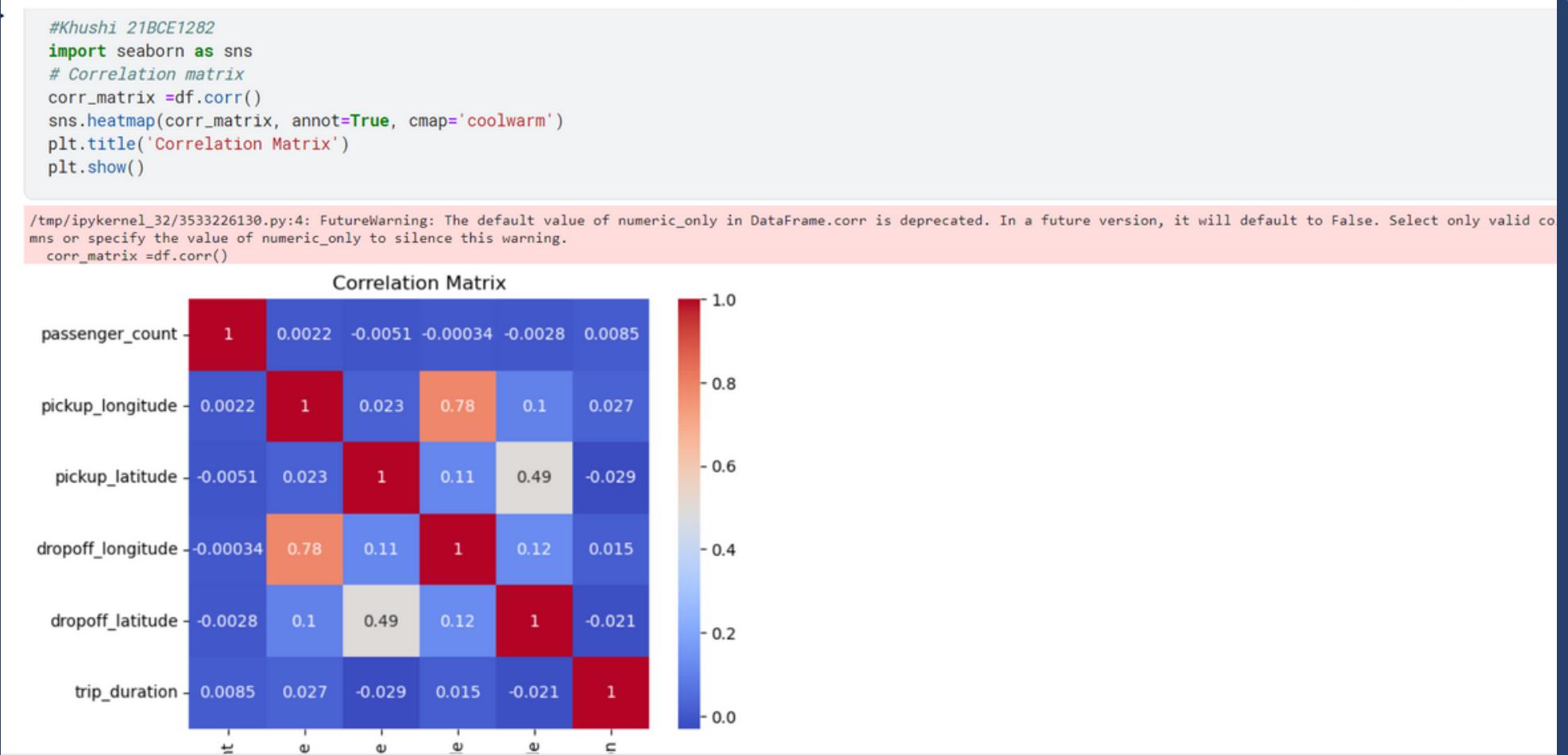
+ Code + Markdown

The output table shows that all columns have a sum of 0, indicating no null elements were found.

- No null elements found in any column

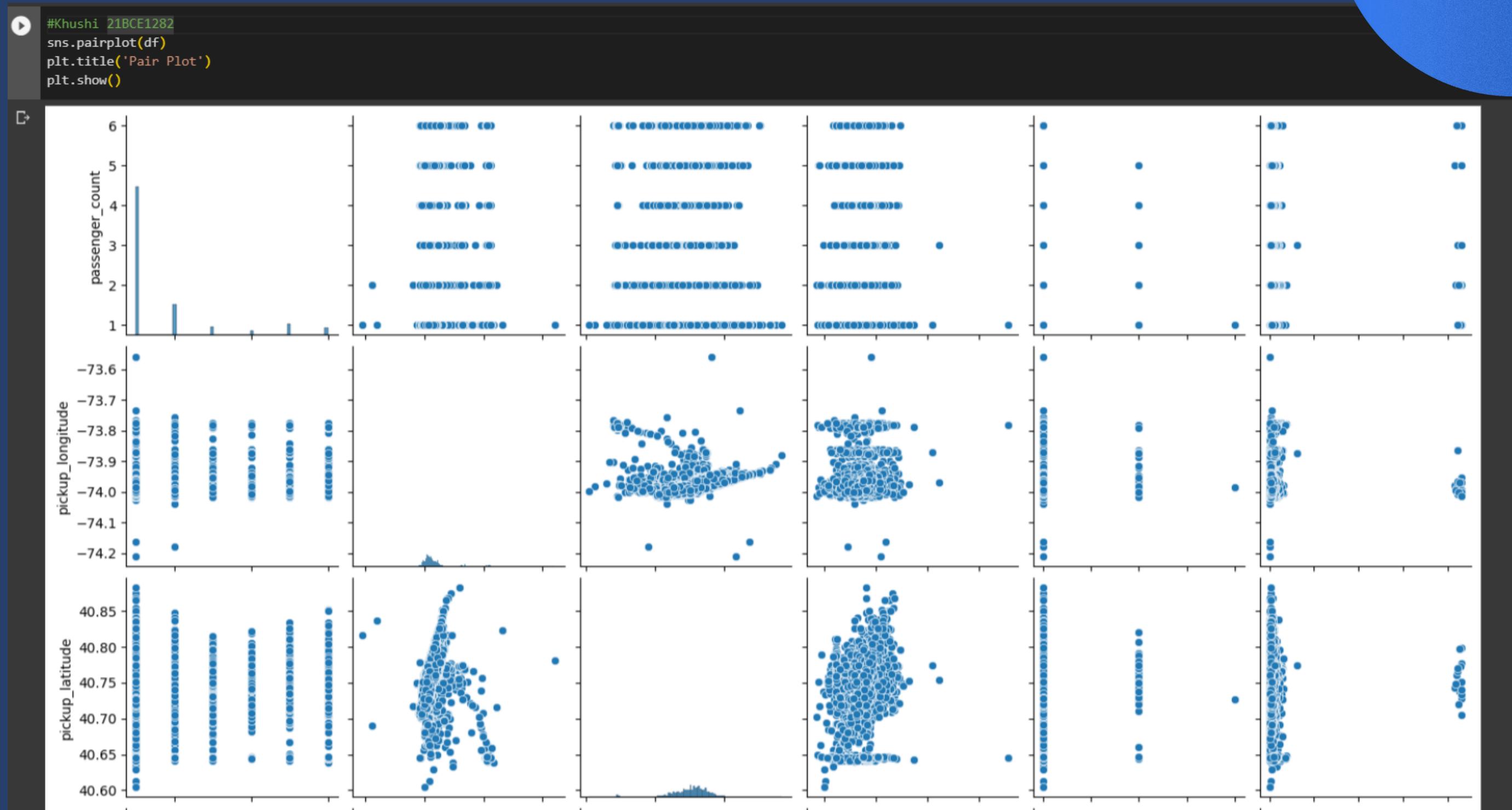
# Exploratory data analysis

A correlation matrix is a table that shows the correlation coefficients between multiple variables.



# Pairplot

A pairplot is a graphical representation that shows the relationships between pairs of variables in a dataset.



# Feature Selection

## Label Encoder

Label Encoder is a preprocessing technique used in machine learning to convert categorical variables into numerical form.

- In the given dataset, the column 'store\_and\_fwd\_flag' contains two variables i.e 'N', 'Y'
- It could be change into numerical values using Label encoder

```
#Khushi 21BCE1282
#Label encoder
from sklearn import preprocessing
label_encoder= preprocessing.LabelEncoder()
# Encode labels in column 'store_and_fwd_flag'.
df['store_and_fwd_flag']= label_encoder.fit_transform(df['store_and_fwd_flag'])
counts = df['store_and_fwd_flag'].value_counts()
print(counts)

0    1450599
1     8045
Name: store_and_fwd_flag, dtype: int64
```

0 0 0 0

+ Code + Markdown

```
[10]: #Khushi 21BCE1282
import seaborn as sns
```

# Seperating Target variable and Features

passeng  
pickup\_l  
pickup\_r  
dropoff\_l  
dropoff\_r  
trip\_

#Khushi 21BCE1282  
X=df.drop(['trip\_duration'],axis=1)  
y=df['trip\_duration']

# Min Max Scaling

- It scales the values of a feature to a fixed range, typically between 0 and 1,

The values in every column differ too much in range

Thus Min- Max scaling is used to range them  
between 0-1

This reduces non-linearity and model fits well to the dataset

Features are converted to data frame after fitting  
with min max scalar for training them model

The screenshot shows a Jupyter Notebook interface with two code cells and one output cell.

**Cell 14:**

```
#Khushi_21BCE1282
from sklearn.preprocessing import MinMaxScaler
min_max=MinMaxScaler()
min_max.fit(X)
```

**Cell 14 Output:**

```
MinMaxScaler()
```

**Cell 15:**

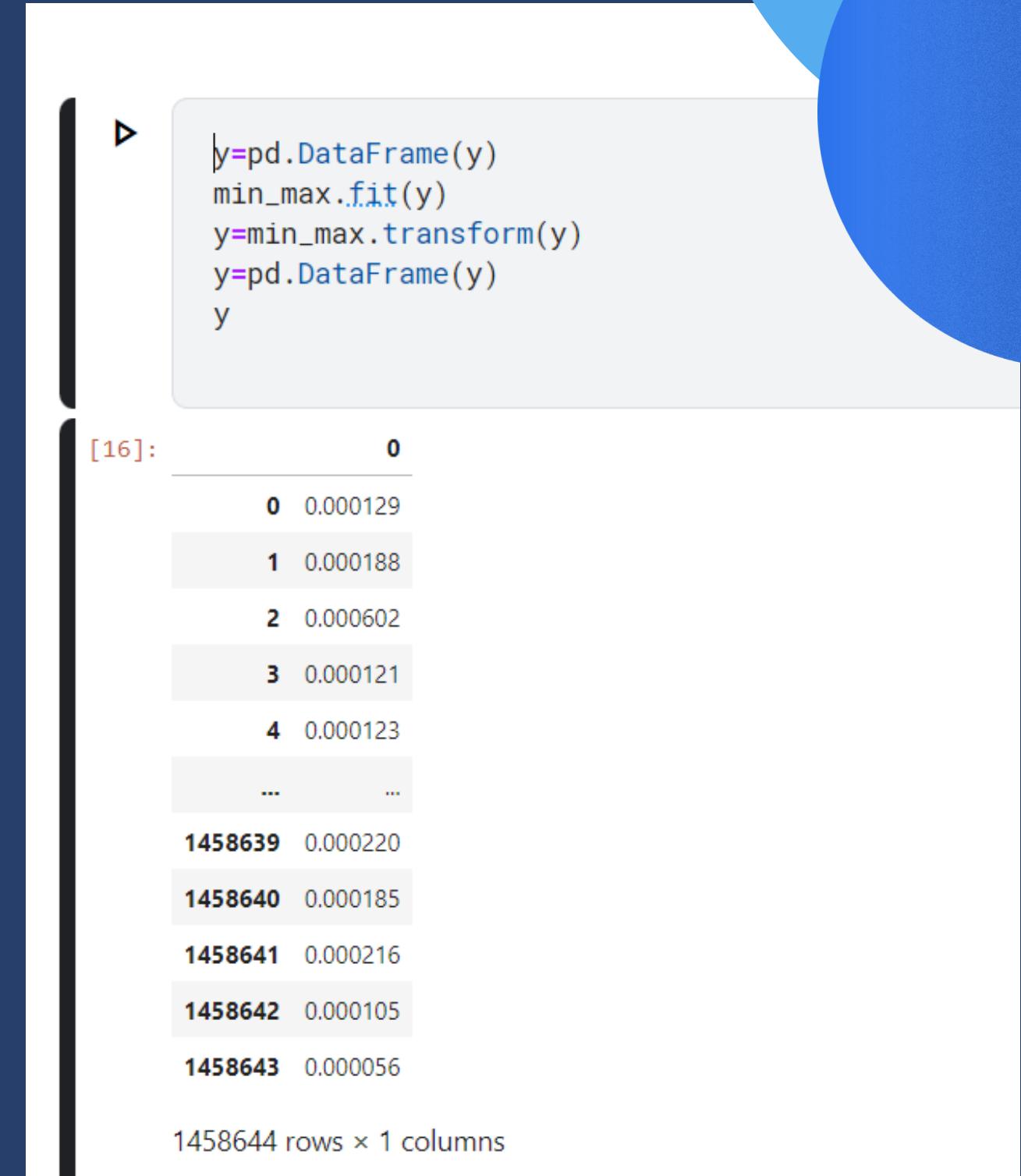
```
X=min_max.transform(X)
X=pd.DataFrame(X)
X
```

**Cell 15 Output:**

	0	1	2	3	4	5
0	0.111111	0.791302	0.365738	0.791591	0.731222	0.0
1	0.111111	0.791331	0.364062	0.791016	0.728287	0.0
2	0.111111	0.791354	0.365510	0.790920	0.726493	0.0
3	0.111111	0.790842	0.363001	0.790805	0.726206	0.0
4	0.111111	0.791452	0.367181	0.791454	0.732663	0.0
...	...	...	...	...	...	...
1458639	0.444444	0.791302	0.364459	0.791092	0.729055	0.0
1458640	0.111111	0.790992	0.364565	0.791500	0.733858	0.0
1458641	0.111111	0.791682	0.365787	0.790935	0.726262	0.0
1458642	0.111111	0.791304	0.364661	0.791426	0.730498	0.0

Target variable was first converted to dataframe since data in 2D format can only go under Min-Max scaling

Target variable are converted to data frame after fitting with min max scalar for training them model



```
y=pd.DataFrame(y)
min_max.fit(y)
y=min_max.transform(y)
y=pd.DataFrame(y)
y
```

[16]:

	0
0	0.000129
1	0.000188
2	0.000602
3	0.000121
4	0.000123
...	...
1458639	0.000220
1458640	0.000185
1458641	0.000216
1458642	0.000105
1458643	0.000056

1458644 rows × 1 columns

# Training the model

- **Linear Regression**

Linear regression is a statistical modeling technique used to analyze the relationship between a dependent variable and one or more independent variables

```
[18]:  
#Khushi 21BCE1282  
#splitting into training and testing data for creating model  
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=3)
```

```
[19]:  
from sklearn.linear_model import LinearRegression  
lr=LinearRegression()  
lr.fit(X_train,y_train)
```

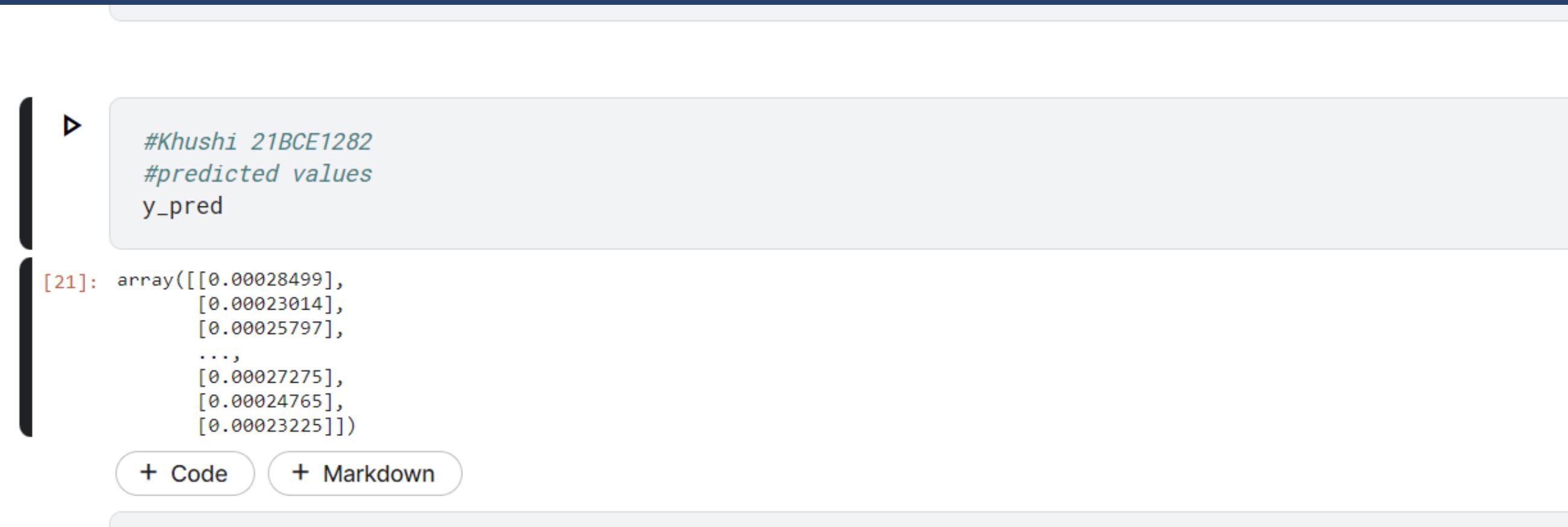
```
[19]: ▾ LinearRegression  
LinearRegression()
```

+ Code + Markdown

```
| ▶ y_pred=lr.predict(X_test)
```



## Predicted values using the Linear Regression based ML model



#Khushi 21BCE1282  
#predicted values  
y\_pred

```
[21]: array([[0.00028499],  
           [0.00023014],  
           [0.00025797],  
           ...,  
           [0.00027275],  
           [0.00024765],  
           [0.00023225]])
```

+ Code + Markdown

# *Testing of the model*

The primary goal of testing an ML model is to assess its performance and evaluate how well it generalizes to unseen data.

## **R2 Score**

It measures the proportion of the variance

## **Mean squared error**

It measures the average squared difference between the predicted values and the actual values

## **Mean absolute error**

It measures the average absolute difference between the predicted values and the actual values

```
[0.00024765],  
[0.00023225])  
  
+ Code + Markdown  
  
▶ #Khushi 21BCE1282  
from sklearn.metrics import r2_score,mean_squared_error,mean_absolute_error  
score = r2_score(y_test,y_pred)  
print("The accuracy of our model is {}%".format(score, 2))  
mae = mean_absolute_error(y_test,y_pred)  
print("Mean absolute error:",mae)  
mse = mean_squared_error(y_test,y_pred)  
print("Mean squared error: ",mse)  
  
The accuracy of our model is 0.0015590513645580284%  
Mean absolute error: 0.0001707185148601693  
Mean squared error: 1.8471075927914924e-06  
  
+ Code + Markdown
```

Error to be minimum.  
Thus our model fits the dataset accurately!!

# Conclusion

- Thus Regression based ML model is trained with the given NYC taxi data set with minimum error.



# THANK YOU

Created by  
**Khushi Singh**  
**21BCE1282**



Vellore Institute of technology,Chennai

7396207675

khushi.singh2021n@vitstudent.ac.in