

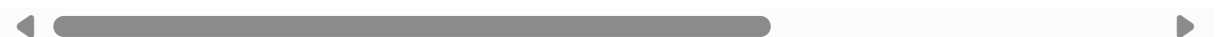
```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [2]: df=pd.read_csv("uber.csv")
df
```

Out[2]:

	Unnamed: 0	key	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude
0	24238194	2015-05-07 19:52:06.0000003	7.5	2015-05-07 19:52:06 UTC	-73.999817	40.738354
1	27835199	2009-07-17 20:04:56.0000002	7.7	2009-07-17 20:04:56 UTC	-73.994355	40.728225
2	44984355	2009-08-24 21:45:00.00000061	12.9	2009-08-24 21:45:00 UTC	-74.005043	40.740770
3	25894730	2009-06-26 08:22:21.0000001	5.3	2009-06-26 08:22:21 UTC	-73.976124	40.790844
4	17610152	2014-08-28 17:47:00.000000188	16.0	2014-08-28 17:47:00 UTC	-73.925023	40.744085
...
199995	42598914	2012-10-28 10:49:00.00000053	3.0	2012-10-28 10:49:00 UTC	-73.987042	40.739361
199996	16382965	2014-03-14 01:09:00.0000008	7.5	2014-03-14 01:09:00 UTC	-73.984722	40.736831
199997	27804658	2009-06-29 00:42:00.00000078	30.9	2009-06-29 00:42:00 UTC	-73.986017	40.756481
199998	20259894	2015-05-20 14:56:25.0000004	14.5	2015-05-20 14:56:25 UTC	-73.997124	40.725451
199999	11951496	2010-05-15 04:08:00.00000076	14.1	2010-05-15 04:08:00 UTC	-73.984395	40.720071

200000 rows × 9 columns



```
In [3]: df.head()
```

Out[3]:

	Unnamed: 0	key	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dr
0	24238194	2015-05-07 19:52:06.0000003	7.5	2015-05-07 19:52:06 UTC	-73.999817	40.738354	
1	27835199	2009-07-17 20:04:56.0000002	7.7	2009-07-17 20:04:56 UTC	-73.994355	40.728225	
2	44984355	2009-08-24 21:45:00.00000061	12.9	2009-08-24 21:45:00 UTC	-74.005043	40.740770	
3	25894730	2009-06-26 08:22:21.0000001	5.3	2009-06-26 08:22:21 UTC	-73.976124	40.790844	
4	17610152	2014-08-28 17:47:00.000000188	16.0	2014-08-28 17:47:00 UTC	-73.925023	40.744085	



In [4]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200000 entries, 0 to 199999
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Unnamed: 0            200000 non-null  int64
1   key                   200000 non-null  object
2   fare_amount           200000 non-null  float64
3   pickup_datetime       200000 non-null  object
4   pickup_longitude      200000 non-null  float64
5   pickup_latitude       200000 non-null  float64
6   dropoff_longitude     199999 non-null  float64
7   dropoff_latitude      199999 non-null  float64
8   passenger_count       200000 non-null  int64
dtypes: float64(5), int64(2), object(2)
memory usage: 13.7+ MB
```

In [5]: `df.columns`

Out[5]: Index(['Unnamed: 0', 'key', 'fare_amount', 'pickup_datetime', 'pickup_longitude', 'pickup_latitude', 'dropoff_longitude', 'dropoff_latitude', 'passenger_count'], dtype='object')

In [6]: `df=df.drop(['Unnamed: 0','key'],axis=1)`

In [7]: `df.head()`

Out[7]:

	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
0	7.5	2015-05-07 19:52:06 UTC	-73.999817	40.738354	-73.999512	40.723217
1	7.7	2009-07-17 20:04:56 UTC	-73.994355	40.728225	-73.994710	40.750325
2	12.9	2009-08-24 21:45:00 UTC	-74.005043	40.740770	-73.962565	40.772647
3	5.3	2009-06-26 08:22:21 UTC	-73.976124	40.790844	-73.965316	40.803349
4	16.0	2014-08-28 17:47:00 UTC	-73.925023	40.744085	-73.973082	40.761247

In [8]: `df.shape`

Out[8]: (200000, 7)

In [9]: `df.dtypes`

Out[9]:

```
fare_amount      float64
pickup_datetime  object
pickup_longitude float64
pickup_latitude  float64
dropoff_longitude float64
dropoff_latitude float64
passenger_count  int64
dtype: object
```

In [10]: df.describe

```
Out[10]: <bound method NDFrame.describe of
pickup_longitude \
0          7.5  2015-05-07 19:52:06 UTC          -73.999817
1          7.7  2009-07-17 20:04:56 UTC          -73.994355
2         12.9  2009-08-24 21:45:00 UTC          -74.005043
3          5.3  2009-06-26 08:22:21 UTC          -73.976124
4         16.0  2014-08-28 17:47:00 UTC          -73.925023
...
199995      3.0  2012-10-28 10:49:00 UTC          -73.987042
199996      7.5  2014-03-14 01:09:00 UTC          -73.984722
199997     30.9  2009-06-29 00:42:00 UTC          -73.986017
199998     14.5  2015-05-20 14:56:25 UTC          -73.997124
199999     14.1  2010-05-15 04:08:00 UTC          -73.984395

      pickup_latitude  dropoff_longitude  dropoff_latitude  passenger_count
0          40.738354          -73.999512          40.723217              1
1          40.728225          -73.994710          40.750325              1
2          40.740770          -73.962565          40.772647              1
3          40.790844          -73.965316          40.803349              3
4          40.744085          -73.973082          40.761247              5
...
199995      40.739367          -73.986525          40.740297              1
199996      40.736837          -74.006672          40.739620              1
199997      40.756487          -73.858957          40.692588              2
199998      40.725452          -73.983215          40.695415              1
199999      40.720077          -73.985508          40.768793              1

[200000 rows x 7 columns]>
```

In [11]: df.isnull().sum()

```
Out[11]: fare_amount          0
pickup_datetime          0
pickup_longitude         0
pickup_latitude          0
dropoff_longitude        1
dropoff_latitude         1
passenger_count          0
dtype: int64
```

In [12]: df['dropoff_latitude'].fillna(value=df['dropoff_latitude'].mean(),inplace=True)

In [13]: df.isnull().sum()

```
Out[13]: fare_amount          0
pickup_datetime          0
pickup_longitude         0
pickup_latitude          0
dropoff_longitude        1
dropoff_latitude         0
passenger_count          0
dtype: int64
```

In [14]: df['dropoff_longitude'].fillna(value=df['dropoff_longitude'].median(),inplace=True)

```
In [15]: df.isnull().sum()
```

```
Out[15]: fare_amount      0
pickup_datetime      0
pickup_longitude      0
pickup_latitude      0
dropoff_longitude      0
dropoff_latitude      0
passenger_count      0
dtype: int64
```

```
In [16]: df.dtypes
```

```
Out[16]: fare_amount      float64
pickup_datetime      object
pickup_longitude      float64
pickup_latitude      float64
dropoff_longitude      float64
dropoff_latitude      float64
passenger_count      int64
dtype: object
```

```
In [17]: df.pickup_datetime=pd.to_datetime(df.pickup_datetime,errors='coerce')
```

```
In [18]: df.dtypes
```

```
Out[18]: fare_amount      float64
pickup_datetime      datetime64[ns, UTC]
pickup_longitude      float64
pickup_latitude      float64
dropoff_longitude      float64
dropoff_latitude      float64
passenger_count      int64
dtype: object
```

```
In [19]: df=df.assign(hour=df.pickup_datetime.dt.hour,
                        day=df.pickup_datetime.dt.day,
                        month=df.pickup_datetime.dt.month,
                        year=df.pickup_datetime.dt.year,
                        dayofweek=df.pickup_datetime.dt.dayofweek)
```

```
In [20]: df.head()
```

```
Out[20]:
```

	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
0	7.5	2015-05-07 19:52:06+00:00	-73.999817	40.738354	-73.999512	40.723217
1	7.7	2009-07-17 20:04:56+00:00	-73.994355	40.728225	-73.994710	40.750325
2	12.9	2009-08-24 21:45:00+00:00	-74.005043	40.740770	-73.962565	40.772647
3	5.3	2009-06-26 08:22:21+00:00	-73.976124	40.790844	-73.965316	40.803349
4	16.0	2014-08-28 17:47:00+00:00	-73.925023	40.744085	-73.973082	40.761247

```
In [21]: df=df.drop('pickup_datetime',axis=1)
```

```
In [22]: df.head
```

```
Out[22]: <bound method NDFrame.head of
fare_amount pickup_longitude pickup_latit
ude dropoff_longitude \
0          7.5      -73.999817      40.738354      -73.999512
1          7.7      -73.994355      40.728225      -73.994710
2         12.9      -74.005043      40.740770      -73.962565
3          5.3      -73.976124      40.790844      -73.965316
4         16.0      -73.925023      40.744085      -73.973082
...         ...         ...         ...         ...
199995        3.0      -73.987042      40.739367      -73.986525
199996        7.5      -73.984722      40.736837      -74.006672
199997       30.9      -73.986017      40.756487      -73.858957
199998       14.5      -73.997124      40.725452      -73.983215
199999       14.1      -73.984395      40.720077      -73.985508

      dropoff_latitude passenger_count hour  day  month  year  dayofweek
0          40.723217             1    19   7     5  2015           3
1          40.750325             1    20  17     7  2009           4
2          40.772647             1    21  24     8  2009           0
3          40.803349             3     8  26     6  2009           4
4          40.761247             5    17  28     8  2014           3
...         ...         ...    ...    ...    ...    ...         ...
199995        40.740297             1    10  28    10  2012           6
199996        40.739620             1     1  14     3  2014           4
199997        40.692588             2     0  29     6  2009           0
199998        40.695415             1    14  20     5  2015           2
199999        40.768793             1     4  15     5  2010           5

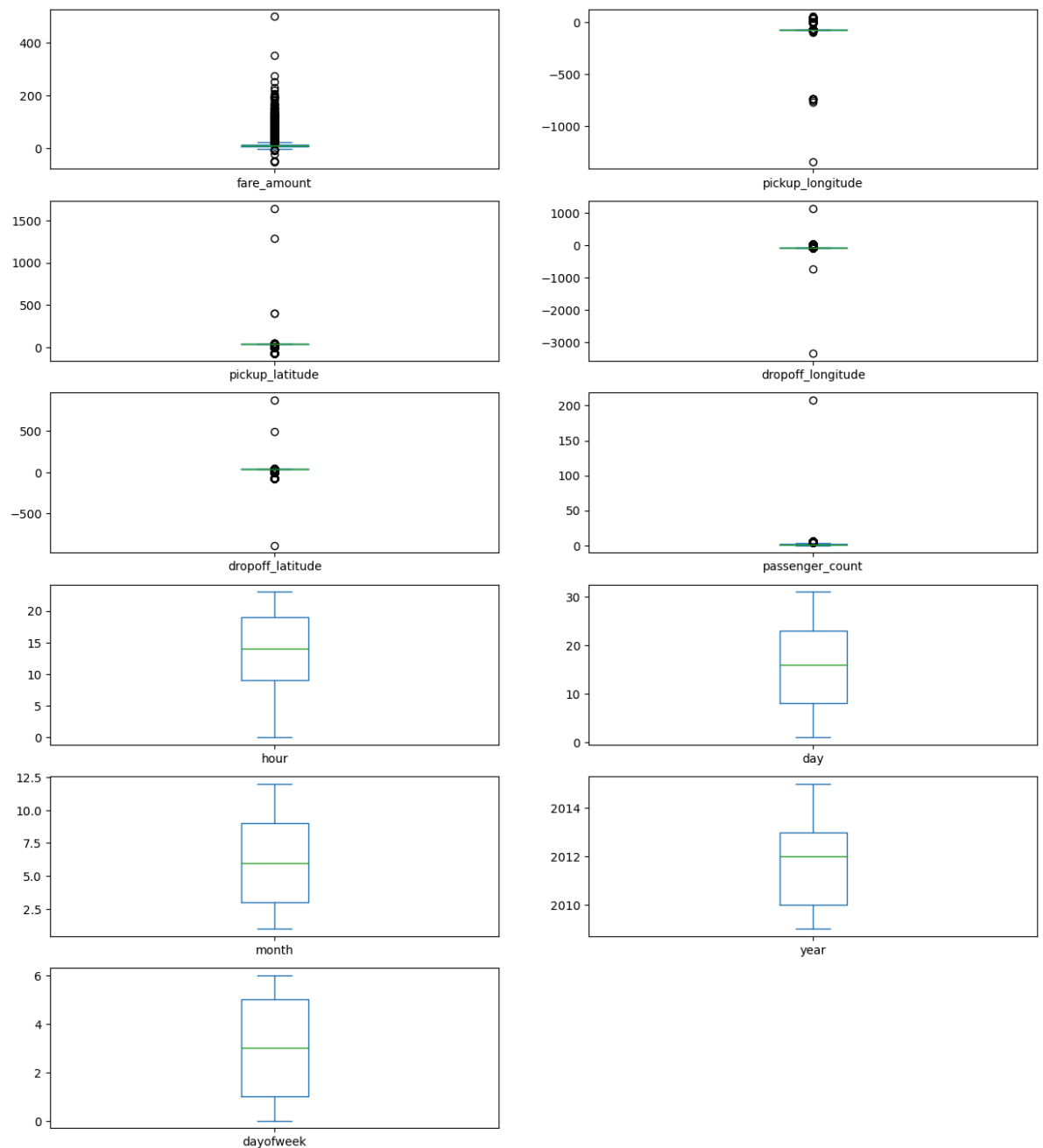
[200000 rows x 11 columns]>
```

```
In [23]: df.dtypes
```

```
Out[23]: fare_amount      float64
pickup_longitude    float64
pickup_latitude     float64
dropoff_longitude   float64
dropoff_latitude    float64
passenger_count      int64
hour                 int32
day                  int32
month                int32
year                 int32
dayofweek            int32
dtype: object
```

```
In [24]: df.plot(kind="box",subplots=True,layout=(7,2),figsize=(15,20))
```

```
Out[24]: fare_amount      Axes(0.125,0.786098;0.352273x0.0939024)
pickup_longitude Axes(0.547727,0.786098;0.352273x0.0939024)
pickup_latitude  Axes(0.125,0.673415;0.352273x0.0939024)
dropoff_longitude Axes(0.547727,0.673415;0.352273x0.0939024)
dropoff_latitude  Axes(0.125,0.560732;0.352273x0.0939024)
passenger_count   Axes(0.547727,0.560732;0.352273x0.0939024)
hour              Axes(0.125,0.448049;0.352273x0.0939024)
day              Axes(0.547727,0.448049;0.352273x0.0939024)
month            Axes(0.125,0.335366;0.352273x0.0939024)
year             Axes(0.547727,0.335366;0.352273x0.0939024)
dayofweek        Axes(0.125,0.222683;0.352273x0.0939024)
dtype: object
```



```
In [25]: def remove_outlier(df1,col):  
        Q1=df1[col].quantile(0.25)  
        Q3=df1[col].quantile(0.75)  
        IQR=Q3-Q1  
        lower_whisker=Q1-1.5*IQR  
        upper_whisker=Q3+1.5*IQR  
        df[col]=np.clip(df1[col],lower_whisker,upper_whisker)  
        return df1  
def treat_outliers_all(df1,col_list):  
    for c in col_list:  
        df1=remove_outlier(df,c)  
    return df1
```

```
In [26]: df = treat_outliers_all(df,df.iloc[:,0::])
```

```
In [29]: pip install haversine
```

Requirement already satisfied: haversine in c:\users\suwasini\anaconda3\lib\site-packages (2.8.1)

Note: you may need to restart the kernel to use updated packages.

```
In [30]: import haversine as hs
travel_dist = []
for pos in range(len(df['pickup_longitude'])):
    long1,lati1,long2,lati2 = [df['pickup_longitude'][pos],df['pickup_latitude'][pos],df['dropoff_longitude'][pos],df['dropoff_latitude'][pos]]
    loc1 =(lati1,long1)
    loc2 =(lati2,long2)
    c = hs.haversine(loc1,loc2)
    travel_dist.append(c)

print(travel_dist)
df['dist_travel_km']= travel_dist
df.head
```

IOPub data rate exceeded.

The notebook server will temporarily stop sending output to the client in order to avoid crashing it.

To change this limit, set the config variable

`--NotebookApp.iopub_data_rate_limit`.

Current values:

NotebookApp.iopub_data_rate_limit=1000000.0 (bytes/sec)

NotebookApp.rate_limit_window=3.0 (secs)


```
Out[30]: <bound method NDFrame.head of
         fare_amount pickup_longitude pickup_latit
         ude dropoff_longitude \
0          7.50      -73.999817      40.738354      -73.999512
1          7.70      -73.994355      40.728225      -73.994710
2         12.90      -74.005043      40.740770      -73.962565
3          5.30      -73.976124      40.790844      -73.965316
4         16.00      -73.929786      40.744085      -73.973082
...         ...         ...         ...         ...
199995        3.00      -73.987042      40.739367      -73.986525
199996        7.50      -73.984722      40.736837      -74.006672
199997       22.25      -73.986017      40.756487      -73.922036
199998       14.50      -73.997124      40.725452      -73.983215
199999       14.10      -73.984395      40.720077      -73.985508

         dropoff_latitude passenger_count hour day month year dayofweek \
0          40.723217          1.0    19   7   5  2015          3
1          40.750325          1.0    20  17   7  2009          4
2          40.772647          1.0    21  24   8  2009          0
3          40.803349          3.0     8  26   6  2009          4
4          40.761247          3.5    17  28   8  2014          3
...         ...         ...    ...   ...   ...   ...         ...
199995        40.740297          1.0    10  28  10  2012          6
199996        40.739620          1.0     1  14   3  2014          4
199997        40.692588          2.0     0  29   6  2009          0
199998        40.695415          1.0    14  20   5  2015          2
199999        40.768793          1.0     4  15   5  2010          5

         dist_travel_km
0          1.683325
1          2.457593
2          5.036384
3          1.661686
4          4.116088
...         ...
199995        0.112210
199996        1.875053
199997        8.919323
199998        3.539720
199999        5.417791
```

[200000 rows x 12 columns]>

```
In [31]: df = df.loc[(df.dist_travel_km>=1) |(df.dist_travel_km<=130) ]
         print("Remaining obervation:" , df.shape)
```

Remaining obervation: (200000, 12)


```
In [32]: incorrect_coordinates =df.loc[(df.pickup_latitude>90) | (df.pickup_latitude< -90)|
         (df.dropoff_latitude>90) | (df.dropoff_latitude< -90)
         (df.pickup_longitude>180) | (df.pickup_longitude< -18
         (df.dropoff_latitude>90) | (df.dropoff_latitude< -90
```

```
In [33]: df.drop(incorrect_coordinates, inplace = True, errors = 'ignore')
```

```
In [34]: df.head()
```

```
Out[34]:
```

	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	passenger_count
0	7.5	-73.999817	40.738354	-73.999512	40.723217	1.0
1	7.7	-73.994355	40.728225	-73.994710	40.750325	1.0
2	12.9	-74.005043	40.740770	-73.962565	40.772647	1.0
3	5.3	-73.976124	40.790844	-73.965316	40.803349	3.0
4	16.0	-73.929786	40.744085	-73.973082	40.761247	3.0

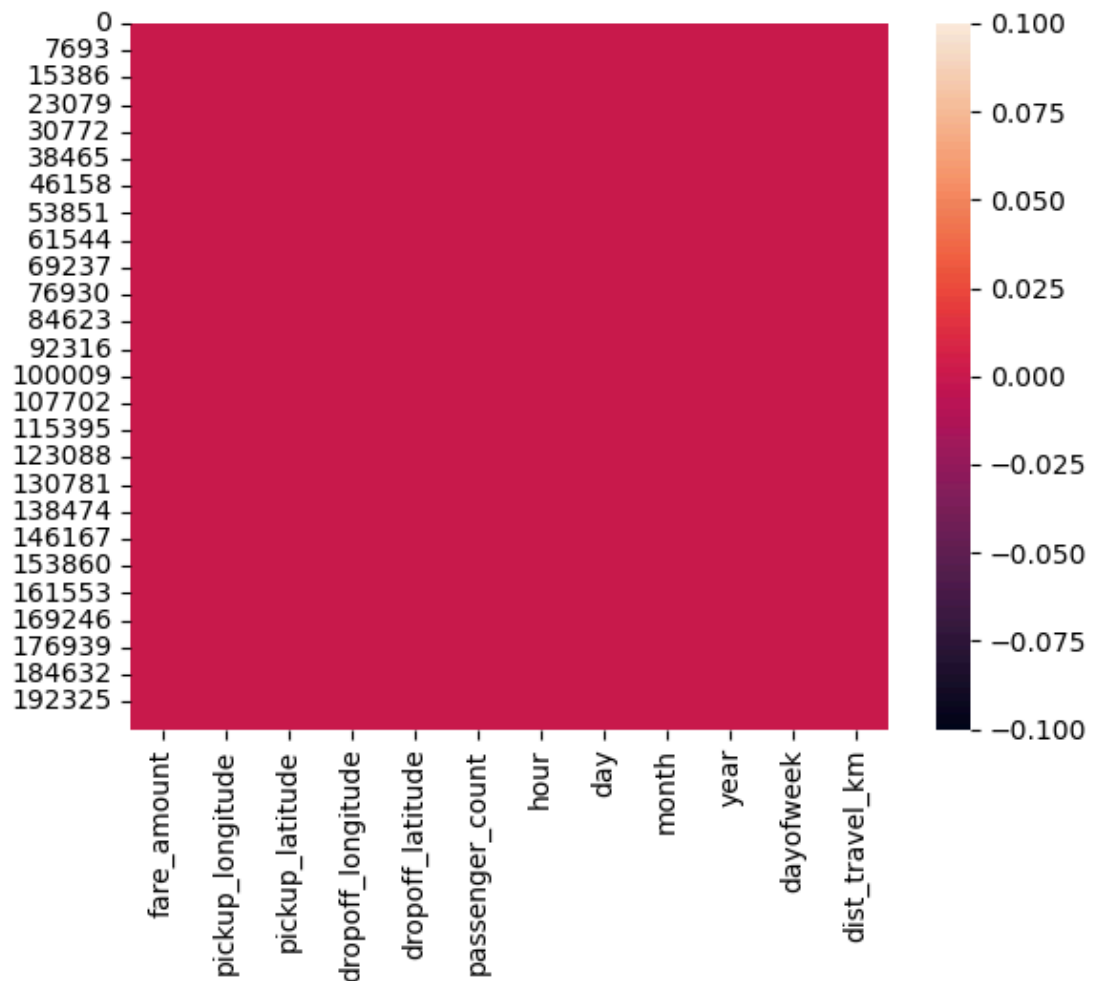


```
In [35]: df.isnull().sum()
```

```
Out[35]: fare_amount      0
pickup_longitude      0
pickup_latitude      0
dropoff_longitude      0
dropoff_latitude      0
passenger_count      0
hour      0
day      0
month      0
year      0
dayofweek      0
dist_travel_km      0
dtype: int64
```

In [36]: `sns.heatmap(df.isnull())`

Out[36]: <Axes: >



In [37]: `corr = df.corr()`

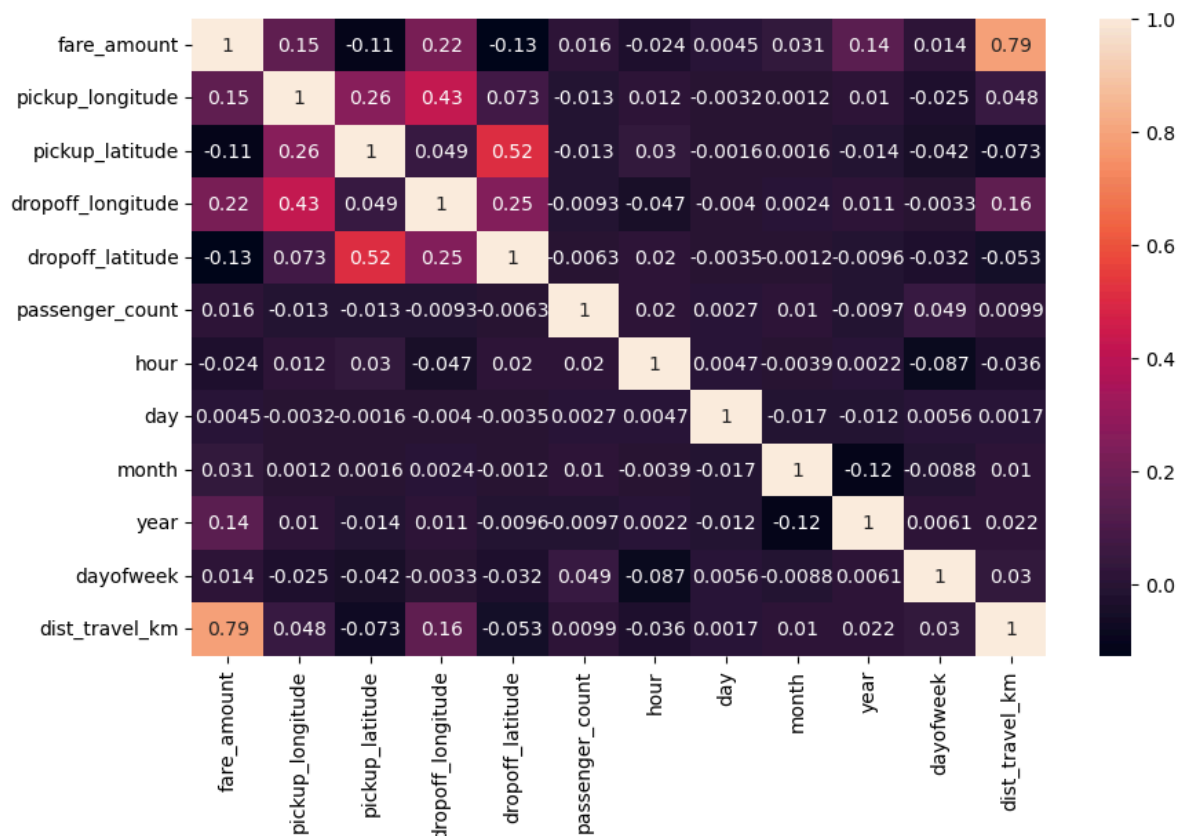
In [38]: `corr`

Out[38]:

	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
fare_amount	1.000000	0.154069	-0.110842	0.218675	-0.125898
pickup_longitude	0.154069	1.000000	0.259497	0.425619	0.073290
pickup_latitude	-0.110842	0.259497	1.000000	0.048889	0.515714
dropoff_longitude	0.218675	0.425619	0.048889	1.000000	0.245667
dropoff_latitude	-0.125898	0.073290	0.515714	0.245667	1.000000
passenger_count	0.015778	-0.013213	-0.012889	-0.009303	-0.006308
hour	-0.023623	0.011579	0.029681	-0.046558	0.019783
day	0.004534	-0.003204	-0.001553	-0.004007	-0.003479
month	0.030817	0.001169	0.001562	0.002391	-0.001193
year	0.141277	0.010198	-0.014243	0.011346	-0.009603
dayofweek	0.013652	-0.024652	-0.042310	-0.003336	-0.031919
dist_travel_km	0.786385	0.048446	-0.073362	0.155191	-0.052701

```
In [39]: fig,axis = plt.subplots(figsize= (10,6))
sns.heatmap(df.corr(), annot = True)
```

Out[39]: <Axes: >



```
In [40]: df.dtypes
```

```
Out[40]: fare_amount      float64
pickup_longitude    float64
pickup_latitude     float64
dropoff_longitude   float64
dropoff_latitude    float64
passenger_count     float64
hour                int32
day                 int32
month               int32
year                int32
dayofweek           int32
dist_travel_km      float64
dtype: object
```

```
In [41]: x = df[['pickup_longitude', 'pickup_latitude', 'dropoff_longitude', 'dropoff_latitude', 'passenger_count', 'hour', 'day', 'month', 'year', 'dayofweek', 'dist_travel_km']]
```

```
In [42]: y = df['fare_amount']
```

```
In [43]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.33)
```

```
In [44]: from sklearn.linear_model import LinearRegression
         regression = LinearRegression()
```

```
In [45]: regression.fit(x_train,y_train)
```

```
Out[45]: 

LinearRegression ⓘ ⓘ  
(https://scikit-learn.org/1.4/modules/generated/sklearn.linear_model.LinearRegression.  
LinearRegression())


```

```
In [46]: regression.intercept_
```

```
Out[46]: 3670.958009297504
```

```
In [47]: regression.coef_
```

```
Out[47]: array([ 2.55467049e+01, -7.07837319e+00,  2.00180481e+01, -1.84664321e+01,  
                7.13082501e-02,  5.03464642e-03,  4.20083396e-03,  6.02329544e-02,  
                3.70408277e-01, -3.40314100e-02,  1.84487240e+00])
```

```
In [48]: prediction = regression.predict(x_test)
```

```
In [49]: print(prediction)
```

```
[ 6.47598495 25.25315024 11.89085438 ...  5.26424499  9.26924157  
 9.92743844]
```

```
In [50]: y_test
```

```
Out[50]: 115141    6.10  
         45538    5.50  
         193115   17.00  
         100757   13.70  
         114572   16.90  
         ...  
         118577   22.25  
         137346    8.00  
         62892    5.30  
         98580    9.50  
         38342    9.30  
Name: fare_amount, Length: 66000, dtype: float64
```

```
In [51]: from sklearn.metrics import r2_score
```

```
In [52]: r2_score(y_test,prediction)
```

```
Out[52]: 0.6659814522016027
```

```
In [53]: from sklearn.metrics import mean_squared_error
```

```
In [54]: MSE = mean_squared_error(y_test,prediction)
```

In [55]: MSE

Out[55]: 9.855757486376126

In []: *#SUWASINI CHABUKSWAR*
#DIV:A
#BATCH:A2
#ROLL NO:14124