

Deena 20104016

Basic Analysis using Numpy and Pandas

Import Libraries

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

Importing Dataset

```
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	40.7
1	27835199	2009-07-17 20:04:56.0000002	7.7	2009-07-17 20:04:56 UTC	-73.994355	40.7
2	44984355	2009-08-24 21:45:00.00000061	12.9	2009-08-24 21:45:00 UTC	-74.005043	40.7
3	25894730	2009-06-26 08:22:21.0000001	5.3	2009-06-26 08:22:21 UTC	-73.976124	40.7
4	17610152	2014-08-28 17:47:00.000000188	16.0	2014-08-28 17:47:00 UTC	-73.925023	40.7
...
199995	42598914	2012-10-28 10:49:00.00000053	3.0	2012-10-28 10:49:00 UTC	-73.987042	40.7
199996	16382965	2014-03-14 01:09:00.0000008	7.5	2014-03-14 01:09:00 UTC	-73.984722	40.7
199997	27804658	2009-06-29 00:42:00.00000078	30.9	2009-06-29 00:42:00 UTC	-73.986017	40.7
199998	20259894	2015-05-20 14:56:25.0000004	14.5	2015-05-20 14:56:25 UTC	-73.997124	40.7
199999	11951496	2010-05-15 04:08:00.00000076	14.1	2010-05-15 04:08:00 UTC	-73.984395	40.7

200000 rows × 9 columns



To display first 10 rows

In [3]: `df.head(10)`

Out[3]:

	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
5	44470845	2011-02-12 02:27:09.0000006	4.9	2011-02-12 02:27:09 UTC	-73.969019	40.755910
6	48725865	2014-10-12 07:04:00.0000002	24.5	2014-10-12 07:04:00 UTC	-73.961447	40.693965
7	44195482	2012-12-11 13:52:00.00000029	2.5	2012-12-11 13:52:00 UTC	0.000000	0.000000
8	15822268	2012-02-17 09:32:00.00000043	9.7	2012-02-17 09:32:00 UTC	-73.975187	40.745767
9	50611056	2012-03-29 19:06:00.000000273	12.5	2012-03-29 19:06:00 UTC	-74.001065	40.741787

To display last 5 rows

In [4]: `df.tail(5)`

Out[4]:

	Unnamed: 0	key	fare_amount	pickup_datetime	pickup_longitude	pickup_lat
199995	42598914	2012-10-28 10:49:00.00000053	3.0	2012-10-28 10:49:00 UTC	-73.987042	40.73
199996	16382965	2014-03-14 01:09:00.0000008	7.5	2014-03-14 01:09:00 UTC	-73.984722	40.73
199997	27804658	2009-06-29 00:42:00.00000078	30.9	2009-06-29 00:42:00 UTC	-73.986017	40.75
199998	20259894	2015-05-20 14:56:25.0000004	14.5	2015-05-20 14:56:25 UTC	-73.997124	40.72
199999	11951496	2010-05-15 04:08:00.00000076	14.1	2010-05-15 04:08:00 UTC	-73.984395	40.72

Satistical Summary

In [5]: `df.describe()`

Out[5]:

	Unnamed: 0	fare_amount	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
count	2.000000e+05	200000.000000	200000.000000	200000.000000	199999.000000	199999.000000
mean	2.771250e+07	11.359955	-72.527638	39.935885	-72.525292	39.935885
std	1.601382e+07	9.901776	11.437787	7.720539	13.117408	7.720539
min	1.000000e+00	-52.000000	-1340.648410	-74.015515	-3356.666300	-88.015515
25%	1.382535e+07	6.000000	-73.992065	40.734796	-73.991407	40.734796
50%	2.774550e+07	8.500000	-73.981823	40.752592	-73.980093	40.752592
75%	4.155530e+07	12.500000	-73.967154	40.767158	-73.963658	40.767158
max	5.542357e+07	499.000000	57.418457	1644.421482	1153.572603	87.418457

To find shape and size

In [6]: `df.shape`

Out[6]: (200000, 9)

In [7]: `df.size`

Out[7]: 1800000

To fill the null values

In [8]:

df.isna()

Out[8]:

	Unnamed: 0	key	fare_amount	pickup_datetime	pickup_longitude	pickup_latitude	dropo
0	False	False	False	False	False	False	
1	False	False	False	False	False	False	
2	False	False	False	False	False	False	
3	False	False	False	False	False	False	
4	False	False	False	False	False	False	
...
199995	False	False	False	False	False	False	
199996	False	False	False	False	False	False	
199997	False	False	False	False	False	False	
199998	False	False	False	False	False	False	
199999	False	False	False	False	False	False	

200000 rows × 9 columns

To fill missing values

In [9]:

df.dropna()

3	25894730	2009-06-26 08:22:21.0000001	5.3	2009-06-26 08:22:21 UTC	-73.976124	4
4	17610152	2014-08-28 17:47:00.000000188	16.0	2014-08-28 17:47:00 UTC	-73.925023	4
...
199995	42598914	2012-10-28 10:49:00.00000053	3.0	2012-10-28 10:49:00 UTC	-73.987042	4
199996	16382965	2014-03-14 01:09:00.0000008	7.5	2014-03-14 01:09:00 UTC	-73.984722	4
199997	27804658	2009-06-29 00:42:00.00000078	30.9	2009-06-29 00:42:00 UTC	-73.986017	4
199998	20259894	2015-05-20 14:56:25.0000004	14.5	2015-05-20 14:56:25 UTC	-73.997124	4
199999	11951496	2010-05-15 04:08:00.00000076	14.1	2010-05-15 04:08:00 UTC	-73.984395	4

199999 rows × 9 columns

columns

```
In [10]: df.columns
```

```
Out[10]: Index(['Unnamed: 0', 'key', 'fare_amount', 'pickup_datetime',  
              'pickup_longitude', 'pickup_latitude', 'dropoff_longitude',  
              'dropoff_latitude', 'passenger_count'],  
              dtype='object')
```

to print a particular coloumn

```
In [11]: data=df[["passenger_count", "fare_amount"]][0:500]  
data
```

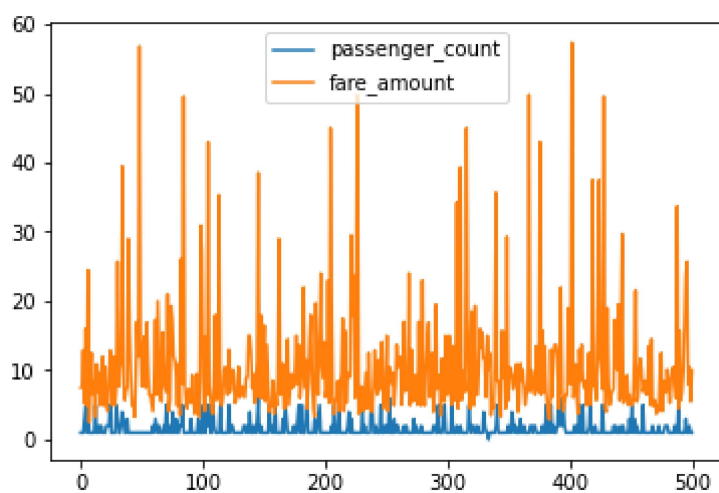
	passenger_count	fare_amount
0	1	7.5
1	1	7.7
2	1	12.9
3	3	5.3
4	5	16.0
...
495	1	25.7
496	1	8.0
497	2	10.5
498	1	5.5
499	1	10.0

500 rows × 2 columns

line plot

```
In [12]: data.plot.line()
```

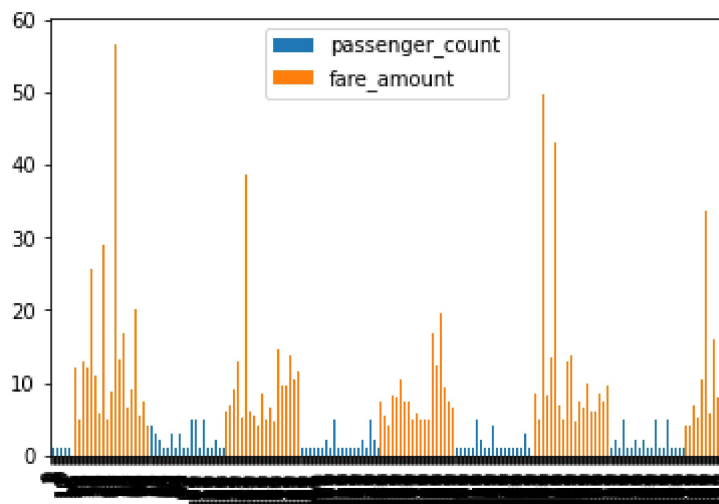
```
Out[12]: <AxesSubplot:>
```



bar plot

```
In [13]: data.plot.bar()
```

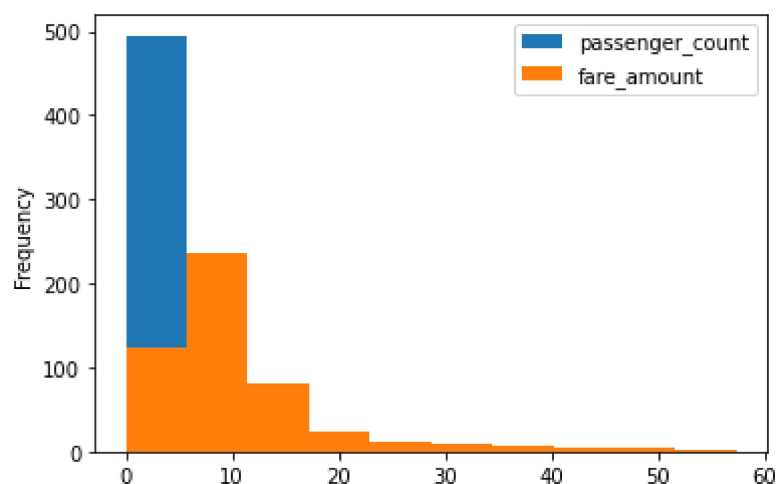
```
Out[13]: <AxesSubplot:>
```



hist plot

```
In [14]: data.plot.hist()
```

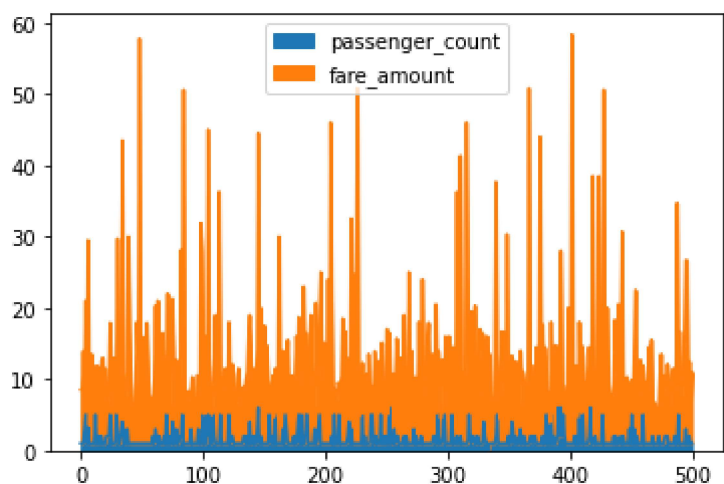
```
Out[14]: <AxesSubplot:ylabel='Frequency'>
```



Area plot

```
In [15]: data.plot.area()
```

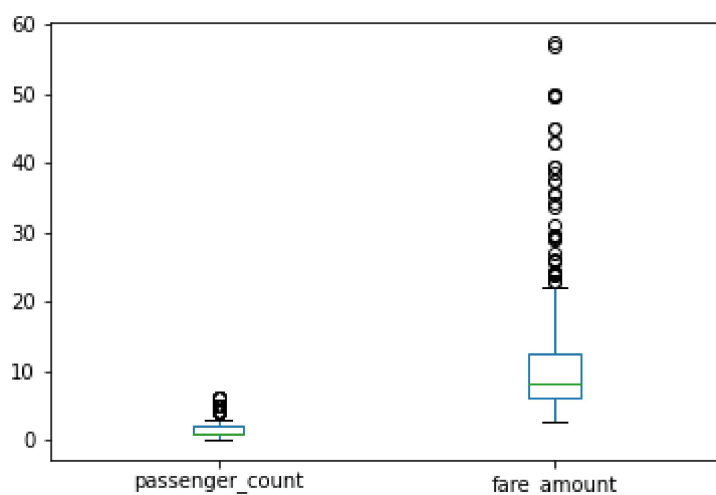
```
Out[15]: <AxesSubplot:>
```



Box plot

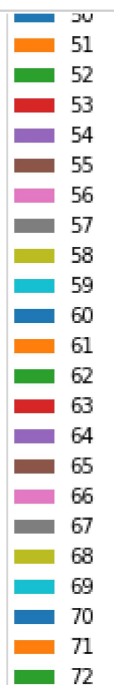
```
In [16]: data.plot.box()
```

```
Out[16]: <AxesSubplot:>
```



pie plot

```
In [17]: data.plot.pie(y="passenger_count")
```




```
In [18]: data.plot.scatter(x="passenger_count",y="fare_amount")
```

```
Out[18]: <AxesSubplot:xlabel='passenger_count', ylabel='fare_amount'>
```

