

Working on Real Project with Python

(A part of Big Data Analysis)

Police Dataset

Here, The data from a Police check post is given The data is available as a csv file.

```
In [2]: import pandas as pd
```

```
In [3]: data = pd.read_csv('3. Police Data.csv')
data.head()
```

Out[3]:	stop_date	stop_time	country_name	driver_gender	driver_age_raw	driver_age	driver_race	violation_raw	violation	search_conducted	search_type	stop_outcome	is_arrested	stop_duration	drugs_related_stop
0	1/2/2005	1:55	NaN	M	1985.0	20.0	White	Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False
1	1/18/2005	8:15	NaN	M	1965.0	40.0	White	Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False
2	1/23/2005	23:15	NaN	M	1972.0	33.0	White	Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False
3	2/20/2005	17:15	NaN	M	1986.0	19.0	White	Call for Service	Other	False	NaN	Arrest Driver	True	16-30 Min	False
4	3/14/2005	10:00	NaN	F	1984.0	21.0	White	Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False

```
In [4]: data.shape
```

```
Out[4]: (65535, 15)
```

For Data Cleaning

Remove the column that only contains missing values

```
In [6]: data.isnull()
```

Out[6]:	stop_date	stop_time	country_name	driver_gender	driver_age_raw	driver_age	driver_race	violation_raw	violation	search_conducted	search_type	stop_outcome	is_arrested	stop_duration	drugs_related_stop
0	False	False	True	False	False	False	False	False	False	False	True	False	False	False	False
1	False	False	True	False	False	False	False	False	False	False	True	False	False	False	False
2	False	False	True	False	False	False	False	False	False	False	True	False	False	False	False
3	False	False	True	False	False	False	False	False	False	False	True	False	False	False	False
4	False	False	True	False	False	False	False	False	False	False	True	False	False	False	False
...
65530	False	False	True	False	False	False	False	False	False	False	True	False	False	False	False
65531	False	False	True	False	False	False	False	False	False	False	True	False	False	False	False
65532	False	False	True	False	False	False	False	False	False	False	True	False	False	False	False
65533	False	False	True	True	True	True	True	True	True	True	True	True	True	True	False
65534	False	False	True	False	False	False	False	False	False	False	True	False	False	False	False

65535 rows × 15 columns

```
In [7]: data.isnull().sum()
```

```
Out[7]: stop_date      0
stop_time      0
country_name    65535
driver_gender   4861
driver_age_raw  4054
driver_age      4387
driver_race     4869
violation_raw   4869
violation       4969
search_conducted 0
search_type     63956
stop_outcome    4969
is_arrested     4869
stop_duration   4869
drugs_related_stop 0
dtype: int64
```

```
In [9]: data.drop(columns = 'country_name', inplace = True)
```

```
In [10]: data
```

2	1/23/2005	23:15	M	1972.0	33.0	White		Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False
3	2/20/2005	17:15	M	1986.0	19.0	White	Call for Service	Other	False	NaN	Arrest Driver	True	16-30 Min	False	False
4	3/14/2005	10:00	F	1984.0	21.0	White		Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False
...
65530	12/6/2012	17:54	F	1987.0	25.0	White		Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False
65531	12/6/2012	22:22	M	1954.0	58.0	White		Speeding	Speeding	False	NaN	Warning	False	0-15 Min	False
65532	12/6/2012	23:20	M	1985.0	27.0	Black	Equipment/Inspection Violation	Equipment	False	NaN	Citation	False	0-15 Min	False	False
65533	12/7/2012	0:23	NaN	NaN	NaN	NaN	NaN	NaN	False	NaN	NaN	NaN	NaN	NaN	False
65534	12/7/2012	0:30	F	1985.0	27.0	White		Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False

65535 rows × 14 columns

Based on Filtering and value plus

For speeding,were Men or Women stopped more often?

```
In [11]: data.head()
```

Out[11]:	stop_date	stop_time	driver_gender	driver_age_raw	driver_age	driver_race	violation_raw	violation	search_conducted	search_type	stop_outcome	is_arrested	stop_duration	drugs_related_stop
0	1/2/2005	1:55	M	1985.0	20.0	White	Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False
1	1/18/2005	8:15	M	1965.0	40.0	White	Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False
2	1/23/2005	23:15	M	1972.0	33.0	White	Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False
3	2/20/2005	17:15	M	1986.0	19.0	White	Call for Service	Other	False	NaN	Arrest Driver	True	16-30 Min	False
4	3/14/2005	10:00	F	1984.0	21.0	White	Speeding	Speeding	False	NaN	Citation	False	0-15 Min	False

```
In [14]: data[data.violation == 'Speeding'].driver_gender.value_counts()
```

```
Out[14]: M    25517
F       11686
Name: driver_gender, dtype: int64
```

Groupby

Does gender affect who gets searched during a stop?

```
In [15]: data.groupby('driver_gender').search_conducted.sum()
```

```
Out[15]: driver_gender
F         366
M        2113
Name: search_conducted, dtype: int64
```

```
In [16]: data.search_conducted.value_counts()
```

```
Out[16]: False    63956
True       2479
Name: search_conducted, dtype: int64
```

Mapping and Data-Type casting

What is the maen stop_duration?

```
In [17]: data.stop_duration.value_counts()
```

```
Out[17]: 0-15 Min    47379
16-30 Min    11448
30+ Min      2647
2              1
Name: stop_duration, dtype: int64
```

```
In [18]: data['stop_duration'] = data['stop_duration'].map({'0-15 Min': 7.5 , '16-30 Min' : 24 , '30+ Min' : 45})
```

```
In [19]: data
```

Out[19]:	stop_date	stop_time	driver_gender	driver_age_raw	driver_age	driver_race	violation_raw	violation	search_conducted	search_type	stop_outcome	is_arrested	stop_duration	drugs_related_stop
0	1/2/2005	1:55	M	1985.0	20.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
1	1/18/2005	8:15	M	1965.0	40.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
2	1/23/2005	23:15	M	1972.0	33.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
3	2/20/2005	17:15	M	1986.0	19.0	White	Call for Service	Other	False	NaN	Arrest Driver	True	24.0	False
4	3/14/2005	10:00	F	1984.0	21.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
...
65530	12/6/2012	17:54	F	1987.0	25.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
65531	12/6/2012	22:22	M	1954.0	58.0	White	Speeding	Speeding	False	NaN	Warning	False	7.5	False
65532	12/6/2012	23:20	M	1985.0	27.0	Black	Equipment/Inspection Violation	Equipment	False	NaN	Citation	False	7.5	False
65533	12/7/2012	0:23	NaN	NaN	NaN	NaN	NaN	NaN	False	NaN	NaN	NaN	NaN	False
65534	12/7/2012	0:30	F	1985.0	27.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False

65535 rows × 14 columns

```
In [20]: data['stop_duration'].mean()
```

```
Out[20]: 12.187420698181345
```

Groupby , Describe

Compare the age distribution for each violation

```
In [21]: data.groupby('violation').driver_age.describe()
```

Out[21]:		count	mean	std	min	25%	50%	75%	max
	violation								
	Equipment	6507.0	31.682957	11.380671	16.0	23.0	28.0	39.0	81.0
	Moving Violation	11276.0	36.736443	13.258390	15.0	25.0	35.0	47.0	86.0
	Other	3477.0	40.362381	12.754423	16.0	30.0	41.0	50.0	86.0
	Registration/plates	2240.0	32.656696	11.150780	16.0	24.0	30.0	40.0	74.0
	Seat belt	3.0	30.333333	10.214369	23.0	24.5	26.0	34.0	42.0
	Speeding	37120.0	33.262981	12.615781	15.0	23.0	30.0	42.0	88.0

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

Matplotlib

```
In [23]: import numpy as np
import matplotlib.pyplot as plt
```

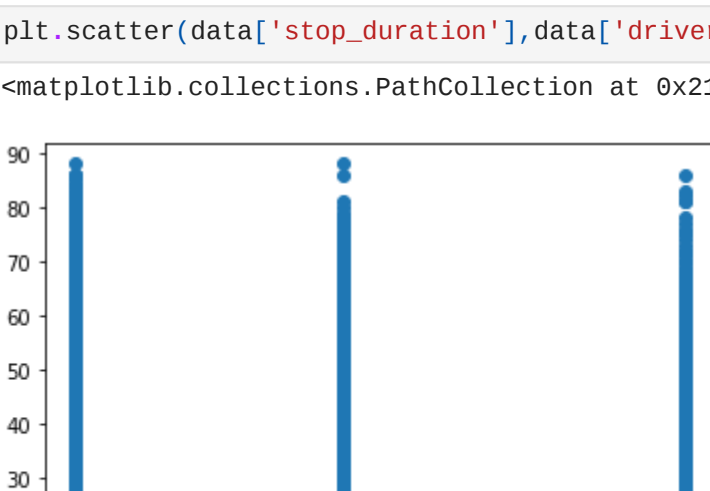
```
In [24]: data.head()
```

Out[24]:	stop_date	stop_time	driver_gender	driver_age_raw	driver_age	driver_race	violation_raw	violation	search_conducted	search_type	stop_outcome	is_arrested	stop_duration	drugs_related_stop
0	1/2/2005	1:55	M	1985.0	20.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
1	1/18/2005	8:15	M	1965.0	40.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
2	1/23/2005	23:15	M	1972.0	33.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
3	2/20/2005	17:15	M	1986.0	19.0	White	Call for Service	Other	False	NaN	Arrest Driver	True	24.0	False
4	3/14/2005	10:00	F	1984.0	21.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False

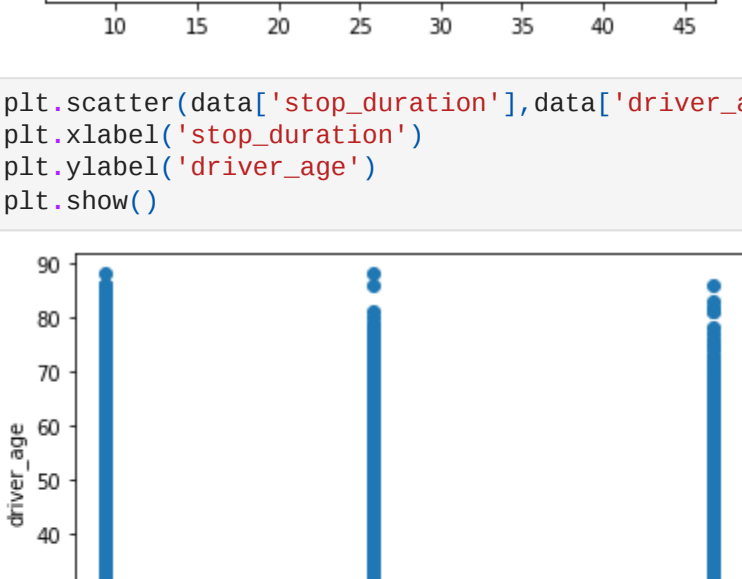
Scatterplot

```
In [38]: plt.scatter(data['stop_duration'],data['driver_age'])
```

```
Out[38]: <matplotlib.collections.PathCollection at 0x2131ffb40>
```

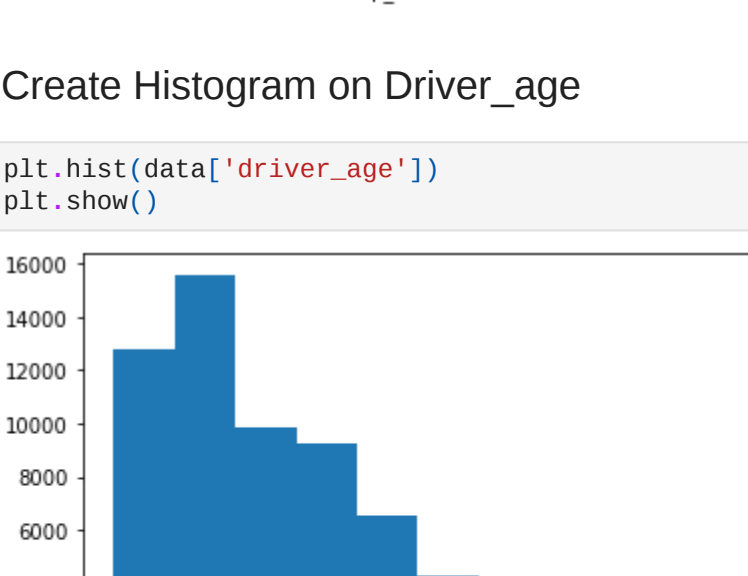


```
In [39]: plt.scatter(data['stop_duration'],data['driver_age'])
plt.xlabel('stop_duration')
plt.ylabel('driver_age')
plt.show()
```



Create Histogram on Driver_age

```
In [35]: plt.hist(data['driver_age'])
plt.show()
```



```
In [61]: data.head()
```

Out[61]:	stop_time	driver_gender	driver_age_raw	driver_age	driver_race	violation_raw	violation	search_conducted	search_type	stop_outcome	is_arrested	stop_duration	drugs_related_stop
0	1:55	M	1985.0	20.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
1	8:15	M	1965.0	40.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
2	23:15	M	1972.0	33.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False
3	17:15	M	1986.0	19.0	White	Call for Service	Other	False	NaN	Arrest Driver	True	24.0	False
4	10:00	F	1984.0	21.0	White	Speeding	Speeding	False	NaN	Citation	False	7.5	False

```
In [62]: data.shape
```

```
Out[62]: (65535, 13)
```

```
In [63]: import seaborn as sns
```

```
In [64]: plt.figure(figsize = (15 , 9))
sns.heatmap(data.corr(), annot = True)
```

```
Out[64]: <AxesSubplot:~>
```



```
In [92]: data.head()
```

Out[92]:	driver_age_raw	driver_age	violation_raw	violation	stop_outcome	is_arrested	stop_duration	drugs_related_stop
0	1985.0	20.0	Speeding	Speeding	Citation	False	7.5	False
1	1965.0	40.0	Speeding	Speeding	Citation	False	7.5	False
2	1972.0	33.0	Speeding	Speeding	Citation	False	7.5	False
3	1986.0	19.0	Call for Service	Other	Arrest Driver	True	24.0	False
4	1984.0	21.0	Speeding	Speeding	Citation	False	7.5	False

```
In [93]: data.drop(columns = 'violation',inplace = True)
```

```
In [94]: data
```

Out[94]:	driver_age_raw	driver_age	violation_raw	stop_outcome	is_arrested	stop_duration	drugs_related_stop
0	1985.0	20.0	Speeding	Citation	False	7.5	False
1	1965.0	40.0	Speeding	Citation	False	7.5	False
2	1972.0	33.0	Speeding	Citation	False	7.5	False
3	1986.0	19.0	Call for Service	Arrest Driver	True	24.0	False
4	1984.0	21.0	Speeding	Citation	False	7.5	False
...