

EDA-Exploitory Data Analysis -Drug Dataset

import libraries

```
In [1]: import pandas as pd
import numpy as np
import matplotlib as pp
```

import dataset

```
In [2]: data=pd.read_csv(r"E:\154\4_drug200.csv")
```

```
In [3]: display(data)
```

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	drugY
...
195	56	F	LOW	HIGH	11.567	drugC
196	16	M	LOW	HIGH	12.006	drugC
197	52	M	NORMAL	HIGH	9.894	drugX
198	23	M	NORMAL	NORMAL	14.020	drugX
199	40	F	LOW	NORMAL	11.349	drugX

To display top 10 rows

```
In [4]: data.head(10)
```

```
Out[4]:
```

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	drugY
5	22	F	NORMAL	HIGH	8.607	drugX
6	49	F	NORMAL	HIGH	16.275	drugY
7	41	M	LOW	HIGH	11.037	drugC
8	60	M	NORMAL	HIGH	15.171	drugY
9	43	M	LOW	NORMAL	19.368	drugY

```
In [ ]:
```

To display last 5 rows

```
In [5]: data.tail()
```

```
Out[5]:
```

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
195	56	F	LOW	HIGH	11.567	drugC
196	16	M	LOW	HIGH	12.006	drugC
197	52	M	NORMAL	HIGH	9.894	drugX
198	23	M	NORMAL	NORMAL	14.020	drugX
199	40	F	LOW	NORMAL	11.349	drugX

```
In [6]: data.dtypes
```

```
Out[6]: Age          int64  
Sex          object  
BP           object  
Cholesterol  object  
Na_to_K      float64  
Drug         object  
dtype: object
```

To view statistical summary

In [7]: `data.describe()`

Out[7]:

	Age	Na_to_K
count	200.000000	200.000000
mean	44.315000	16.084485
std	16.544315	7.223956
min	15.000000	6.269000
25%	31.000000	10.445500
50%	45.000000	13.936500
75%	58.000000	19.380000
max	74.000000	38.247000

To Print no of elements

In [8]: `data.size`

Out[8]: 1200

In [9]: `data.ndim`

Out[9]: 2

To print no of rows and columns

In [10]: `data.shape`

Out[10]: (200, 6)

To find missing values

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

Out[11]:

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
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
...
195	False	False	False	False	False	False
196	False	False	False	False	False	False
197	False	False	False	False	False	False
198	False	False	False	False	False	False
199	False	False	False	False	False	False

200 rows × 6 columns

To drop null values with constatns

In [12]: `data.fillna(5)`

Out[12]:

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	drugY
...
195	56	F	LOW	HIGH	11.567	drugC
196	16	M	LOW	HIGH	12.006	drugC
197	52	M	NORMAL	HIGH	9.894	drugX
198	23	M	NORMAL	NORMAL	14.020	drugX
199	40	F	LOW	NORMAL	11.349	drugX

200 rows × 6 columns

```
In [13]: data.dropna()
```

```
Out[13]:
```

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC
2	47	M	LOW	HIGH	10.114	drugC
3	28	F	NORMAL	HIGH	7.798	drugX
4	61	F	LOW	HIGH	18.043	drugY
...
195	56	F	LOW	HIGH	11.567	drugC
196	16	M	LOW	HIGH	12.006	drugC
197	52	M	NORMAL	HIGH	9.894	drugX
198	23	M	NORMAL	NORMAL	14.020	drugX
199	40	F	LOW	NORMAL	11.349	drugX

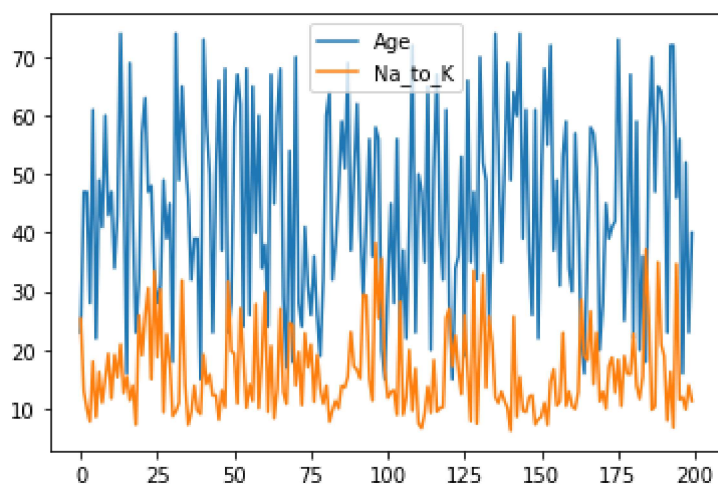
200 rows × 6 columns

Line Plot

Type *Markdown* and LaTeX: α^2

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

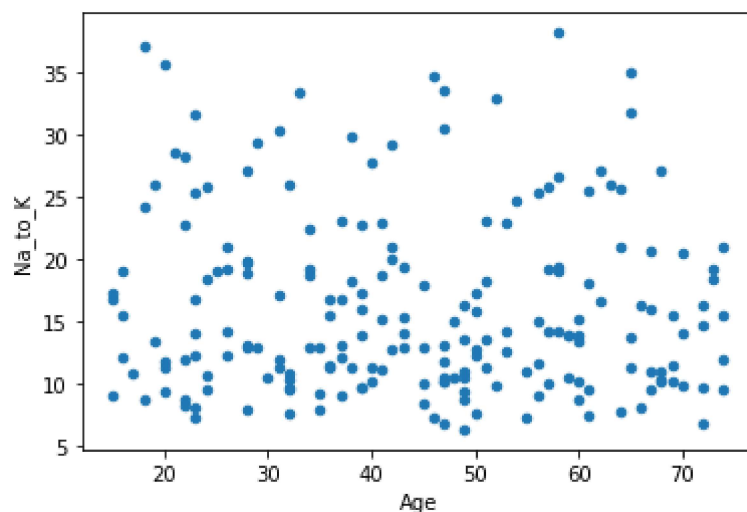
```
Out[14]: <AxesSubplot:>
```



Scatter Plot

```
In [15]: data.plot.scatter(x='Age',y='Na_to_K')
```

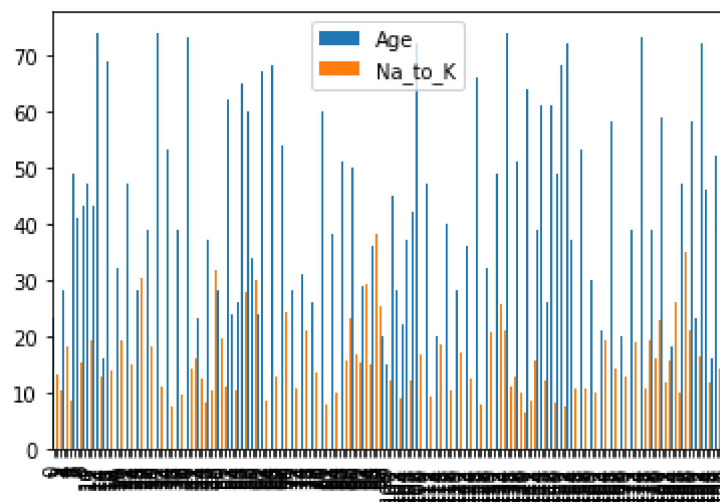
```
Out[15]: <AxesSubplot:xlabel='Age', ylabel='Na_to_K'>
```



Bar Chart

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

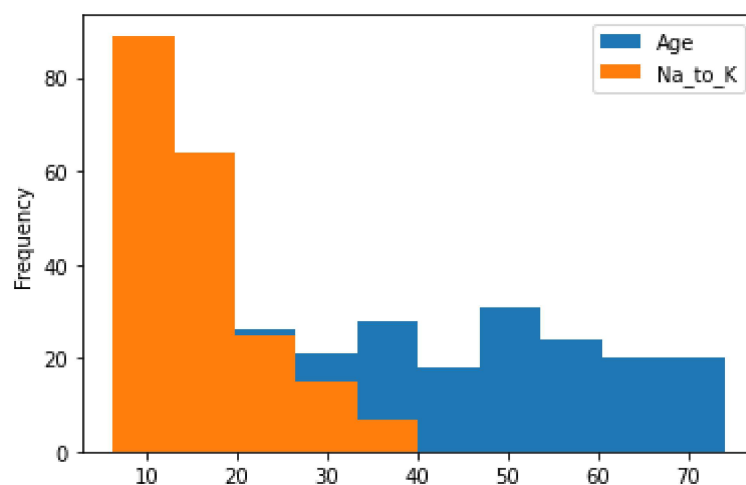
```
Out[16]: <AxesSubplot:>
```



Histogram

```
In [17]: data.plot.hist()
```

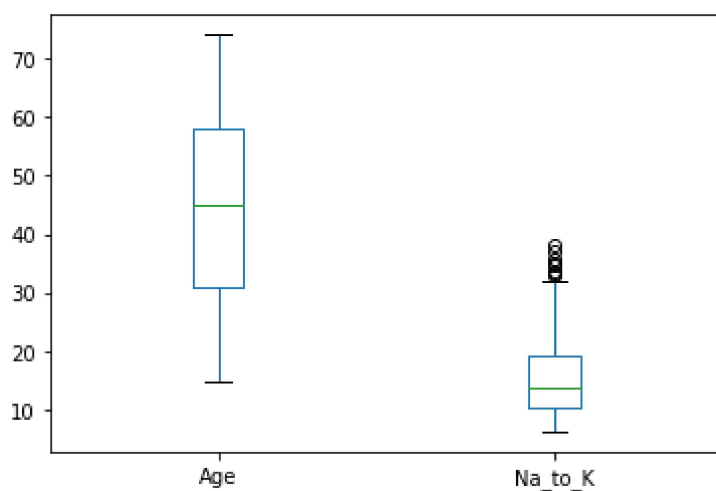
```
Out[17]: <AxesSubplot:ylabel='Frequency'>
```



Box Plot

```
In [18]: data.plot.box()
```

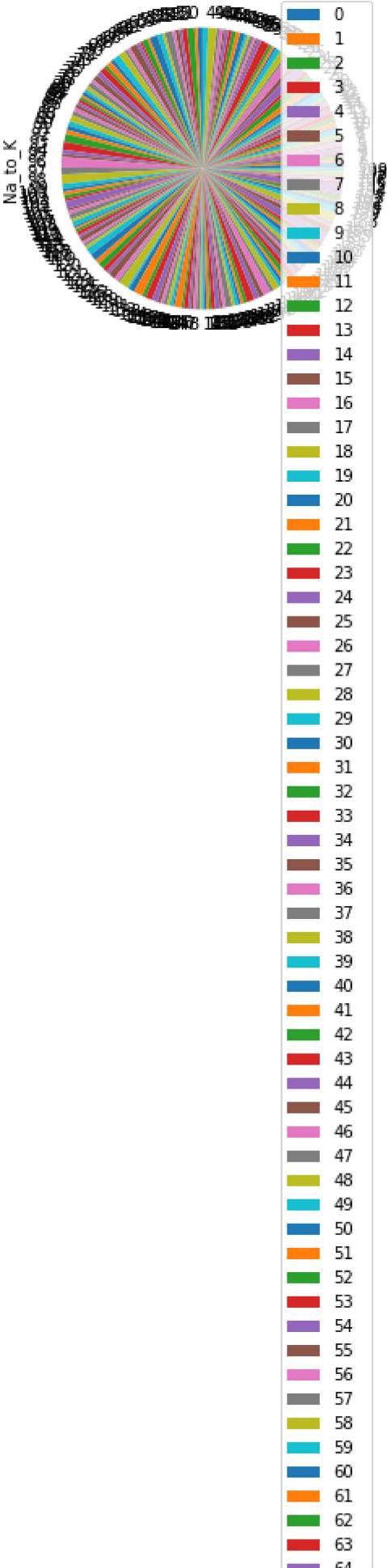
```
Out[18]: <AxesSubplot:>
```

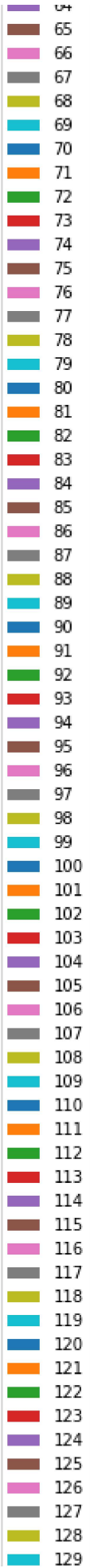


Pie Chart

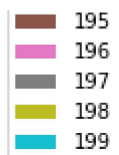
```
In [32]: data.plot.pie(y="Na_to_K")
```

```
Out[32]: <AxesSubplot:ylabel='Na_to_K'>
```



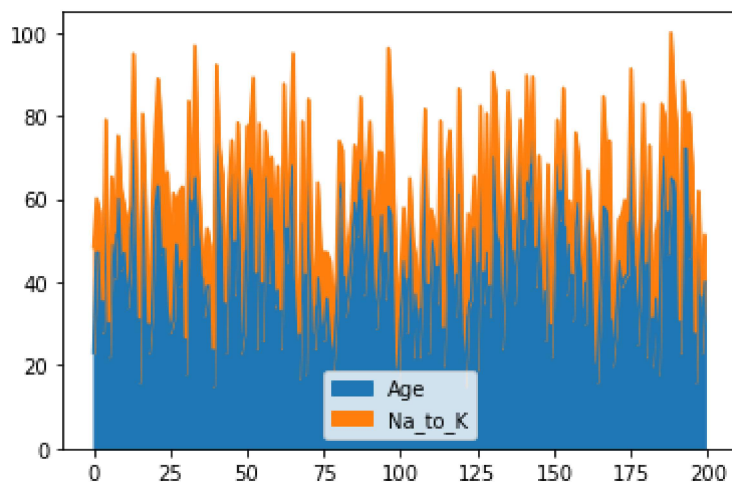
<div></div>	130
<div></div>	131
<div></div>	132
<div></div>	133
<div></div>	134
<div></div>	135
<div></div>	136
<div></div>	137
<div></div>	138
<div></div>	139
<div></div>	140
<div></div>	141
<div></div>	142
<div></div>	143
<div></div>	144
<div></div>	145
<div></div>	146
<div></div>	147
<div></div>	148
<div></div>	149
<div></div>	150
<div></div>	151
<div></div>	152
<div></div>	153
<div></div>	154
<div></div>	155
<div></div>	156
<div></div>	157
<div></div>	158
<div></div>	159
<div></div>	160
<div></div>	161
<div></div>	162
<div></div>	163
<div></div>	164
<div></div>	165
<div></div>	166
<div></div>	167
<div></div>	168
<div></div>	169
<div></div>	170
<div></div>	171
<div></div>	172
<div></div>	173
<div></div>	174
<div></div>	175
<div></div>	176
<div></div>	177
<div></div>	178
<div></div>	179
<div></div>	180
<div></div>	181
<div></div>	182
<div></div>	183
<div></div>	184
<div></div>	185
<div></div>	186
<div></div>	187
<div></div>	188
<div></div>	189
<div></div>	190
<div></div>	191
<div></div>	192
<div></div>	193
<div></div>	194



Area

```
In [20]: data.plot.area()
```

```
Out[20]: <AxesSubplot:>
```



To Find Mean

```
In [21]: data.mean()
```

```
Out[21]: Age          44.315000  
Na_to_K      16.084485  
dtype: float64
```

To Find Median

```
In [22]: data.median()
```

```
Out[22]: Age          45.00000  
Na_to_K      13.9365  
dtype: float64
```

To Find Mode

In [23]:

data.mode()

Out[23]:

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	47.0	M	HIGH	HIGH	12.006	drugY
1	NaN	NaN	NaN	NaN	18.295	NaN

Describe

In [24]:

data.describe()

Out[24]:

	Age	Na_to_K
count	200.000000	200.000000
mean	44.315000	16.084485
std	16.544315	7.223956
min	15.000000	6.269000
25%	31.000000	10.445500
50%	45.000000	13.936500
75%	58.000000	19.380000
max	74.000000	38.247000

Sum

In [25]:

data.sum()

Out[25]:

Age	8863
Sex	FMMFFFFMMMFMMFFMMFMFFMMFFMMFMFFMMFFMMFFMMFF...
BP	HIGHLOWLOWNORMALLOWNORMALNORMALLOWNORMALLOWLOW...
Cholesterol	HIGHHHIGHHHIGHHHIGHHHIGHHHIGHHHIGHHHIGHHHIGHNORMALHIGH...
Na_to_K	3216.897
Drug	drugYdrugCdrugCdrugXdrugYdrugXdrugYdrugCdrugYd...
dtype:	object

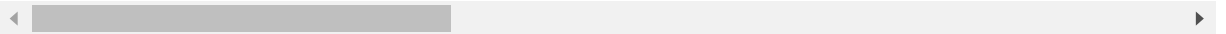
Cumulative Sum

In [26]: `data.cumsum()`

Out[26]:

	Age	Sex
0	23	F
1	70	FM
2	117	FMM
3	145	FMMF
4	206	FMMFF
...
195	8732	FMMFFFFMMMFMMFFMMFFMMFFMFMMFMMMMFMFFMMFF... HIGHLOWLOWNOF
196	8748	FMMFFFFMMMFMMFFMMFFMMFFMFMMFMMMMFMFFMMFF... HIGHLOWLOWNOF
197	8800	FMMFFFFMMMFMMFFMMFFMMFFMFMMFMMMMFMFFMMFF... HIGHLOWLOWNOF
198	8823	FMMFFFFMMMFMMFFMMFFMMFFMFMMFMMMMFMFFMMFF... HIGHLOWLOWNOF
199	8863	FMMFFFFMMMFMMFFMMFFMMFFMFMMFMMMMFMFFMMFF... HIGHLOWLOWNOF

200 rows × 6 columns



Minimum Values

In [27]: `data.min()`

Out[27]:

Age	15
Sex	F
BP	HIGH
Cholesterol	HIGH
Na_to_K	6.269
Drug	drugA

dtype: object

Maximum Values

In [28]: `data.max()`

Out[28]:

Age	74
Sex	M
BP	NORMAL
Cholesterol	NORMAL
Na_to_K	38.247
Drug	drugY

dtype: object

Correlation

```
In [29]: from scipy.stats import spearmanr  
print(spearmanr(data['Age'],data['Na_to_K']))
```

SpearmanrResult(correlation=-0.047273882688479915, pvalue=0.5062200581387418)

Covariance

```
In [30]: from scipy.stats import pearsonr  
print(pearsonr(data['Age'],data['Na_to_K']))
```

(-0.06311949726772592, 0.3745756399034559)

Count

```
In [31]: data.count()
```

```
Out[31]: Age           200  
Sex           200  
BP            200  
Cholesterol    200  
Na_to_K        200  
Drug           200  
dtype: int64
```