

# SUMESH R - 20104169

## Basic Analysis using NumPy and Pandas

### Import Libraries

```
In [1]: import pandas as pd
```

```
In [2]: import numpy as np
```

### Import Dataset

```
In [3]: data = pd.read_csv("4_drug200.csv")
```

```
In [4]: 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

200 rows × 6 columns

### To display top 10 rows

```
In [5]: data.head(10)
```

```
Out[5]:
```

	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

## to display last 5 rows

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

```
Out[6]:
```

	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

## 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 number of elements

```
In [8]: data.size
```

Out[8]: 1200

# to print number of row and cols

```
In [9]: data.shape
```

Out[9]: (200, 6)

# to find missing values

```
In [10]: data.isna()
```

Out[10]:

	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

# fill null values with a constant

```
In [11]: data.fillna(5)
```

Out[11]:

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
0	23	F	HIGH	HIGH	25.355	drugY
1	47	M	LOW	HIGH	13.093	drugC

	Age	Sex	BP	Cholesterol	Na_to_K	Drug
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

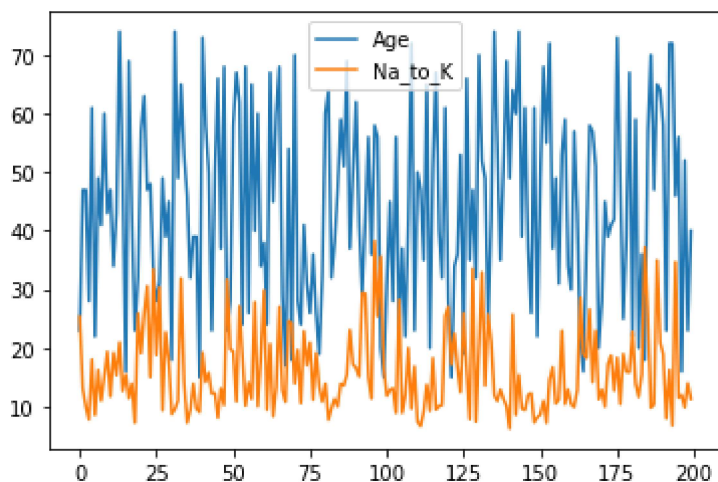
## to select a particular columns

```
In [12]: df=pd.DataFrame(data[['Age','Na_to_K']])
import matplotlib.pyplot as plt
```

## line plot

```
In [13]: df.plot.line()
```

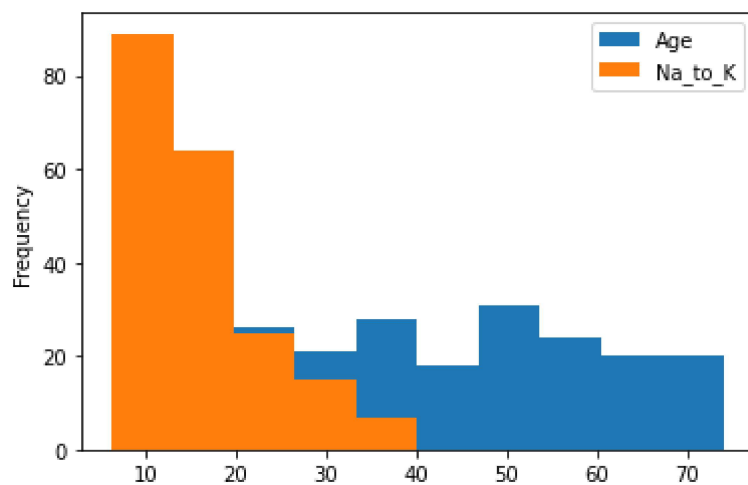
Out[13]: <AxesSubplot:>



## histogram

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

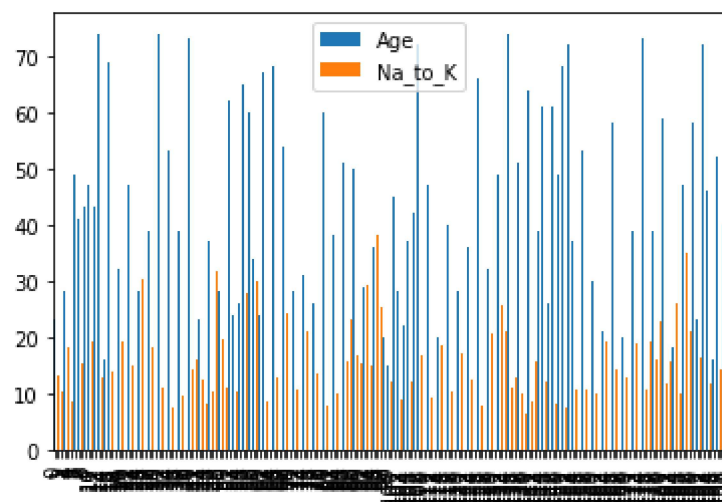
Out[14]: <AxesSubplot:ylabel='Frequency'>



## bar chart

```
In [15]: df.plot.bar()
```

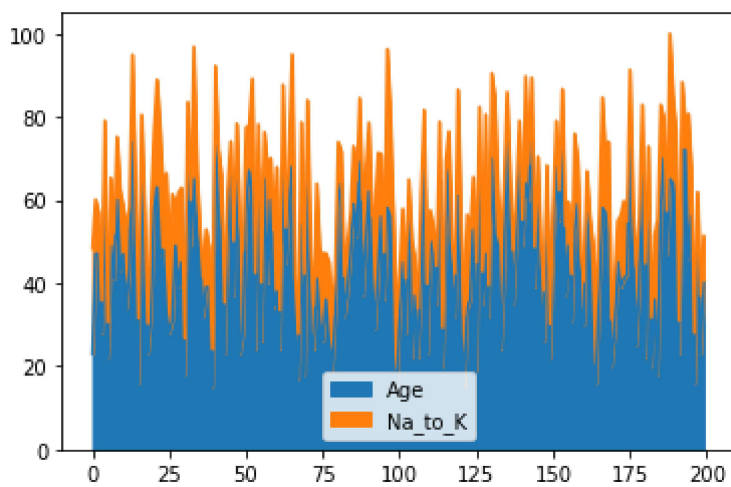
```
Out[15]: <AxesSubplot:>
```



## area plot

```
In [16]: df.plot.area()
```

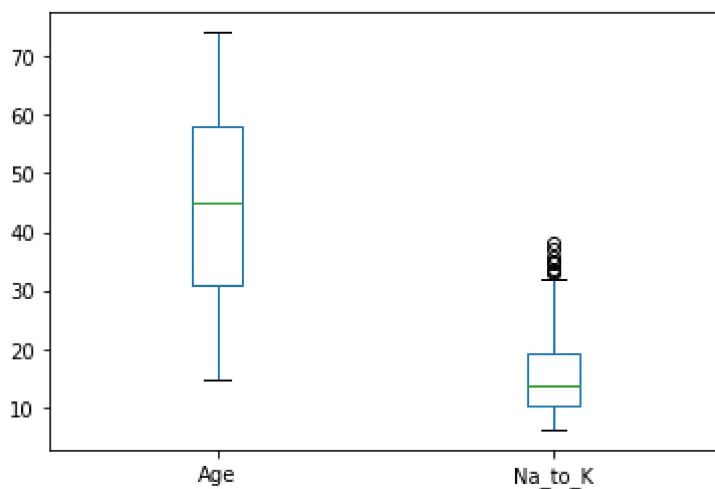
```
Out[16]: <AxesSubplot:>
```



## box plot

```
In [17]: df.plot.box()
```

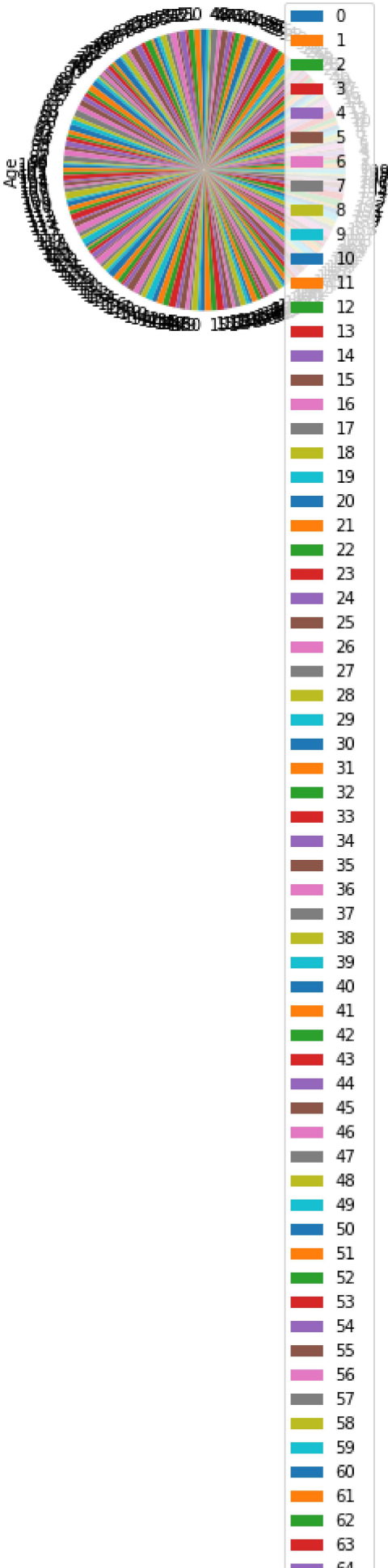
```
Out[17]: <AxesSubplot:>
```

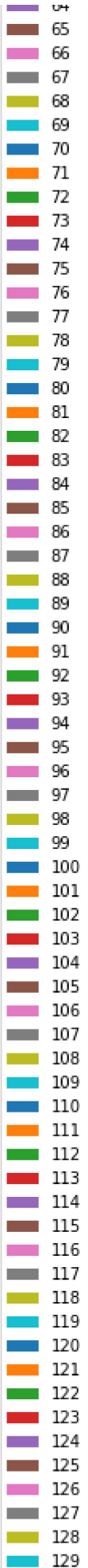


## pie plot

```
In [18]: df.plot.pie(y="Age")
```

```
Out[18]: <AxesSubplot:ylabel='Age'>
```







# scatter plot

- 130
- 131
- 132
- 133

```
In [19]: df.plot.scatter(x="Age",y="Na_to_K")
```

Out[19]: <AxesSubplot:xlabel='Age', ylabel='Na\_to\_K'>

