

DATA SPECIALIZATION

#Name : Shreya Sharma #Roll no. : 46 #Sectin : 3B #Date : 20/07/2024

In [11]: #Aim : To Perform Data specialization/statical description on data

In [1]: import pandas as pd

In [3]: import os

In [5]: os.getcwd()

Out[5]: 'C:\\Users\\pravi'

In [7]: os.chdir('C:\\Users\\pravi\\Desktop')

In [9]: df=pd.read_csv("framingham.csv")

In [13]: df.head()

Out[13]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp	diabetes	totChol	sysBP	diaBP	
0	1	39	4.0	0	0.0	0.0	0	0	0	195.0	106.0	70.0	26
1	0	46	2.0	0	0.0	0.0	0	0	0	250.0	121.0	81.0	26
2	1	48	1.0	1	20.0	0.0	0	0	0	245.0	127.5	80.0	26
3	0	61	3.0	1	30.0	0.0	0	1	0	225.0	150.0	95.0	26
4	0	46	3.0	1	23.0	0.0	0	0	0	285.0	130.0	84.0	26

In [15]: df.head(100)

Out[15]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp	diabetes	totChol	sysBP	diaBP	
0	1	39	4.0	0	0.0	0.0	0	0	0	195.0	106.0	70.0	26
1	0	46	2.0	0	0.0	0.0	0	0	0	250.0	121.0	81.0	26
2	1	48	1.0	1	20.0	0.0	0	0	0	245.0	127.5	80.0	26
3	0	61	3.0	1	30.0	0.0	0	1	0	225.0	150.0	95.0	26
4	0	46	3.0	1	23.0	0.0	0	0	0	285.0	130.0	84.0	26
...
95	0	65	3.0	0	0.0	0.0	0	0	0	193.0	123.0	76.5	26
96	0	63	4.0	1	20.0	0.0	0	0	1	239.0	134.0	80.0	26
97	0	40	2.0	0	0.0	0.0	0	0	0	205.0	100.0	60.0	26
98	0	56	1.0	0	0.0	0.0	0	1	0	296.0	180.0	90.0	26
99	0	56	1.0	1	15.0	0.0	0	0	0	269.0	121.0	75.0	26

100 rows × 16 columns

In [17]: df.tail()

Out[17]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp	diabetes	totChol	sysBP	diaBP	
4233	1	50	1.0	1	1.0	0.0	0	1	0	313.0	179.0	92.0	26
4234	1	51	3.0	1	43.0	0.0	0	0	0	207.0	126.5	80.0	26
4235	0	48	2.0	1	20.0	NaN	0	0	0	248.0	131.0	72.0	26
4236	0	44	1.0	1	15.0	0.0	0	0	0	210.0	126.5	87.0	26
4237	0	52	2.0	0	0.0	0.0	0	0	0	269.0	133.5	83.0	26

In [19]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4238 entries, 0 to 4237
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   male                  4238 non-null   int64
1   age                   4238 non-null   int64
2   education             4133 non-null   float64
3   currentSmoker        4238 non-null   int64
4   cigsPerDay           4209 non-null   float64
5   BPMeds               4185 non-null   float64
6   prevalentStroke      4238 non-null   int64
7   prevalentHyp         4238 non-null   int64
8   diabetes              4238 non-null   int64
9   totChol              4188 non-null   float64
10  sysBP                4238 non-null   float64
11  diaBP                4238 non-null   float64
12  BMI                  4219 non-null   float64
13  heartRate            4237 non-null   float64
14  glucose              3850 non-null   float64
15  TenYearCHD           4238 non-null   int64
dtypes: float64(9), int64(7)
memory usage: 529.9 KB
```

In [21]: df.shape

Out[21]: (4238, 16)

In [23]: df.size

Out[23]: 67808

In [25]: df.ndim

Out[25]: 2

In [27]: df.tail(10)

Out[27]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp	diabetes	totChol	sysBP	diaBP
4228	0	50	1.0	0	0.0	0.0	0	1	1	260.0	190.0	130.0
4229	0	51	3.0	1	20.0	0.0	0	1	0	251.0	140.0	80.0
4230	0	56	1.0	1	3.0	0.0	0	1	0	268.0	170.0	102.0
4231	1	58	3.0	0	0.0	0.0	0	1	0	187.0	141.0	81.0
4232	1	68	1.0	0	0.0	0.0	0	1	0	176.0	168.0	97.0
4233	1	50	1.0	1	1.0	0.0	0	1	0	313.0	179.0	92.0
4234	1	51	3.0	1	43.0	0.0	0	0	0	207.0	126.5	80.0
4235	0	48	2.0	1	20.0	NaN	0	0	0	248.0	131.0	72.0
4236	0	44	1.0	1	15.0	0.0	0	0	0	210.0	126.5	87.0
4237	0	52	2.0	0	0.0	0.0	0	0	0	269.0	133.5	83.0

In [29]: df.describe()

Out[29]:

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp	diabetes
count	4238.000000	4238.000000	4133.000000	4238.000000	4209.000000	4185.000000	4238.000000	4238.000000	4238.000000
mean	0.429212	49.584946	1.978950	0.494101	9.003089	0.029630	0.005899	0.310524	0.025722
std	0.495022	8.572160	1.019791	0.500024	11.920094	0.169584	0.076587	0.462763	0.158311
min	0.000000	32.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.000000	42.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	0.000000	49.000000	2.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	1.000000	56.000000	3.000000	1.000000	20.000000	0.000000	0.000000	1.000000	0.000000
max	1.000000	70.000000	4.000000	1.000000	70.000000	1.000000	1.000000	1.000000	1.000000

In []: