

1. Mean, Median & Mode (MMM) Table Inference

	sl_no	ssc_p	hsc_p	degree_p	etest_p	mba_p	salary
Mean	108	67.3034	66.3332	66.3702	72.1006	62.2782	288655
Median	108	67	65	66	71	62	265000
Mode	1	62	63	65	60	56.7	300000

Inference Tabulation:

Index	Mean	Median	Outlier Status	Mode
Ssc_p	The average ssc pass mark of 215 students = 67.3034	The average ssc pass mark of 215 students excluding outliers = 67	Not present	Most repeated pass mark secured in ssc among 215 students= 62
Take away	Average performance in Ssc		-	-
Hsc_p	The average hsc pass mark of 215 students = 66.3332	The average hsc pass mark of 215 students excluding outliers = 65	Present (Lesser & Greater)	Most repeated pass mark secured in hsc among 215 students= 63
Take away	Average performance in Hsc		-	-
Degree_p	The average degree pass mark of 215 students= 66.3702	The average degree pass mark of 215 students excluding outliers = 66	Present (Greater)	Most repeated pass mark secured in degree among 215 students= 65
Take away	Average performance in degree graduation		-	-
Etest_p	The average e-test pass mark of 215 students = 72.1006	The average e-test pass mark of 215 students excluding outliers = 71	Not Present	Most repeated pass mark secured in e-test among 215 students= 60
Take away	Above average performance in Etest		-	-
Mba_p	The average MBA pass mark of 215 students = 62.2782	The average MBA pass mark of 215 students excluding outliers = 62	Not present	Most repeated pass mark secured in MBA among 215 students= 56.7
Take away	Average performance in MBA		-	-
Salary	The average salary of 215 students = RS.2,88,655	The average salary of 215 students excluding outliers= RS.2,65,000	Present (Greater)	Most repeated salary package received among 215 students= RS.3,00,000
Take away	Among 215 students, most of the students receive RS.3,00,000 as their salary.			

முத்து வசமதி

2. Formula Reasoning:

$$IQR = Q3 - Q1$$

Lesser Outlier

Less Than

Outlier range = $Q1 - 1.5 * IQR$

Greater Outlier

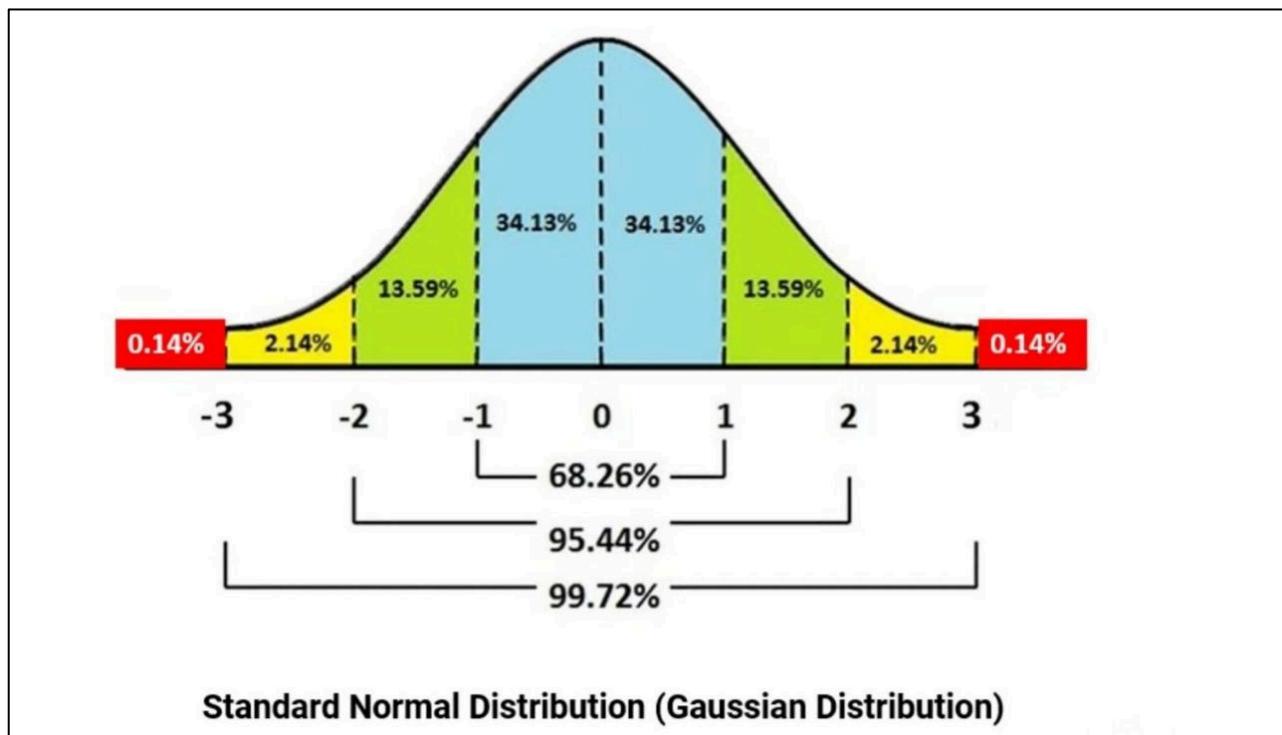
Greater Than

Outlier range = $Q3 + 1.5 * IQR$

Question : Why 1.5 scale is multiplied with IQR for outliers?

Reason:

Consider the Gaussian distribution decision range for outlier detection as shown below.



Here we are looking for the data which is useful for us which lies within the range mean plus or minus 3 sigma data points. Those data points which lies beyond this range is considered as outlier.

To achieve this let's try to calculate the approximate scale value which has to be multiplied with IQR so that we end up with $\pm 3\sigma$ range.

Consider Scale = 1

$$\text{Lower Bound} = Q1 - 1 * \text{IQR}$$

$$= Q1 - 1 * (Q3 - Q1)$$

$$= -0.675\sigma - 1 * (0.675 - [-0.675])\sigma$$

$$= -0.675\sigma - 1 * 1.35\sigma$$

$$= -2.025\sigma$$

$$\text{Upper Bound} = Q3 + 1 * \text{IQR}$$

$$= Q3 + 1 * (Q3 - Q1)$$

$$= 0.675\sigma + 1 * (0.675 - [-0.675])\sigma$$

$$= 0.675\sigma + 1 * 1.35\sigma$$

$$= 2.025\sigma$$

Conclusion:

If the scale is taken as 1, then according to IQR Method any data which lies beyond 2.025σ from the mean (μ), on either side, shall be considered an outlier. But as we know, **up to 3σ , on either side of the μ , the data is useful**. So, we cannot take scale = 1, because the decision range gets so small (compared to 3σ) that it considers some data points as outliers, which is **not desirable**

Consider Scale = 1.5

$$\text{Lower Bound} = Q1 - 1.5 * \text{IQR}$$

$$= Q1 - 1.5 * (Q3 - Q1)$$

$$= -0.675\sigma - 1.5 * (0.675 - [-0.675])\sigma$$

$$= -0.675\sigma - 1.5 * 1.35\sigma$$

$$= -2.7\sigma$$

$$\text{Upper Bound} = Q3 + 1.5 * \text{IQR}$$

$$= Q3 + 1.5 * (Q3 - Q1)$$

$$= 0.675\sigma + 1.5 * (0.675 - [-0.675])\sigma$$

$$= 0.675\sigma + 1.5 * 1.35\sigma$$

$$= 2.7\sigma$$

Conclusion:

If the scale is taken as 1.5, then according to IQR Method any data which lies beyond 2.7σ from the mean (μ), on either side, shall be considered an outlier. But as we know, **up to 3σ , on either side of the μ , the data is useful**. So, if we take scale = 1.5, this makes the decision range closest to what Gaussian Distribution considers for outlier detection, and this is exactly what we wanted. So, it is **desirable**.

Consider Scale = 1.7

$$\text{Lower Bound} = Q1 - 1.7 * \text{IQR}$$

$$= Q1 - 1.7 * (Q3 - Q1)$$

$$= -0.675\sigma - 1.7 * (0.675 - [-0.675])\sigma$$

$$= -0.675\sigma - 1.7 * 1.35\sigma$$

$$= -2.97\sigma$$

$$\text{Upper Bound} = Q3 + 1.7 * \text{IQR}$$

$$= Q3 + 1.7 * (Q3 - Q1)$$

$$= 0.675\sigma + 1.7 * (0.675 - [-0.675])\sigma$$

$$= 0.675\sigma + 1.7 * 1.35\sigma$$

$$= 2.97\sigma$$

Conclusion:

If the scale is taken as 1.7, then according to IQR Method any data which lies beyond 2.97σ from the mean (μ), on either side, shall be considered an outlier. But as we know, up to 3σ , on either side of the μ , the data is useful. So, if we take scale = 1.7, this makes the decision range the same as that of Gaussian Distribution considers for outlier detection. So, it is desirable again. But IQR is not the perfect method to detect outliers.

To keep the tolerance we choose 1.5 scale on safer side though 1.7 scale suits the best.

Consider Scale = 2

$$\text{Lower Bound} = Q1 - 2 * \text{IQR}$$

$$= Q1 - 2 * (Q3 - Q1)$$

$$= -0.675\sigma - 2 * (0.675 - [-0.675])\sigma$$

$$= -0.675\sigma - 2 * 1.35\sigma$$

$$= -3.375\sigma$$

$$\text{Upper Bound} = Q3 + 2 * \text{IQR}$$

$$= Q3 + 2 * (Q3 - Q1)$$

$$= 0.675\sigma + 2 * (0.675 - [-0.675])\sigma$$

$$= 0.675\sigma + 2 * 1.35\sigma$$

$$= 3.375\sigma$$

Conclusion:

If the scale is taken as 2, then according to IQR Method any data which lies beyond 3.375σ from the mean (μ), on either side, shall be considered an outlier. But as we know, up to 3σ , on either side of the μ , the data is useful. So, we cannot take scale = 2, because the decision range gets so big (compared to 3σ) that it considers some outliers as data points, which is not desirable again.



3. Problem on IQR

Question :

- a. The interquartile range. Compare the two interquartile ranges.
- b. Any outliers in either set.

The five number summary for the day and night classes is

	Minimum	Q_1	Median	Q_3	Maximum
Day	32	56	74.5	82.5	99
Night	25.5	78	81	89	98

Solution :

Classes	Q_1	Q_3	$IQR = Q_3 - Q_1$	Lesser outlier range = $Q_1 - 1.5 * IQR$	Greater outlier range = $Q_3 + 1.5 * IQR$	Upper & Lower Range
Day	56	82.5	26.5	16.25	122.25	Range:(16.25, 122.25) No values in the row named day is outside the range. So no outliers are present.
Night	78	89	11	61.5	105.5	Range:(61.5, 105.5) Value 25.5 in the row named night is outside the range. So outlier is present. Type of outlier : Low range outlier

4.Kurtosis & Skewness Table Inference

	sl_no	ssc_p	hsc_p	degree_p	etest_p	mba_p	salary
kurtosis	-1.2	-0.60751	0.0869008	-0.0974897	-1.08858	-0.470723	-0.239837
skew	0	-0.132649	0.162611	0.204164	0.282308	0.313576	0.8067

Inference:

Kurtosis: (Tells about the Tailness)

Kurtosis value for all the columns starting from ssc_p to salary is <3 which means the type of kurtosis is platykurtic(Negative Excess Kurtosis).

Which in turn gives the information that the number of students falling under below average & above average (to some extent) whereas number of students falling under average (to better extent) in terms of exam performance & receiving salary, where the most population of students fall under average criteria.

Skewness: (Tells us about Asymmetry from the normal distribution)

Ssc_p alone exhibits negative skewness.i.e, Above average scorers are more in count when compared to below average scorers.

Hsc_p, Degree_p, Etest_p & MBA _p variables exhibits positive skewness. i.e Above average scorers are lesser in count when compared to below average scorers.

Salary variable also exhibits positive skewness. i.e Above average salary receivers are lesser in count when compared to below average salary receivers.

முத்து வசமதி

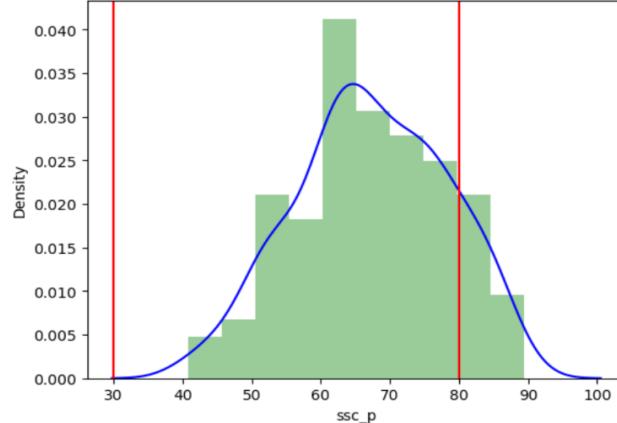
5. Probability Distribution Function (PDF) Code Explanation:

```
In [1]: 1 #Generating a function for PDF such that getting the arguments from the user
2 def get_pdf_probability(dataset,startrange,endrange):
3
4     #Importing necessary libraries for visual representation
5     from matplotlib import pyplot
6     from scipy.stats import norm
7     import seaborn as sns
8
9
10    # Plotting distribution plot with density curve in blue color & histogram in Green colour
11    ax = sns.distplot(dataset,kde=True,kde_kws={'color':'blue'},color='Green')
12
13
14    # Making initial & end range axis vertical line visible
15    pyplot.axvline(startrange,color='Red')
16    pyplot.axvline(endrange,color='Red')
17
18
19    # Generating a sample
20    sample = dataset
21
22
23    # Calculating parameters like mean & std and displaying the calculated values using print statement
24    sample_mean =sample.mean()
25    sample_std = sample.std()
26    print('Mean=% .3f, Standard Deviation=% .3f' % (sample_mean, sample_std))
27
28
29    # Defining the distribution by giving mean & std as input
30    dist = norm(sample_mean, sample_std)
31
32    # Sample probabilities for a range of outcomes
33    values = [value for value in range(startrange, endrange)]
34    probabilities = [dist.pdf(value) for value in values]
35
36    # For whole probability
37    prob= sum(probabilities)
38    print("The area between range({},{}) : {}".format(startrange,endrange,sum(probabilities)))
39
40    return prob
```

Graph

```
In [8]: 1 get_pdf_probability(dataset["ssc_p"],30,80)
Mean=67.303, Standard Deviation=10.827
The area between range(30,80):0.8698639973987944
/var/folders/07/ykqp85052b1lh5kz22ghn8l40000gn/T/ipykernel_43901/2842244316.py:5: Us
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
ax = sns.distplot(dataset,kde=True,kde_kws={'color':'blue'},color='Green')
```

Out[8]: 0.8698639973987944

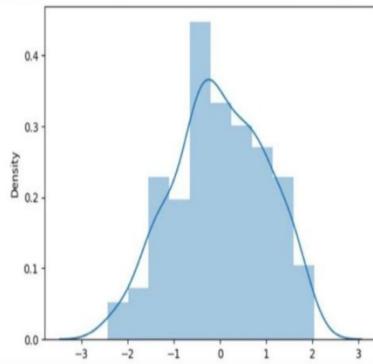


6. Standard Normal Distribution (SND) Code Explanation:

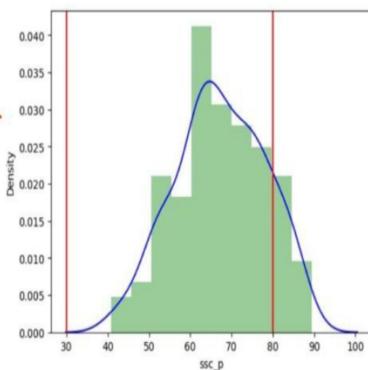
```
In [1]: 1 # Creating function for SND graph
2 def stdNBgraph(dataset):
3
4
5     # Converted to standard Normal Distribution
6     #Importing necessary library
7     import seaborn as sns
8
9     # Calculating parameters like mean & std the required coloumns
10    mean=dataset.mean()
11    std=dataset.std()
12
13
14    #Converting the required column of the dataset into list
15    values=[i for i in dataset]
16
17
18    #Calculating Z Score for the Values obtained in the previous step and storing in the list
19    z_score=[((j-mean)/std) for j in values]
20
21
22    #Feeding the Z score value into the distance plot for creating visual representation
23    sns.distplot(z_score,kde=True)
24
25    x = sum(z_score)/len(z_score)
26    print(x)
27    #z_score.std()
28
```

Graph Comparison (SND & PDF)

```
In [21]: 1 stdNBgraph(dataset["ssc_p"])
C:\Anaconda3\envs\aiml\lib\site-packages\ipykernel_launcher.py:11: UserWarning:
'distplot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either 'distplot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
# This is added back by InteractiveShellApp.init_path()
```



```
In [8]: 1 get_pdf_probability(dataset["ssc_p"],30,80)
Mean=67.303, Standard Deviation=10.827
The area between range[30,80]:0.8698639973987944
/var/folders/07/ykgp85052b1h5kz22ghn8l40000gn/T/ipykernel_43901/2842244316.py:5: UserWarning:
'distplot' is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either 'distplot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
ax = sns.distplot(dataset,kde=True,kde_kws={'color':'blue'},color='Green')
Out[8]: 0.8698639973987944
```



முத்து வசமதி

7. For Loop in Traditional Way

One Liner For Loop:

```
# sample probabilities for a range of outcomes
values = [value for value in range(startrange, endrange)]
probabilities = [dist.pdf(value) for value in values]
```

Tradition Way For Loop:

```
#values = [value for value in range(startrange, endrange)]
values=[]
for value in range(startrange, endrange):
    values.append(value)

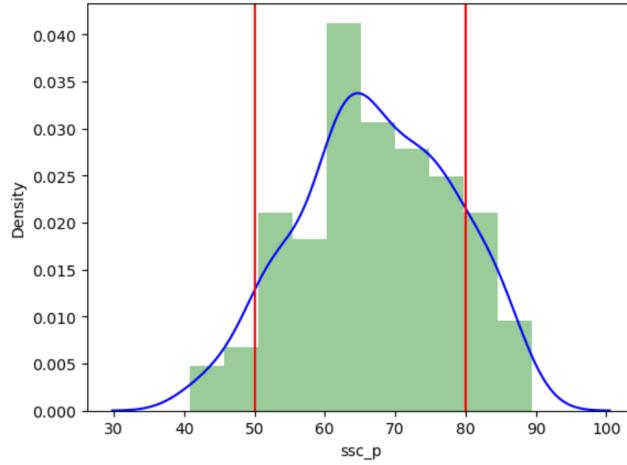
#probabilities = [dist.pdf(value) for value in values]
probabilities = []
for value in values:
    probabilities.append(dist.pdf(value))
```

Proof:

```
In [13]: 1 def get_pdf_probability(dataset,startrange, endrange):
2     from matplotlib import pyplot
3     from scipy.stats import norm
4     import seaborn as sns
5     ax = sns.distplot(dataset,kde=True,kde_kws={'color':'blue'},color='Green')
6     pyplot.axvline(startrange,color='Red')
7     pyplot.axvline(endrange,color='Red')
8     # generate a sample
9     sample = dataset
10    # calculate parameters
11    sample_mean =sample.mean()
12    sample_std = sample.std()
13    print('Mean=% .3f, Standard Deviation=% .3f' % (sample_mean, sample_std))
14    # define the distribution
15    dist = norm(sample_mean, sample_std)
16
17    # sample probabilities for a range of outcomes
18
19    #values = [value for value in range(startrange, endrange)]
20
21    values=[]
22    for value in range(startrange, endrange):
23        values.append(value)
24
25    #probabilities = [dist.pdf(value) for value in values]
26
27    probabilities = []
28    for value in values:
29        probabilities.append(dist.pdf(value))
30
31    prob=sum(probabilities)
32    print("The area between range({},{}) : {}".format(startrange,endrange,sum(probabilities)))
33
34    return prob
```

```
In [14]: 1 get_pdf_probability(dataset["ssc_p"],50,80)
Mean=67.303, Standard Deviation=10.827
The area between range(50,80): 0.8201083468562449
/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_98036/2713336809.py:5: UserWarning:
'distplot` is a deprecated function and will be removed in seaborn v0.14.0.
Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751
ax = sns.distplot(dataset,kde=True,kde_kws={'color':'blue'},color='Green')
```

Out[14]: 0.8201083468562449



முத்து வசமதி

8.CKD File Preprocessing

jupyter Preprocessed Kidney Disease File Last Checkpoint: 6 minutes ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel) O

In [1]:

```
1 import numpy as np
2 import pandas as pd
```

In [2]:

```
1 dataset=pd.read_csv("kidney_disease.csv")
2 dataset
```

Out[2]:

	id	age	F	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appet	pe	ane	classification
0	0	48.0	80.0	1.020	1.0	0.0	NaN	normal	notpresent	notpresent	...	44.0	7800.0	5.2	yes	yes	no	good	no	no	ckd
1	1	7.0	50.0	1.020	4.0	0.0	NaN	normal	notpresent	notpresent	...	38.0	6000.0	NaN	no	no	no	good	no	no	ckd
2	2	62.0	80.0	1.010	2.0	3.0	normal	normal	notpresent	notpresent	...	31.0	7500.0	NaN	no	yes	no	poor	no	yes	ckd
3	3	48.0	70.0	1.005	4.0	0.0	normal	abnormal	present	notpresent	...	32.0	6700.0	3.9	yes	no	no	poor	yes	yes	ckd
4	4	51.0	80.0	1.010	2.0	0.0	normal	normal	notpresent	notpresent	...	35.0	7300.0	4.6	no	no	no	good	no	no	ckd
...	
395	395	55.0	80.0	1.020	0.0	0.0	normal	normal	notpresent	notpresent	...	47.0	6700.0	4.9	no	no	no	good	no	no	notckd
396	396	42.0	70.0	1.025	0.0	0.0	normal	normal	notpresent	notpresent	...	54.0	7800.0	6.2	no	no	no	good	no	no	notckd
397	397	12.0	80.0	1.020	0.0	0.0	normal	normal	notpresent	notpresent	...	49.0	6600.0	5.4	no	no	no	good	no	no	notckd
398	398	17.0	60.0	1.025	0.0	0.0	normal	normal	notpresent	notpresent	...	51.0	7200.0	5.9	no	no	no	good	no	no	notckd
399	399	58.0	80.0	1.025	0.0	0.0	normal	normal	notpresent	notpresent	...	53.0	6800.0	6.1	no	no	no	good	no	no	notckd

400 rows x 26 columns

In [3]:

```
1 dataset.info()
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 26 columns):
 # Column Non-Null Count Dtype
--- ---
 0 id 400 non-null int64
 1 age 391 non-null float64
 2 F 388 non-null float64
 3 sg 353 non-null float64
 4 al 354 non-null float64
 5 su 351 non-null float64
 6 rbc 248 non-null object
 7 pc 335 non-null object
 8 pcc 396 non-null object
 9 ba 396 non-null object
 10 bgr 356 non-null float64
 11 bu 381 non-null float64
 12 sc 383 non-null float64
 13 sod 313 non-null float64
 14 pot 312 non-null float64
 15 hemo 348 non-null float64
 16 pcv 329 non-null float64
 17 wc 294 non-null float64
 18 rc 269 non-null float64
 19 htn 398 non-null object
 20 dm 398 non-null object
 21 cad 398 non-null object
 22 appet 399 non-null object
 23 pe 399 non-null object
 24 ane 399 non-null object
 25 classification 400 non-null object
dtypes: float64(14), int64(1), object(11)
memory usage: 81.4+ KB

முத்து வசமதி

```
In [4]: 1 dataset.isnull().sum().T  
2 #print(a,end='')
```

```
Out[4]: id          0  
age         9  
F          12  
sg         47  
al         46  
su         49  
rbc        152  
pc          65  
pcc         4  
ba          4  
bgr        44  
bu         19  
sc          17  
sod         87  
pot         88  
hemo        52  
pcv         71  
wc          106  
rc         131  
htn         2  
dm          2  
cad         2  
appet       1  
pe          1  
ane         1  
classification  0  
dtype: int64
```

Deleting Certain Rows with Missing Values

```
In [5]: 1 dataset.dropna(subset=['age','pcc', 'ba','htn','dm','cad','appet','pe','ane'], inplace=True)  
2 dataset
```

```
Out[5]:   id  age   F   sg   al   su   rbc     pc     pcc     ba ...   pcv     wc   rc  htn  dm   cad  appet  pe  ane classification  
0    0  48.0  8.0  1.020  1.0  0.0   NaN  normal  notpresent  notpresent ...  44.0  7800.0  5.2  yes  yes  no  good  no  no  ckd  
1    1  7.0   50.0  1.020  4.0  0.0   NaN  normal  notpresent  notpresent ...  38.0  6000.0  NaN  no  no  no  good  no  no  ckd  
2    2  62.0  8.0   1.010  2.0  3.0  normal  normal  notpresent  notpresent ...  31.0  7500.0  NaN  no  yes  no  poor  no  yes  ckd  
3    3  48.0  70.0  1.005  4.0  0.0  normal  abnormal  present  notpresent ...  32.0  6700.0  3.9  yes  no  no  poor  yes  yes  ckd  
4    4  51.0  8.0   1.010  2.0  0.0  normal  normal  notpresent  notpresent ...  35.0  7300.0  4.6  no  no  no  good  no  no  ckd  
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  
395  395  55.0  8.0   1.020  0.0  0.0  normal  normal  notpresent  notpresent ...  47.0  6700.0  4.9  no  no  no  good  no  no  notckd  
396  396  42.0  70.0  1.025  0.0  0.0  normal  normal  notpresent  notpresent ...  54.0  7800.0  6.2  no  no  no  good  no  no  notckd  
397  397  12.0  8.0   1.020  0.0  0.0  normal  normal  notpresent  notpresent ...  49.0  6600.0  5.4  no  no  no  good  no  no  notckd  
398  398  17.0  60.0  1.025  0.0  0.0  normal  normal  notpresent  notpresent ...  51.0  7200.0  5.9  no  no  no  good  no  no  notckd  
399  399  58.0  8.0  1.025  0.0  0.0  normal  normal  notpresent  notpresent ...  53.0  6800.0  6.1  no  no  no  good  no  no  notckd
```

384 rows × 26 columns

Filling Missing Numerical Values with Median Values

```
In [6]: 1 dataset.fillna(dataset.median(), inplace=True)  
2 dataset
```

```
/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_3973/3113556935.py:1: FutureWarning: The default  
of numeric_only in DataFrame.median is deprecated. In a future version, it will default to False. In addition  
specifying 'numeric_only=None' is deprecated. Select only valid columns or specify the value of numeric_only to  
ce this warning.  
dataset.fillna(dataset.median(), inplace=True)
```

```
Out[6]:   id  age   F   sg   al   su   rbc     pc     pcc     ba ...   pcv     wc   rc  htn  dm   cad  appet  pe  ane classification  
0    0  48.0  8.0  1.020  1.0  0.0   NaN  normal  notpresent  notpresent ...  44.0  7800.0  5.2  yes  yes  no  good  no  no  ckd  
1    1  7.0   50.0  1.020  4.0  0.0   NaN  normal  notpresent  notpresent ...  38.0  6000.0  4.7  no  no  no  good  no  no  ckd  
2    2  62.0  8.0   1.010  2.0  3.0  normal  normal  notpresent  notpresent ...  31.0  7500.0  4.7  no  yes  no  poor  no  yes  ckd  
3    3  48.0  70.0  1.005  4.0  0.0  normal  abnormal  present  notpresent ...  32.0  6700.0  3.9  yes  no  no  poor  yes  yes  ckd  
4    4  51.0  8.0   1.010  2.0  0.0  normal  normal  notpresent  notpresent ...  35.0  7300.0  4.6  no  no  no  good  no  no  ckd  
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  
395  395  55.0  8.0   1.020  0.0  0.0  normal  normal  notpresent  notpresent ...  47.0  6700.0  4.9  no  no  no  good  no  no  notckd  
396  396  42.0  70.0  1.025  0.0  0.0  normal  normal  notpresent  notpresent ...  54.0  7800.0  6.2  no  no  no  good  no  no  notckd  
397  397  12.0  8.0   1.020  0.0  0.0  normal  normal  notpresent  notpresent ...  49.0  6600.0  5.4  no  no  no  good  no  no  notckd  
398  398  17.0  60.0  1.025  0.0  0.0  normal  normal  notpresent  notpresent ...  51.0  7200.0  5.9  no  no  no  good  no  no  notckd  
399  399  58.0  8.0  1.025  0.0  0.0  normal  normal  notpresent  notpresent ...  53.0  6800.0  6.1  no  no  no  good  no  no  notckd
```

384 rows × 26 columns

முத்து வசமதி

```
In [7]: 1 dataset.isna().sum()
```

```
Out[7]: id          0  
age         0  
F           0  
sg          0  
al          0  
su          0  
rbc        146  
pc          60  
pcc         0  
ba          0  
bgr         0  
bu          0  
sc          0  
sod         0  
pot         0  
hemo        0  
pcv         0  
wc          0  
rc          0  
htn         0  
dm          0  
cad         0  
appet       0  
pe          0  
ane         0  
classification  0  
dtype: int64
```

Finding Mode to Replace Missing Categorical Values

```
In [8]: 1 categorical_columns = ['rbc','pc']  
2 modes = dataset[categorical_columns].mode()  
3 modes
```

```
Out[8]: rbc      pc  
0  normal  normal
```

Replacing the Missing Categorical Values With Calculated Mode

```
In [9]: 1 dataset[dataset.select_dtypes(include=['object']).columns] = dataset.select_dtypes(include=['object']).ap  
2 dataset
```

```
Out[9]:   id  age    F    sg    al    su    rbc    pc    pcc    ba    ...    pcv    wc    rc    htn    dm    cad    appet    pe    ane    classification  
0   0  48.0  80.0  1.020  1.0  0.0  normal  normal  notpresent  notpresent  ...  44.0  7800.0  5.2  yes  yes  no  good  no  no  ckd  
1   1  7.0   50.0  1.020  4.0  0.0  normal  normal  notpresent  notpresent  ...  38.0  6000.0  4.7  no   no  no  good  no  no  ckd  
2   2  62.0  80.0  1.010  2.0  3.0  normal  normal  notpresent  notpresent  ...  31.0  7500.0  4.7  no   yes  no  poor  no  yes  ckd  
3   3  48.0  70.0  1.005  4.0  0.0  normal  abnormal  present  notpresent  ...  32.0  6700.0  3.9  yes  no   no  poor  yes  yes  ckd  
4   4  51.0  80.0  1.010  2.0  0.0  normal  normal  notpresent  notpresent  ...  35.0  7300.0  4.6  no   no  no  good  no  no  ckd  
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  
395 395  55.0  80.0  1.020  0.0  0.0  normal  normal  notpresent  notpresent  ...  47.0  6700.0  4.9  no   no  no  good  no  no  notckd  
396 396  42.0  70.0  1.025  0.0  0.0  normal  normal  notpresent  notpresent  ...  54.0  7800.0  6.2  no   no  no  good  no  no  notckd  
397 397  12.0  80.0  1.020  0.0  0.0  normal  normal  notpresent  notpresent  ...  49.0  6600.0  5.4  no   no  no  good  no  no  notckd  
398 398  17.0  60.0  1.025  0.0  0.0  normal  normal  notpresent  notpresent  ...  51.0  7200.0  5.9  no   no  no  good  no  no  notckd  
399 399  58.0  80.0  1.025  0.0  0.0  normal  normal  notpresent  notpresent  ...  53.0  6800.0  6.1  no   no  no  good  no  no  notckd
```

384 rows × 26 columns

முத்து வசமதி

```
In [10]: 1 dataset.isnull().sum()
```

```
Out[10]: id          0  
age         0  
F           0  
sg          0  
al          0  
su          0  
rbc         0  
pc          0  
pcc         0  
ba          0  
bgr         0  
bu          0  
sc          0  
sod         0  
pot         0  
hemo        0  
pcv         0  
wc          0  
rc          0  
htn         0  
dm          0  
cad         0  
appet       0  
pe          0  
ane         0  
classification 0  
dtype: int64
```

Defining Quan & Qual for Outlier Removal

```
In [11]: 1 def quanQual(dataset):  
2     quan=[]  
3     qual=[]  
4     for columnName in dataset.columns:  
5         #print(columnName)  
6         if(dataset[columnName].dtype=='O'): #print("qual")  
7             qual.append(columnName)  
8         else:  
9             #print("quan")  
10            quan.append(columnName)  
11    return quan,qual
```

```
In [12]: 1 quan,qual=quanQual(dataset)  
2 quan
```

```
Out[12]: ['id',  
          'age',  
          'F',  
          'sg',  
          'al',  
          'su',  
          'bgr',  
          'bu',  
          'sc',  
          'sod',  
          'pot',  
          'hemo',  
          'pcv',  
          'wc',  
          'rc']
```

முத்து வசமதி

```
In [13]: 1 dataset[qual]
```

	rbc	pc	pcc	ba	htn	dm	cad	appet	pe	ane	classification
0	normal	normal	notpresent	notpresent	yes	yes	no	good	no	no	ckd
1	normal	normal	notpresent	notpresent	no	no	no	good	no	no	ckd
2	normal	normal	notpresent	notpresent	no	yes	no	poor	no	yes	ckd
3	normal	abnormal	present	notpresent	yes	no	no	poor	yes	yes	ckd
4	normal	normal	notpresent	notpresent	no	no	no	good	no	no	ckd
...
395	normal	normal	notpresent	notpresent	no	no	no	good	no	no	notckd
396	normal	normal	notpresent	notpresent	no	no	no	good	no	no	notckd
397	normal	normal	notpresent	notpresent	no	no	no	good	no	no	notckd
398	normal	normal	notpresent	notpresent	no	no	no	good	no	no	notckd
399	normal	normal	notpresent	notpresent	no	no	no	good	no	no	notckd

384 rows × 11 columns

Outlier Removal

```
In [14]: 1 descriptive=pd.DataFrame(index=["Mean","Median","Mode","Q1:25%","Q2:50%",
2                                     "Q3:75%","99%","Q4:100%","IQR","1.5rule","Lesser","Greater","Min","Max"],
3                                     for columnName in quan:
4                                         descriptive[columnName]["Mean"]=dataset[columnName].mean()
5                                         descriptive[columnName]["Median"]=dataset[columnName].median()
6                                         descriptive[columnName]["Mode"]=dataset[columnName].mode()[0]
7                                         descriptive[columnName]["Q1:25%"]=dataset.describe()[columnName]["25%"]
8                                         descriptive[columnName]["Q2:50%"]=dataset.describe()[columnName]["50%"]
9                                         descriptive[columnName]["Q3:75%"]=dataset.describe()[columnName]["75%"]
10                                        descriptive[columnName]["99%"]=np.percentile(dataset[columnName],99)
11                                        descriptive[columnName]["Q4:100%"]=dataset.describe()[columnName]["max"]
12                                        descriptive[columnName]["IQR"]=descriptive[columnName]["Q3:75%"]-descriptive[columnName]["Q1:25%"]
13                                        descriptive[columnName]["1.5rule"] = 1.5*descriptive[columnName]["IQR"]
14                                        descriptive[columnName]["Lesser"] = descriptive[columnName]["Q1:25%"]-descriptive[columnName]["1.5rule"]
15                                        descriptive[columnName]["Greater"] = descriptive[columnName]["Q3:75%"]+descriptive[columnName]["1.5rule"]
16                                        descriptive[columnName]["Min"] = dataset[columnName].min()
17                                        descriptive[columnName]["Max"] = dataset[columnName].max()
```

```
In [15]: 1 lesser=[]
2 greater=[]
3
4 for columnName in quan:
5     if(descriptive[columnName]["Lesser"]>descriptive[columnName]["Min"]):
6         lesser.append(columnName)
7     if(descriptive[columnName]["Greater"]<descriptive[columnName]["Q4:100%"]):
8         greater.append(columnName)
9
```

```
In [16]: 1 lesser
Out[16]: ['age', 'F', 'sg', 'bgr', 'sod', 'pot', 'hemo', 'pcv', 'wc', 'rc']
```

```
In [17]: 1 greater
Out[17]: ['F', 'su', 'bgr', 'bu', 'sc', 'sod', 'pot', 'wc', 'rc']
```

```
In [18]: 1 for column in lesser:
2     dataset[column][dataset[column]<descriptive[column]["Lesser"]]=descriptive[column]["Lesser"]
3 for column in greater:
4     dataset[column][dataset[column]>descriptive[column]["Greater"]]=descriptive[column]["Greater"]
```

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#using-a-view-versus-a-copy
dataset[column][dataset[column]<descriptive[column]["Lesser"]]=descriptive[column]["Lesser"]
/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_3973/3400726572.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#using-a-view-versus-a-copy
dataset[column][dataset[column]<descriptive[column]["Lesser"]]=descriptive[column]["Lesser"]
/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_3973/3400726572.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#using-a-view-versus-a-copy
dataset[column][dataset[column]<descriptive[column]["Lesser"]]=descriptive[column]["Lesser"]
/var/folders/07/ykgp85052b11h5kz22ghn8l40000gn/T/ipykernel_3973/3400726572.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

```
In [19]: 1 descriptive=pd.DataFrame(index=["Mean","Median","Mode","Q1:25%","Q2:50%",
2                                     "Q3:75%","99%","Q4:100%","IQR","1.5rule","Lesser","Greater","Min","Max"],
3                                     for columnName in quan:
4                                         descriptive[columnName]["Mean"]=dataset[columnName].mean()
5                                         descriptive[columnName]["Median"]=dataset[columnName].median()
6                                         descriptive[columnName]["Mode"]=dataset[columnName].mode()[0]
7                                         descriptive[columnName]["Q1:25%"]=dataset.describe()[columnName]["25%"]
8                                         descriptive[columnName]["Q2:50%"]=dataset.describe()[columnName]["50%"]
9                                         descriptive[columnName]["Q3:75%"]=dataset.describe()[columnName]["75%"]
10                                        descriptive[columnName]["99%"]=np.percentile(dataset[columnName],99)
11                                        descriptive[columnName]["Q4:100%"]=dataset.describe()[columnName]["max"]
12                                        descriptive[columnName]["IQR"]=descriptive[columnName]["Q3:75%"]-descriptive[columnName]["Q1:25%"]
13                                        descriptive[columnName]["1.5rule"] = 1.5*descriptive[columnName]["IQR"]
14                                        descriptive[columnName]["Lesser"] = descriptive[columnName]["Q1:25%"]-descriptive[columnName]["1.5rule"]
15                                        descriptive[columnName]["Greater"] = descriptive[columnName]["Q3:75%"]+descriptive[columnName]["1.5rule"]
16                                        descriptive[columnName]["Min"] = dataset[columnName].min()
17                                        descriptive[columnName]["Max"] = dataset[columnName].max()
```

```
In [20]: 1 descriptive
```

```
Out[20]:
```

	id	age	F	sg	al	su	bgr	bu	sc	sod	pot	hemo	pcv
Mean	198.979167	51.574219	76.002604	1.017728	0.90625	0.0	134.552083	49.758984	2.061328	138.205729	4.378385	12.544434	39.114583
Median	198.5	55.0	80.0	1.02	0.0	0.0	122.0	42.0	1.3	138.0	4.4	12.6	40.0
Mode	0	60.0	80.0	1.02	0.0	0.0	224.25	110.125	5.4	138.0	4.4	12.6	40.0
Q1:25%	97.75	42.0	70.0	1.015	0.0	0.0	101.75	27.0	0.9	135.0	4.0	10.9	34.0
Q2:50%	198.5	55.0	80.0	1.02	0.0	0.0	122.0	42.0	1.3	138.0	4.4	12.6	40.0
Q3:75%	301.25	65.0	80.0	1.02	2.0	0.0	150.75	60.25	2.7	141.0	4.8	14.525	44.0
99%	395.17	80.17	95.0	1.025	4.0	0.0	224.25	110.125	5.4	150.0	6.0	17.617	53.17
Q4:100%	399.0	90.0	95.0	1.025	5.0	0.0	224.25	110.125	5.4	150.0	6.0	17.8	54.0
IQR	203.5	23.0	10.0	0.005	2.0	0.0	49.0	33.25	1.8	6.0	0.8	3.625	10.0
1.5rule	305.25	34.5	15.0	0.0075	3.0	0.0	73.5	49.875	2.7	9.0	1.2	5.4375	15.0
Lesser	-207.5	7.5	55.0	1.0075	-3.0	0.0	28.25	-22.875	-1.8	126.0	2.8	5.4625	19.0
Greater	606.5	99.5	95.0	1.0275	5.0	0.0	224.25	110.125	5.4	150.0	6.0	19.9625	59.0

```
2 greater = []
3
4 for columnName in quan:
5     if(descriptive[columnName]["Lesser"]>descriptive[columnName]["Min"]):
6         lesser.append(columnName)
7     if(descriptive[columnName]["Greater"]<descriptive[columnName]["Q4:100%"]):
8         greater.append(columnName)
```

```
In [22]: 1 lesser
```

```
Out[22]: []
```

```
In [23]: 1 greater
```

```
Out[23]: []
```

```
In [24]: 1 dataset
```

```
Out[24]:
```

	id	age	F	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appet	pe	ane	classification
0	0	48.0	8.00	1.0200	1.0	0.0	normal	normal	notpresent	notpresent	...	44.0	7800.0	5.2000	yes	yes	no	good	no	no	ckd
1	1	7.5	55.0	1.0200	4.0	0.0	normal	normal	notpresent	notpresent	...	38.0	6000.0	4.7000	no	no	no	good	no	no	ckd
2	2	62.0	80.0	1.0100	2.0	0.0	normal	normal	notpresent	notpresent	...	31.0	7500.0	4.7000	no	yes	no	poor	no	yes	ckd
3	3	48.0	70.0	1.0075	4.0	0.0	normal	abnormal	present	notpresent	...	32.0	6700.0	3.9000	yes	no	no	poor	yes	yes	ckd
4	4	51.0	80.0	1.0100	2.0	0.0	normal	normal	notpresent	notpresent	...	35.0	7300.0	4.6000	no	no	no	good	no	no	ckd
...	
395	395	55.0	80.0	1.0200	0.0	0.0	normal	normal	notpresent	notpresent	...	47.0	6700.0	4.9000	no	no	no	good	no	no	notckd
396	396	42.0	70.0	1.0250	0.0	0.0	normal	normal	notpresent	notpresent	...	54.0	7800.0	5.8125	no	no	no	good	no	no	notckd
397	397	12.0	80.0	1.0200	0.0	0.0	normal	normal	notpresent	notpresent	...	49.0	6600.0	5.4000	no	no	no	good	no	no	notckd
398	398	17.0	60.0	1.0250	0.0	0.0	normal	normal	notpresent	notpresent	...	51.0	7200.0	5.8125	no	no	no	good	no	no	notckd
399	399	58.0	80.0	1.0250	0.0	0.0	normal	normal	notpresent	notpresent	...	53.0	6800.0	5.8125	no	no	no	good	no	no	notckd

384 rows x 26 columns

Saving the Preprocessed dataset as a New File

```
In [25]: 1 dataset.to_csv("PreProcessed_kidney_disease.csv",index=False)
```

---- The End----

