

02/24/22
8:50pm

* Data Preprocessing / Data Preparation

⇒ Feature Engineering

Data cleaning

Missing Values

Outliers

Data Wrangling

↳ (Data Transformation)

Encoding

Scaling

Transformation.

Us.
(Car Insurance) # Data cleaning (In spyder)

import numpy as np

import pandas as pd

why two dataset

Because to check old
and new dataset.

⇒ Load Data Set

df_raw = pd.read_csv ("claimants.csv")

df = pd.read_csv ("claimants.csv")

[Out]: df_raw | DataFrame | 1340, 7 | column names : CASENUM,
df | DataFrame | 1340, 7 | column names : "

⇒ Check if there are null values.

df.isnull().sum()

out :-	CASENUM	0
	CLMSEX	12
	CLMINSUR	41
	SEATBELT	48
	CLMAGE	189
	LOSS	0
	ATTORNEY	0

MISSING
values.

Data Understanding

- ~~~~* ~~~
- Case Num :- Number of case (insurance)
- Discrete Clm sex :- 0, 1 [Male] [Female] Person who
- Discrete Clm Insur :- 1 (Yes), 0 (No) # having Insurance is
or not claiming
- Discrete Seat belt :- 0 (No), 1 (Yes) # wearing seat belt Insurance is called
"claimants" During The accident
- Continuous Clm age :- Age of the person
- Loss :- Loss / Damage happened
- Attorney :- Whether they hire a lawyer
(or) not, [0] No, [1] Yes.

There are various methods for filling null values

* For Continuous Data

mean (no outliers)

Because Mean is

* For Discrete Data

Mode

(Less affected) Impacted with Outliers

Replacing null values For **discrete** Variables

~~~~ \* ~~~\* ~~~~  
So, we have to Replace Them with "**mode**"

# df[ "CLMSEX" ]. value\_counts()  
(mode)

[out] :- 1.0 742 → Highest, so we  
0.0 586 Replace with "1"

# df[ "CLMSEX" ]. fillna(1, inplace=True)  
↓ changing in original  
Data = TRUE

[out] Filled with "1"

# df[ "CLMINSUR" ]. Value Counts()

[out] 1.0 1820 → MODE  
0.0 120

# df[ "CLMINSUR" ]. fillna(1, inplace=True)

[out] Filled with "1", so replaced with  
"1" value

```
# df[["SEATBELT"]].valuecounts()
```

[Out]: 0.0 1270 → Mode.  
1.0 22

```
# df[["SEATBELT"]].fillna(0, inplace=True)
```

[Out]: Filled with "0"

Time  
11:30pm

# Replacing null Values For Continuous Variables  
By  
# Pandas

```
# df[["CLMAGE"]].mean()
```

[Out]: 28.414

```
# df[["CLMAGE"]].median()
```

[Out]: 30.00

```
# df[["CLMAGE"]].fillna(28.414, inplace=True)
```

[Out]: Filled with 28.414

If we ask again what is median value.

```
# df[["CLMAGE"]].median()
```

[Out]: 28.414, so it changes 30.00 To 28.414

Because we Filled with null values  
28.414.

# # sklearn (scikit Learn) (One stop shop For ML)

## \* SIMPLE IMPUTER

# Load data set again  
Library → (module) (Function)

From

Sklearn.impute

import SimpleImputer

Ex: missing values = "???"

calculating mean Value  
 $\frac{\sum}{n}$

# mean\_imputer = SimpleImputer (strategy = "mean")  
(Function) → (Argument)

name  
as

Function. (or) Store in  
mean\_imputer

Ex:- Core Python

[Keyword  
argument]

If we don't  
write.  
This time  
it gives

Same : a = 4

# Converting array  
To data frame  
To Remove  
array. From mean value.

# df[["CLMAGE"]] = pd.DataFrame (mean\_imputer.fit\_transform

array([50.,  
18.,  
5.,  
...])

storing

- Trans

form (df[["CLMAGE"]]))

Fill in This column  
(or) Replacing

Calculate  
The mean  
value

Out

Filled with 28.414

For Median Only change mean To median

# median\_imputer = SimpleImputer (strategy = "Median")

# df[["CLMAGE"]] = pd.DataFrame (median\_imputer. fit -  
Transform (df[["CLMAGE"]]))

For Mode [Discrete Data]

mode\_imputer = SimpleImputer (strategy = "mode")

# df[["CLMSEX"]] = pd.DataFrame (mode\_imputer. fit\_transform (df[["CLMSEX"]]))

# df[["CLMINSUR"]] = pd.DataFrame (mode\_imputer. fit\_transform  
(df[["CLMINSUR"]]))

# df[["SEATBELT"]] = pd.DataFrame (mode\_imputer. fit\_transform  
(df[["SEATBELT"]]))

# **Output** Fills Directly The mode  
Value in to the dataset.

# MODE = it is for Discrete Data.

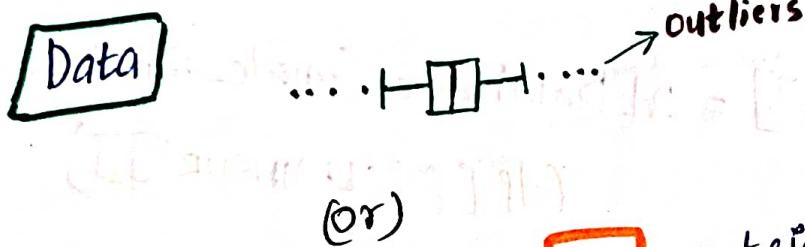
# MEAN → Continuous Data  
# MEDIAN

3/4/22 2:08 AM

# Outliers → # Applicable Only For  
Continuous Data

1 Que :- What is Outlier?

Ans :- An Outlier is a data point in a dataset.  
that is distant from all other observation  
which is significantly different from the remaining



(or)

- The data point that lies outside the overall distribution of the dataset.

2 Que :- What are Impacts having Outliers in a DataSet?

Ans :- It will cause various problems in statistical analysis. (it may cause a significant impact on mean and standard deviation)

Ex :-  $\bar{x} = \text{mean (affect)}$

$$\frac{\sum (x - \bar{x})^2}{n-1} = \text{variance}$$

$$\sqrt{\frac{\sum (x - \bar{x})^2}{n-1}} \quad (\text{standard deviation})$$

all are affected by outliers

\* In addition, Some Machine learning models are Sensitive to Outliers, which may Decrease There Performance

3 Que : Reasons For Outliers ?

- Ams :
1. Data Entry Errors (Ex:- Salary = 100,000, But 10,000 Entering)
  2. Measurement Errors (Ex:- Measuring in meters instead of KM)
  3. Instrumental Error

4 Que :- Types of Outliers ?

- Ams :-
- \* **Univariate Outlier** ---> Identifying Outlier For Single Variable
  - \* **Bivariate Outlier** ---> Identifying as Outlier by Analyzing 2 Variables at a Time

# Importing Libraries

- # Import numpy as np
- # Import pandas as pd
- # Import matplotlib as plt.
- # Import Seaborn as sns.

## Real Time Case study

### boston - House Price Dataset

```
# boston = pd.read_csv("D:\\data science \\Shubham\\Krishna class \\6. Data cleaning \\boston.csv")
```

```
boston.head()
```

out

|   | CRIM    | ZN   | INDUS | CHAS | NOX   | RM    | AGE  | DIS   | RAD | TAX  | PTRATIO |
|---|---------|------|-------|------|-------|-------|------|-------|-----|------|---------|
| 0 | 0.00632 | 18.0 | 2.31  | 0    | 0.538 | 6.575 | 65.2 | 4.090 | 1   | 29.6 | 15.3    |
| 1 | 0.02731 | 0.0  | 7.07  | 0    | 0.469 | 6.421 | 78.9 | 4.967 | 2   | 24.2 | 17.8    |
| 2 | 0.02729 | 0.0  | 7.07  | 0    | 0.469 | 7.185 | 61.1 | 4.967 | 2   | 24.2 | 17.8    |
| 3 | 0.03237 | 0.0  | 2.18  | 0    | 0.458 | 6.998 | 45.8 | 6.062 | 3   | 22.2 | 18.7    |
| 4 | 0.06905 | 0.0  | 2.18  | 0    | 0.458 | 7.147 | 54.2 | 6.062 | 3   | 22.2 | 18.7    |

| B      | LSTAT | PRICE |
|--------|-------|-------|
| 396.90 | 4.98  | 24.0  |
| 396.90 | 9.14  | 21.6  |
| 392.83 | 4.03  | 34.7  |
| 394.63 | 2.94  | 33.4  |
| 396.90 | 5.33  | 36.2  |

```
# boston.info()
```

506 rows, 14 columns

out

|    |         |         |
|----|---------|---------|
| 0  | CRIM    | float64 |
| 1  | ZN      | float64 |
| 2  | INDUS   | float64 |
| 3  | CHAS    | → int64 |
| 4  | NOX     | float64 |
| 5  | RM      | float64 |
| 6  | AGE     | float64 |
| 7  | DIS     |         |
| 8  | RAD     | → int64 |
| 9  | TAX     | → int64 |
| 10 | PTRATIO | float64 |
| 11 | B       | float64 |
| 12 | LSTAT   | float64 |
| 13 | PRICE   | float64 |

(float64(11), int64(3))

## Various ways of finding The Outlier

Changing Every Value to Z score.

option ①

Z score

$$|Z\text{score}| > 2$$

$$\left[ \begin{array}{l} \therefore z < -2 \\ z > +2 \end{array} \right] \mid \text{modulus} \mid$$

Ex :-

| X | Z <sub>score</sub> |
|---|--------------------|
| 6 | 4.56               |
| 4 | -1                 |
| 9 | 9 - 5.6            |
| 7 | 7 - 5.6            |
| 2 | 2 - 5.6            |

$$Z = \frac{X - \mu}{\sigma}$$

$$I = \frac{6 - 5.6}{1}$$

$$\therefore \mu = \frac{6 + 4 + 9 + 7 + 2}{5}$$

$$\mu = 28/5 = 5.6$$

$$EX/\sigma = 1$$

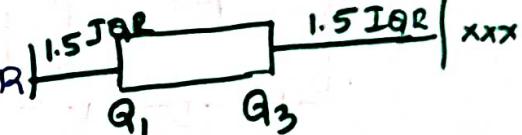
option ②

IQR

(Inter Quartile Range)

any value  $< Q_1 - 1.5 \text{ IQR}$

any value  $> Q_3 + 1.5 \text{ IQR}$



# any value Less than  $Q_1 - 1.5 \text{ IQR}$  is Outlier

# any value More Than  $Q_3 + 1.5 \text{ IQR}$  is Outlier

option ③

Visualization

\* Box plot, Histogram For Univariate Outliers

\* Scatter plot For bivariate Outliers.

option③

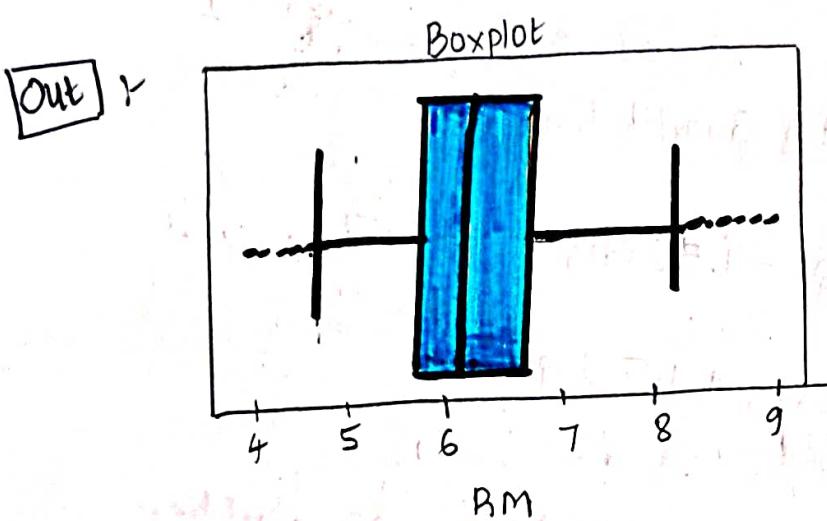
Detection of Outliers of "RM" column (Based on Boxplot)

# sns. boxplot (boston.RM)

plt. title ("Boxplot")

plt. show()

\* Overall view of Outlier



Outliers in Both Tails of "RM".

option①

Detection of Outliers of "RM" column (Based on Zscore)

\* For Exact Outliers in Data

$$Z\text{score} = \frac{x - \mu}{\sigma}$$

# boston["RM"].mean(), boston["RM"].std()

[out] (6.284, 0.702)

# boston["RM\_Zscore"] = ([boston["RM"] - boston["RM"].mean()]) / boston["RM"].std()

creating new column  
and storing Zscores

$$\frac{\underline{[RM].mean()}}{\underline{\mu}} \quad \underline{\%} \quad \underline{\sigma}$$

Output :-

clarity :-

$$\frac{x - \mu}{\sigma} = \frac{(\text{boston}['RM']) - \text{boston}['RM'].mean())}{\text{boston}['RM'].std()}$$

$$\frac{x - \mu}{\sigma}$$

$$\frac{x - \mu}{\sigma}$$

E Answers Ami

$$\text{boston}['RM\_Zscore'] =$$

↑ Endulo store age + thage.

ZSCORE

|     | CRIM    | ZN   | INDUS | CHAS | NOX   | RM | AGE | DIS | RAD | TAX | PTRATIO | B | ... RM SCORE |
|-----|---------|------|-------|------|-------|----|-----|-----|-----|-----|---------|---|--------------|
| 0   | 0.0632  | 18.0 | 2.31  |      | 6.515 |    |     |     |     |     |         |   | 0.413263     |
| 1   | 0.02731 | 0.0  | 7.07  |      | 6.421 |    |     |     |     |     |         |   | 0.194082     |
| 2   | 0.02729 | 0.0  | 2.18  |      | 7.185 |    |     |     |     |     |         |   | 1.281446     |
| ..  |         |      |       |      |       |    |     |     |     |     |         |   |              |
| 505 | 0.04741 | 0.0  | 11.93 |      | 6.030 |    |     |     |     |     |         |   | -0.362408    |

506 x 15 columns

Now, Identify Outliers in this column ↑

# Outlier =  $\text{boston}[(\text{boston}['RM\_Zscore'] > +2) \text{ or } (\text{boston}['RM\_Zscore'] < -2)]$

2 conditions

Greater Than +2  
Less Than -2

[Out] :- Outlier

| CRIM | RM SCORE |
|------|----------|
| 97   | 2.539    |
| 98   | 2.185    |
| 162  | 3.47     |
| 163  | -2.121   |
| 66   | -3.055   |

3/4/22 4:30 AM

Time  
1:00pm

## # Outlier. Shape

[Out] :- (31, 15)

option 2  
# Detection of Outlier of RM (based on IQR)

- Calculate First ( $Q_1$ ) and third quartile ( $Q_3$ )
  - Find interquartile range ( $Q_3 - Q_1$ )
  - Find Lower bound  $Q_1 - 1.5$  & Find upper bound  $Q_3 + 1.5$
- for Exact no. of range + - of outliers

#  $Q_3 = \text{boston}["RM"].quantile(0.75)$

#  $Q_1 = \text{boston}["RM"].quantile(0.25)$

$$Q_1 = 0.25(w)$$

$$Q_3 = 0.75(w)$$

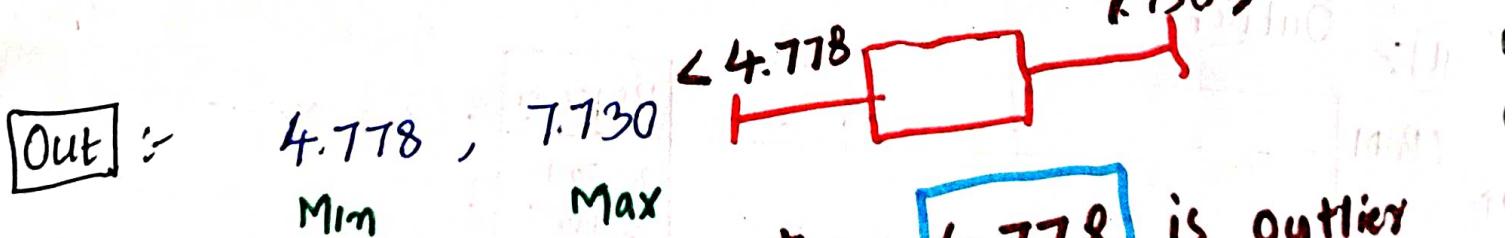
$$IQR = Q_3 - Q_1$$

$$\text{Lower-limit} = Q_1 - (IQR * 1.5) \quad IQR = Q_3 - Q_1$$

$$\text{Upper-limit} = Q_3 + (IQR * 1.5) \quad \Rightarrow Q_3 + IQR * 1.5$$

$$\Rightarrow Q_1 - IQR * 1.5$$

Lower-limit, upper-Limit



- \* Which is having less than 4.778 is outlier
- \* Which is having greater than 7.730 is outlier

## Methods of deal The Outliers

- 3R → \* Remove → Deleting The Outliers
- Technique \* Replace → changing values / Data Manipulation
- \* Retain → Treat Them Separately.  
Work without outlier / with outlier  
"Work" Separately.

Outliers should be detected and "REMOVED" Only From Training data set, NOT from The Test Set. So we should first divide our data set into train and test. and remove Outliers in the train set, but keep those in test set, and Measure how well our model is doing.

option ① Removing \* (Let's trimm The dataset) 74.77  
# boston\_trimmed = boston[(boston["RM"] > Lower-limit) & (boston["RM"] < Upper-limit)] 7.73

boston\_trimmed  
out :- which ever Data is Greater than 4.77 and Less Than 7.73 Consider That data Only, Remove Remaining data.

|   | RM    | RM - Zscore |
|---|-------|-------------|
| 0 | 6.515 | 0.413263    |
| 1 | 6.421 | 0.194082    |
| 2 | 7.185 | 1.2814416   |
| 3 | 6.794 | 0.724955    |
| 4 | 6.030 | -0.362408   |

∴ From 506 records, Removed 30 records

that Means 476 having (or) there.

⇒ it's about 6% of the data is Removed

⇒ it is Only in One Variable (1 column) = 6%.

⇒ If Total Variable (all columns)  $\frac{1,2,3,4,\dots}{(14) \rightarrow [10] \text{ continuous data}}$  Remove May Delete 100 rows  
20% of Data

So, In real Life, we don't use outliers "Remove"

Even, Though, If we Want Remove, still it Contains Outliers.

⇒ Outliers in The Trimmed dataset

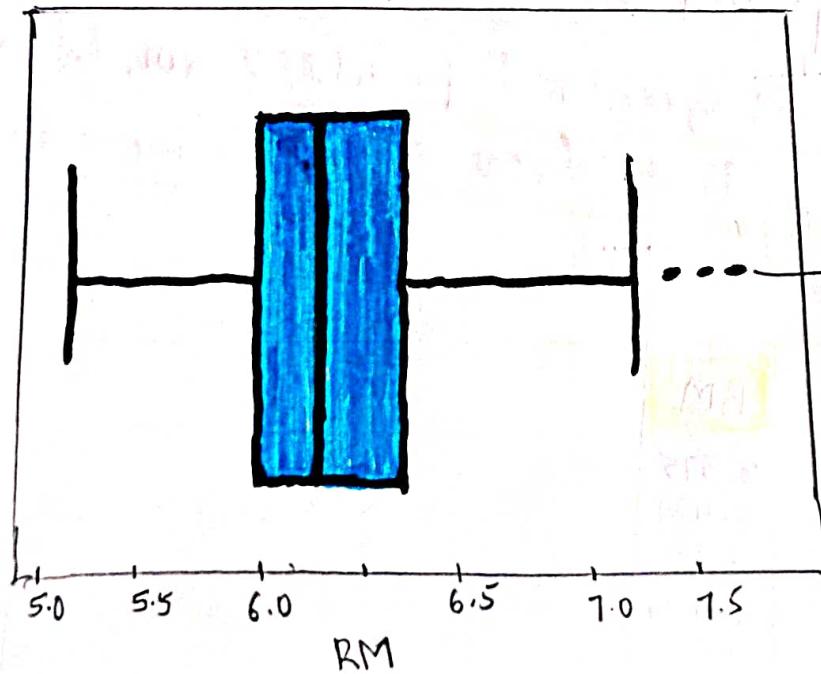
# sns.boxplot(boston\_trimmed.RM)

plt.title("Boxplot")

plt.show()

Boxplot

Out



still, we have Outliers, after Removing The Outlier records

Q: Why? it showing Outliers after Removing records of outliers

A:- When we remove datapoints from our dataset, all the parameters of the distribution are re-calculated.

For 476 records,  
Again, it's calculating ( $Q_1$ ,  $Q_3$ , IQR)  
*(again)*

for  
506

$$G_1 = 4.77, \quad G_1 = \text{New Value (again)} \\ G_3 = 7.730, \quad G_3 = \text{New value (again)}$$

- \* Therefore, In The new-trimmed variables Values that before were not considered Outliers.

That's why, it showing new Outliers after removing some records also.

option ②. Replace

# Replace with  $\frac{\text{Upper Limit} * \text{Based IQR}}{\text{Lower Limit}}$  (Calculated)

Replacing with  
 $4.778$  (min value)

min. value  
10

Every Outlier

130

$$Q_1 - [IQR * 1.5]$$

$$77.730 = Q_3 + [IGR^* 1.5]$$

↓ Every Outlier Replacing with

7.730 (max. value)

⇒ More Detailness

Replacing with  $Q_1, Q_3$  Values

| CRIM | ...                                                                               | RM                               | ...                              |
|------|-----------------------------------------------------------------------------------|----------------------------------|----------------------------------|
| 97   | More Than<br>7.730, so, they<br>are Outliers.<br>Reduce with<br>$\bar{x} = 7.730$ | 8.069<br>7.820<br>7.802<br>8.375 | 7.730<br>7.730<br>7.730<br>7.730 |
| 98   |                                                                                   |                                  |                                  |
| 162  |                                                                                   |                                  |                                  |
| 163  |                                                                                   |                                  |                                  |

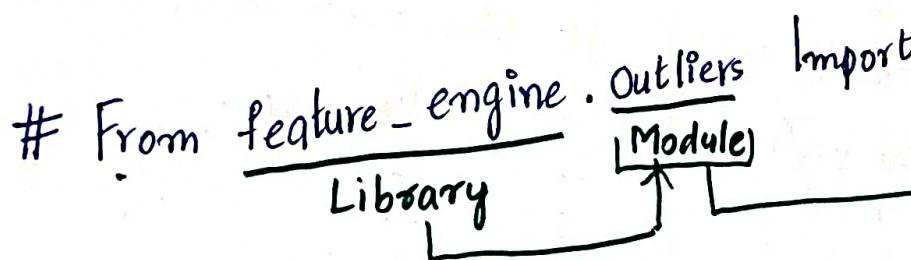
$$Q_3 - IQR * 1.5 < [4.778] \rightarrow \begin{array}{l} \text{less than } 4.778 \\ \text{= outlier} \end{array}$$

$$Q_1 + IQR * 1.5 > [7.730] \rightarrow \begin{array}{l} \text{greater than } 7.730 \\ \text{= outlier} \end{array}$$

\* Winsorizer

# PIP install

feature-engine



In Winsorizer (capping method)  
Function argument

- \* Capping method = 1. IQR
- 2. "Gaussian"

# In feature-engine. Outliers Import ArbitraryOutlierCapper

```
graph LR; FE[feature-engine] --> Outliers[Outliers module]; Outliers --> Function[ArbitraryOutlierCapper Function];
```

We can give our own Lower, upper values.

1. By IGR method

## In Depth of **EEG**

Firstly,

# we calculate  $G_1$  value

```
# Q1 = boston["RM"].quantile(0.25)
```

506  
Records

$$\therefore Q_1 = 5.885$$

$$Q_3 = 6.623$$

**Out** :- calculate  $G_3$  value.

# we calculate Q3 value.  
#  $Q_3 = \text{boston}["RM"]$ . quantile(0.75)

# 93      Q3  
Out      :- 6.623

$$IQR = Q_3 - Q_1 \\ = 0.737$$

Distance between Q<sub>3</sub> and Q<sub>1</sub>

## # IQR Value

$$\# IGR = Q_3 - Q_1$$

IQR

out

```
CODE :-  
# From feature_engine.Outliers import Winsorizer  
from feature_engine.outliers import Winsorizer  
w = Winsorizer(method="iqr", tail="both", fold=1.5, variables  
                = ["RM"])
```

# Winf = Winsozier

```
# Winf = Winsorizer (capping=10, random_state=42)
# Winf.store_in(winf)

# boston_t = Winf.fit_transform(boston[['RM']])
# boston_t = Winf.transform(boston[['RM']])
```

# boston ->  
store in  
# Print (wm. left-tail-caps-, wm. right-tail-caps-)  
4.7184 7.7304

# Sms. boxplot (boston-t . PM)

plt.title ("boxplot")

# # plt.show()

7.730E

AM) ~~IGR~~ Edited (replaced IGR) RM

again

From feature\_engine.Outliers import Winsorizer

win = Winsorizer(capping\_method="lqr", fold=1.5, Tail="both")

Variables = ["RM"])

boston\_t = win.fit\_transform(boston[["RM"]])

Print(win.left\_tail\_caps\_, win.right\_tail\_caps\_)

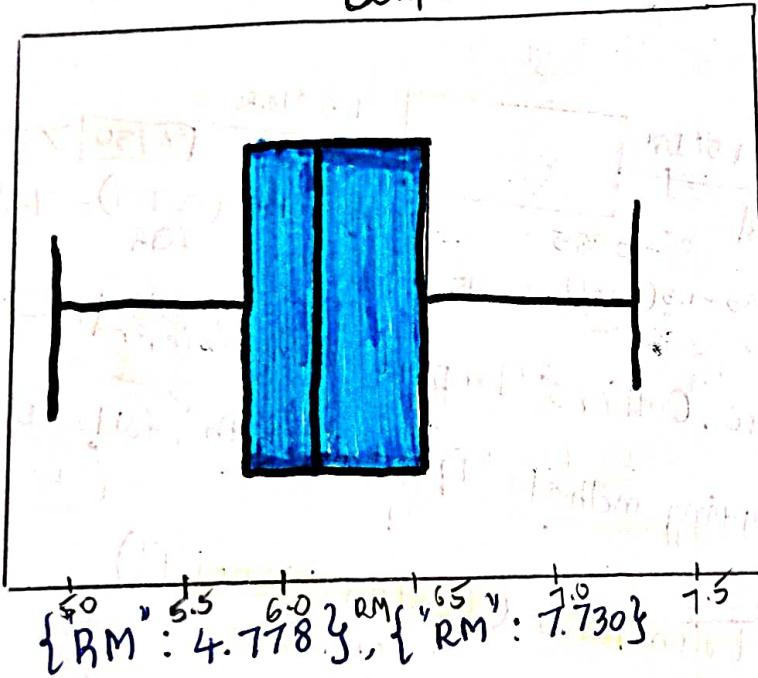
Sns.boxplot(boston\_t.RM)

plt.title("boxplot")

plt.show()

Boxplot

Out



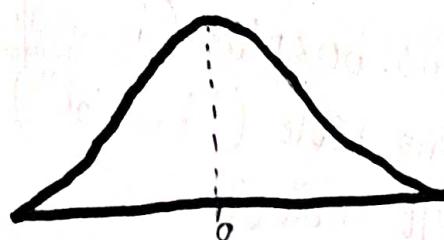
2, Replacing outliers using Winsorizer (min, max, automatically taken by "Gaussian method")

# Gaussian method

= "Normal"

Same

Distribution method



```
# from feature_engine.outliers import Winsorizer
```

```
# Wim = Winsorizer(capping_method = "gaussian", tail = "both")
```

```
fold = 1.5, variables = ['RM'])
```

```
boston_t = Wim.fit_transform(boston[['RM']])
```

```
print(Wim.left_tail_caps_, Wim.right_tail_caps_)
```

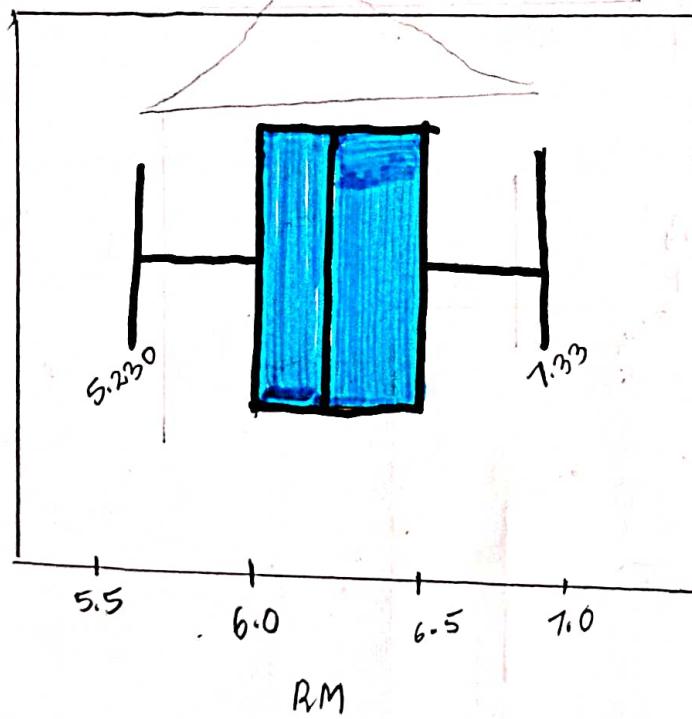
```
Sns.boxplot(boston_t.RM)
```

```
plt.title("Boxplot")
```

```
plt.show()
```

```
{'RM': 5.230} {'RM': 7.33}
```

BOXPLOT



Out

3. Replace arbitrary Outlier Capper ("The min, max, by user")  
~~~~~  
we get the **min** and **max** values based on "Domain Expert"
(or) by **Our Own Research**

from feature_engine.Outliers import ArbitraryOutlierCapper

Capper = ArbitraryOutlierCapper (**max_capping_dict** = { "RM": 7.5 },
min_capping_dict = { "RM": 4.8 })

boston_c = Capper.fit_transform(boston[["RM"]])

Print (Capper.right_tail_caps_, Capper.left_tail_caps_)

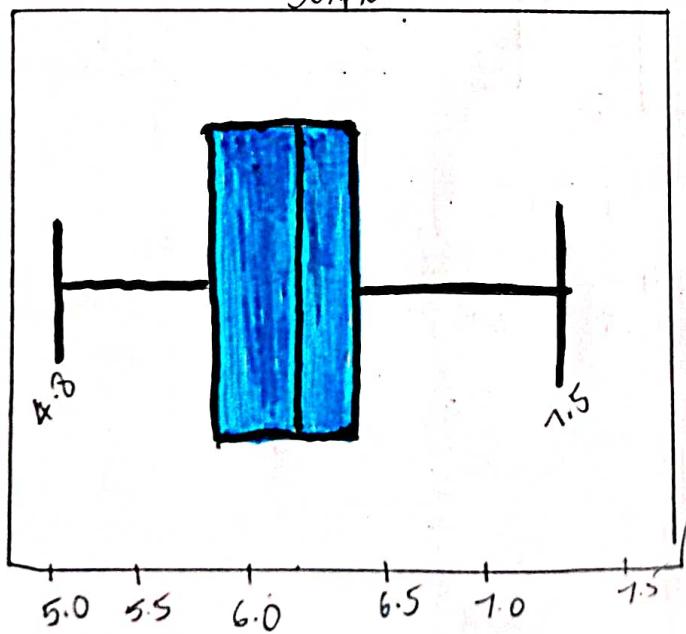
Sns.boxplot(boston_c.RM)

plt.title("Boxplot")

plt.show()

{ "RM": 4.8 }, { "RM": 7.5 }
Boxplot

Out



3/4/2022 5:20 PM

5/4/22
10:30pm

3

Discretization

Converting Continuous Data To Discrete Data

also called as "Binning"

Code :-

Import Pandas as pd

```
# stroke = pd.read_csv("stroke prediction.csv")  
# stroke.head()
```

Out	id	gender	age	hyper_tension	heart_disease	Ever_married	work_type	Residence_Type	avg_glucose_level	bmi	smoke_Status	Stroke
	30669	Male	3.0	0	0	No	children	Rural	95.12	18.0	Nan	0
	30468	Male	58.0	1	0	Yes	private	urban	87.96	39.2	Never	0
	16523	Female	8.0	0	0	No	private	urban	110.89	17.6	NAN	0
	56543	Female	70.0	0	0	Yes	private	Rural	69.04	35.9	Formerly_smoked	0
	46136	male	14.0	0	0	No	Never_Worked	Rural	161.28	19.1	NAN	0

↓ discrete (cat) ↓ continue ↓ discrete (cont.) ↓ discrete (cont.) ↓ discrete ↓ discrete ↓ discrete ↓ discrete ↓ continuous ↓ continuous ↓ continous ↓ discrete

```
# stroke.shape
```

Out : (43400, 12)

```
# stroke.info()
```

information about stroke data.

Out

\Rightarrow Create Rins

$$\# \text{ intervals} = [0, 12, 19, 30, 60, 90] \rightarrow \# \text{ bins} \Rightarrow \begin{matrix} \text{Each} \\ \text{interval} \end{matrix}$$

Intervals = 5
Categories = ["child", "teenager", "young adult", "middle aged",
"senior citizen"]

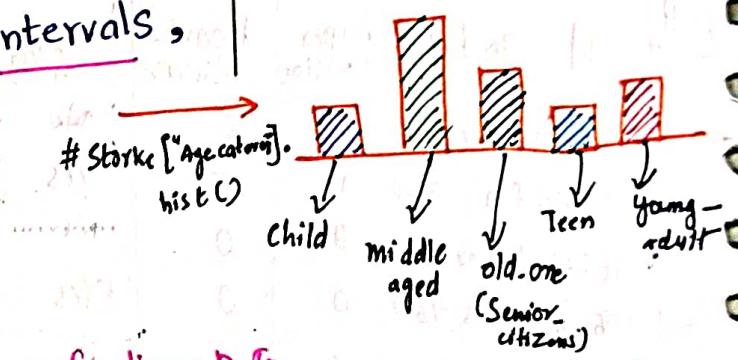
• Creating New Column & Storing The data of
"Age" Divide.

Størke ["Age-category"] = Pd. cut [x]

\downarrow
 $X = \text{Stroke } [\text{"Age"}]$, bins = intervals,
Labels = categories] # Stärke ["his"]

The Age column is
Continuous Data. So,
We converting into
Discrete Variable
intervals categorized

- * 0-12 → child
- * 12-19 → teenager
- * 19-30 → young adult
- * 30-60 → Middle aged
- * 60-90 → Senior citizen



stroke.head()

changed Continuous Data

Discrete Categorical data

Newly-arrived

4

* Encoding *

↓
Applicable For Categorical Variables

Discrete

- * Machine Can't Understand Text Data. So, We Convert the "discrete categorical" to "discrete count" Data

- * There are two types of Categorical Data

1. Nominal

(No Natural Sequence)

2. Ordinal

(Natural Sequence)

we should treat
Every variable
Equally.

⇒ To Convert Categorical Data

To Numeric:

Ordinal Data

	Grades
O	10
A+	9
A	8
B+	7
B	6
F	0

⇒ O > A+ >
10 > 9

citynames	
Hyd	0
Delhi	1
Mumba	2
Kerala	3
Gujarat	4

In this case, we can't
decide by numbering
given to state.

Nominal Data → (dummies)

which is
greater

1. get Dummies (Pandas)

2. One hot Encoding (sklearn)

Ordinal Data

1. map (pandas)

2. Label Encoder (sklearn)

# One Hot Encoding	# Label Encoding
* Nominal	* Ordinal

```

import Pandas as pd
import Numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
# df = pd.read_csv("homeprices.csv")
# df.head()

```

Out :-

	Town	Area	Price
0	chennai	2600	5500000
1	chennai	3000	5650000
2	chennai	3200	6100000
3	chennai	3600	6800000
4	Banglore	2600	5850000

Town

- 1. Chennai
 - 2. Bangalore
 - 3. Hyderabad
- Nominal Data → Discrete → Categorical Data

→ **Pd.get_dummies** [Nominal Variable Encoding using pandas]

dummies = pd.get_dummies(df["Town"])

dummies

Out :-

	Banglore	Chennai	Hyderabad
0	0	1	0
1	0	1	0
2	0	1	0
3	0	1	0
4	1	0	0
5	2	0	0
6	1	0	0
7	1	0	0
8	0	0	0
9	0	0	0
10	0	0	1

Out :-

	Banglore	Chennai	Hyderabad
0	0	1	0
1	0	1	0
2	0	1	0
3	0	1	0
4	1	0	0
5	1	0	0
6	1	0	0
7	0	0	1
8	0	0	1
9	0	0	1
10	0	0	1
11	0	0	1

Dummies

which ever record is there
it shows = 1

* For other columns it
shows = 0

df_dummies = pd.concat ([df, dummies], axis = "columns")

df_dummies

Out :-

	Town	Area	Price	Banglore	Chennai	Hyderabad
0	Chennai	2600	5500000	0	1	0
1	Chennai	3000	5650000	0	1	0
2	Chennai	3200	6100000	0	1	0
3	Chennai	3600	6800000	0	1	0
4	Banglore	3600	5850000	1	0	0
5	Banglore	2600	6150000	1	0	0
6	Banglore	2800	6500000	1	0	0
7	Banglore	3300	6100000	1	0	0
8	Hyderabad	3600	7100000	0	0	1
9	Hyderabad	2600	5750000	0	0	1
10	Hyderabad	2900	6000000	0	0	1
11	Hyderabad	3100	6200000	0	0	1

Delete The Original column

df_dummies . drop ("town", axis = "columns", inplace = True)

df_dummies

Out:

	Area	Price	Banglore	chennai	Hyderabad
0	2600	5500000	0	1	0
1	3000	5650000	0	1	0
2	3200	6100000	0	1	0
3	3800	6800000	0	1	0
4	2600	5850000	1	0	0
5	2800	6150000	1	0	0
6	3300	7100000	1	0	0
7	3600	5750000	1	0	0
8	2600	6000000	0	0	1
9	2900	6200000	0	0	1
10	3100	6900000	0	0	1
11	3600	6500000	0	0	1

after deleting original column in New Data (ie. df_dummies)

Old_data (df) is still remains same.



Dummy Variable Trap



	Bang	chenn	Hyder
0	1	0	0
0	1	0	0
1	0	0	0
1	0	0	0
0	0	1	0
0	0	1	0



	Bang	chenn
0	1	0
0	1	0
1	0	0
1	0	0
0	0	0
0	0	0

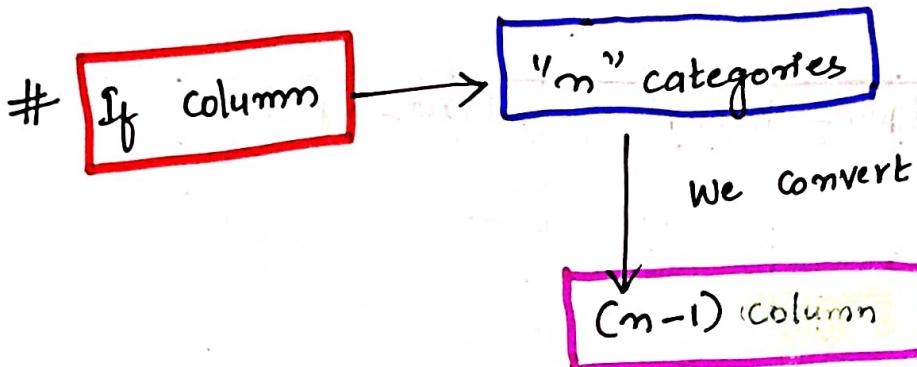
→ hyderabad

If we Remove One
of The column
also, still we can
Identify, col name by
(0,0)

Town - 3 categories

↓ Dummies

3 columns [Multi-collinearity / collation problems]



df_dummies.drop("chennai", axis="columns", inplace=True)

df_dummies

out

Area	Price	Bangalore	Hyderabad
0	2600	5500000	0
1	-	-	0
2	-	-	0
3	-	-	0
4	-	-	0
5	-	-	0
6	-	-	0
7	-	-	0
8	-	-	0
9	-	-	0
10	-	-	0
11	36000	6500000	0

where ever
Both column
having equal
zero's, it is
another column's
data.

Chennai = 0 0

We Can Write One line of Code To do all these

df_dum = pd.get_dummies (df, drop_first=True)

df_dum

out :	Area	Price	town-chennai	town-hyderabad
0			1	0
1			1	0
2			1	0
10			1	0
11			1	0

This Diagram is Exact of (Previous) last page diagram.
But deleted banglore Replaced chennai
town-banglore [deleted] → Because, B < C (Alphabet order)
Replaced with (0 0) ← hga (chennai)

⇒ One Hot Encoding → Nominal using sklearn

From sklearn.preprocessing import OneHotEncoder

enc = OneHotEncoder (drop = "First")
↑
Stores in enc

enc_df = pd.DataFrame (enc.fit_transform (df[["town"]]).toarray())
↑
Converting
↓
Output

enc-df

[Out]

	0	1
0	1.0	0.0
1	1.0	0.0
2	1.0	0.0
3	1.0	0.0
4	0.0	0.0
5	0.0	0.0
6	0.0	0.0
7	0.0	0.0
8	0.0	1.0
9	0.0	1.0
10	0.0	1.0
11	0.0	1.0

merge with main df

df_ohe = df.join(enc-df)

df_ohe.drop("town", axis="columns", inplace=True)

df_ohe

[Out]

	area	Price	0	1
0	2600	5500000	1.0	0.0
1	3000	5650000	1.0	0.0
2	-	-	1.0	0.0
3	-	-	1.0	0.0
4	-	-	0.0	0.0
5	-	-	0.0	0.0
6	-	-	0.0	0.0
7	-	-	0.0	0.0
8	-	-	0.0	0.0
9	-	-	0.0	1.0
10	-	-	0.0	1.0
11	-	-	0.0	1.0

Alphabetical Order

⇒ **Label Encoder** → **Ordinal Variable** [sklearn]

Used for binary category Variable

dfnew = df.copy() → creating "copy" from original Data.

dfnew

Out

	town	Area	Price
0	Chennai	2600	5500000
1	Chennai	3000	5650000
2	Chennai	3250	61000000
:	:	:	:
11	Hyderabad	3600	6950000

From sklearn.preprocessing import LabelEncoder

Le = LabelEncoder()

dfnew.town = Le.fit_transform(dfnew.town)

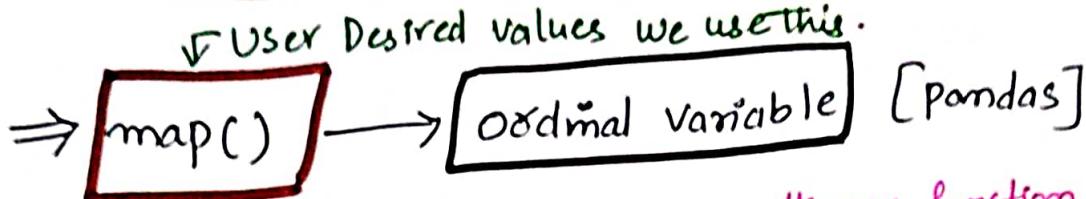
dfnew

	town	Area	Price
0	1	6.0	
1	1	6.0	
2	1	3.0	
3	1	3.0	
4	0	6.0	
5	0	6.0	
6	0	6.0	
7	0	6.0	
8	2	6.0	
9	2	6.0	
10	2	6.0	
11	2	6.0	

Out :

Alphabetical Order

Banglore [B] → 0
Chennai [C] → 1
Hyderabad [H] → 2



map function.

df_m = df.copy()

df_m ["town"] = df_m ["town"].map({ "Chennai":0, "Bangalore":2, "Hyderabad":1 })

df_m

Out :-

	town	Area	Price
0	0		
1	0		
2	0		
3	0		
4	2		
5	2		
6	2		
7	2		
8	1		
9	1		
10	1		

Which ever
order we
write it
↑ takes
in the order

Taken by our choices

Here

Chennai = 0, 20

Bangalore = 2, 15

Hyderabad = 1, 25

We can
Replace
with any
Value



From sklearn.preprocessing import OrdinalEncoder

df_new1 = df.copy()

oe = OrdinalEncoder (categories = [["Bangalore", "Hyderabad", "Chennai"]])

df_new1.town = oe.fit_transform (df_new1 [["town"]])

df_new1

Out :-

Town	0	0	0	0	0	0	0	0	0	0	0
0 - 20	2	2	2	0	0	0	0	0	0	0	0

inf/ 6/4/22
3:00 AM

Dt: 3/4/22
11:30pm

"1"

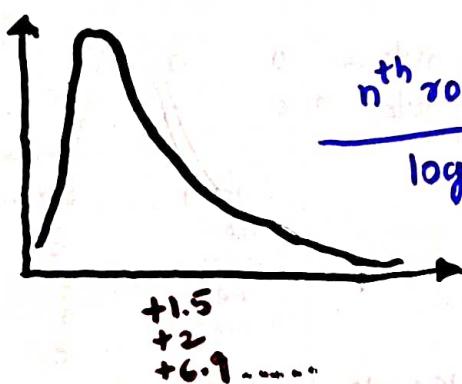
Feature Transformation

Only For
Continuous
Data

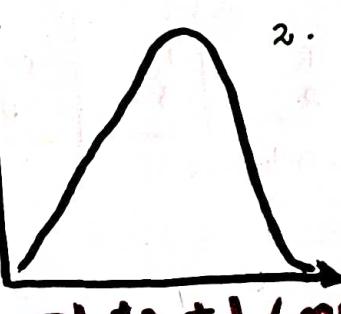
Feature processing : Transformation.

columns / variables / features (same)

Right Skewed

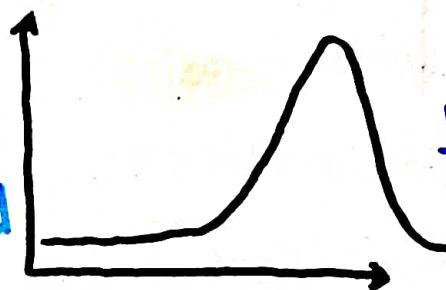


n^{th} root (or)
 $\log(x)$

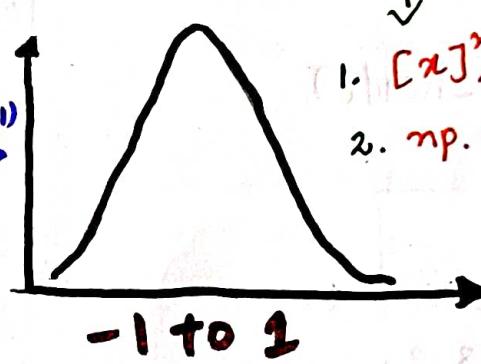


-1 to +1 (normal) Distribution

Left Skewed



n^{th} power (p)
 \exp



-1 to 1

import Libraries / packages.

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

import matplotlib inline

%matplotlib inline

import seaborn as sns

1. Root Transformations
2. log Transformations

3. Reciprocal Transformations
4. Boxcox Transformations

$x^{0.5}$
 $x^{1/2}$
 $x^{1/3}$
 $x^{1/4}$

1. $[x]^r$, $[x]^3$, $[x]^4$
2. np.log

1. $[x]^r$, $[x]^3$, $[x]^4$
2. np.exp

```
# titanic = pd.read_csv("titanic-csv")  
# titanic = pd.read_csv("titanic-csv")  
# titanic = pd.read_csv(location of csv)
```

titanic

Passenger Id	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
1	0	3	Owen Harris	male	22.0	1	0	A/52711	7.2500
2	1	1	cumming	Female	38.0	1	0		71.2833
3	1	3	Larina	Female	26.0	0	0		7.9250
4	1	1	Futrelle	Female	35.0	1	0		53.1000
889	890	1	Behr	male	26.0	0	0		30,0000
890	891	0	Doolay	male	32.0	0	0		7.7500

891 rows x 12 columns.

Cabin	Embarked
NAN	S
C85	C
NAN	S
C123	S
:	
C148	C

```
# titanic = pd.read_csv("titanic.csv", usecols=[["Fare"]])
```

titanic_head()

out

Fare
7.2500
71.2833
7.9250
53.1000
8.0500

Particular column

usecols = ["Fare"]
col. Name.

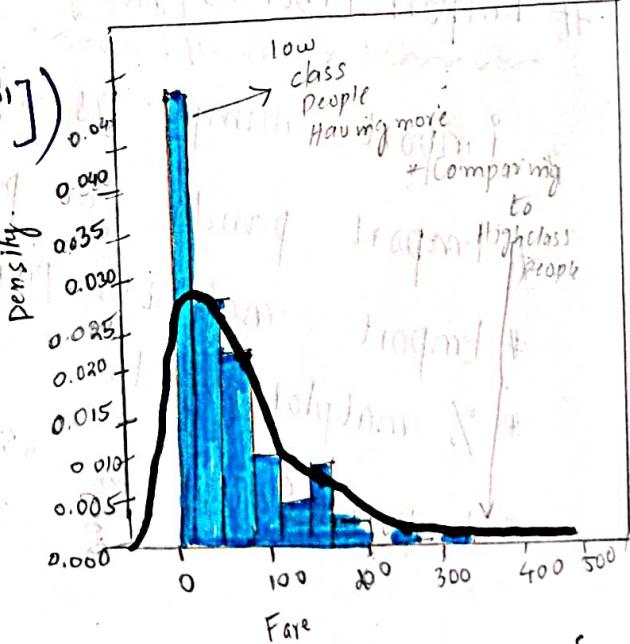
```
# sns.distplot(titanic["Fare"])
```

plt.show()

```
titanic[["Fare"]].skew()
```

out

+ 4.7873



another way

```
t - ["Fare"].hist()  
plt.show()
```

```
titanic["Fare"].skew()
```

Out :- +4.7873

Values

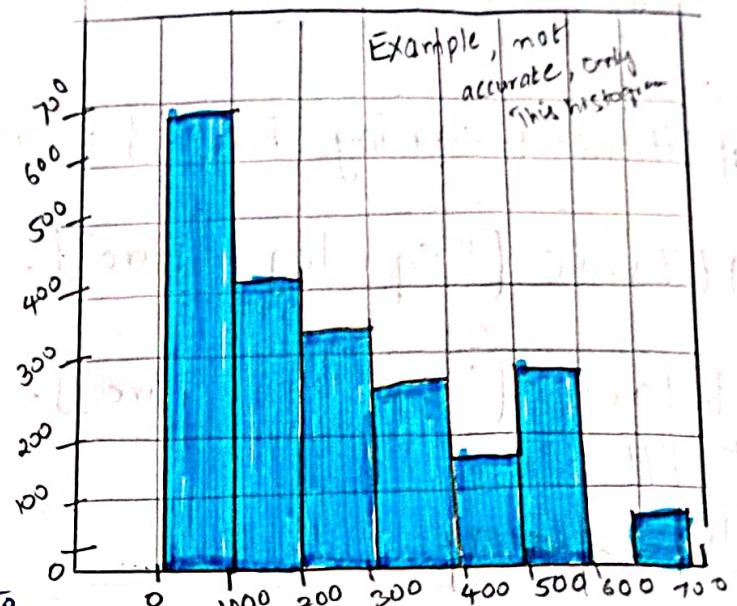
$$\text{Skew} = +0.085 \left(2^{\text{th}} \text{root} \right)$$

$$\text{Skew} = +0.5 \left(4^{\text{th}} \text{root} \right)$$

$$\text{Skew} = +0.21 \left(5^{\text{th}} \text{root} \right)$$

$$\text{Skew} = -0.95 \left(6^{\text{th}} \text{root} \right)$$

consider,
This
which is
close To
"0".



* ROOT Transformation # For Right Skewed

```
# titanic["Sqr-Fare"] = titanic["Fare"] ** (1/2)
```

```
# titanic["Sqr-Fare"].skew()
```

```
# titanic["Sqr-Fare"].hist()
```

Out :- +2.085 More Than +1
or -1

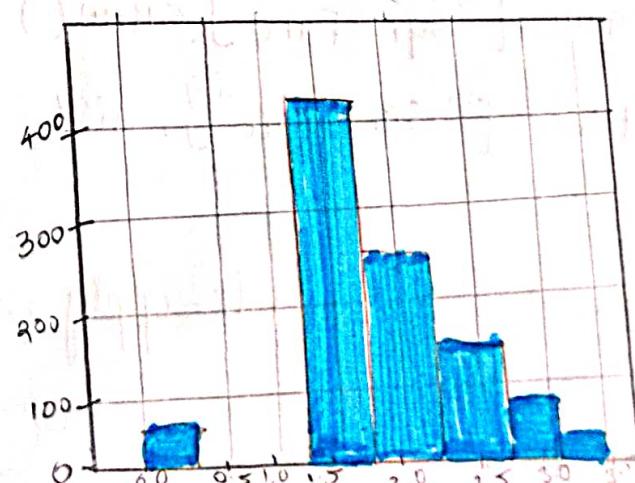
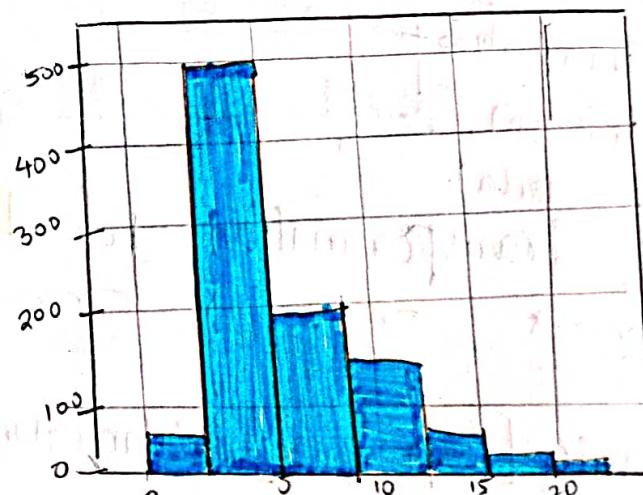
* Try with [$\sqrt[5]{}$], [$\sqrt[6]{}$], ...

```
# titanic["Sqr-Fare"] = titanic["Fare"] ** (1/5)
```

```
# titanic["Sqr-Fare"].hist()
```

```
# titanic["Sqr-Fare"].skew()
```

Out :- -0.212676 Less Than
-1 to +1
Final



→ Log Transformation

Q. Is This mandatory To write Every time?
 - No, Only if we "0" in The Data.

#titanic["Sqr log - "Fare"] = $n \cdot \log(titanic["Fare"] + 0.01)$

#titanic["Sqr log - "Fare"].skew() Why?

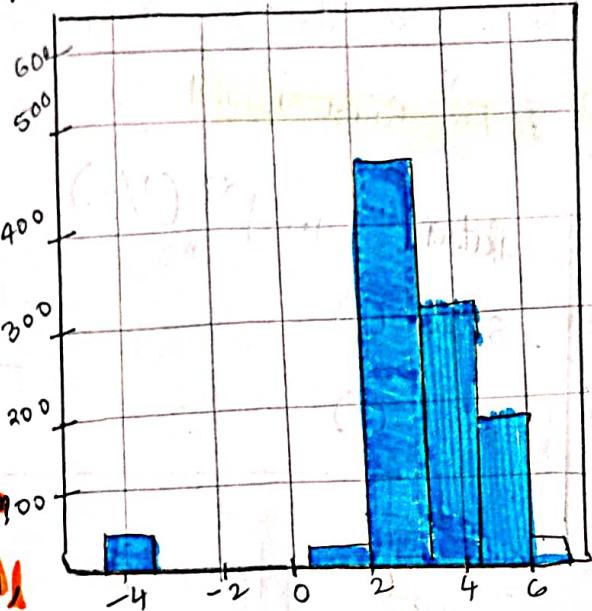
#titanic["Sqr log - "Fare"].hist()

plt.show()

Out

-2.41004

Here, The Value, In This Have more, So we don't use log Transformation For This (0) Data)



Because In Data "Fare" column, we have some records "0" values. So,

If we calculate with log. it give a (infinity) value, so, we add +0.01 to "0" column.

we use "n" power

Root Transformation For (left skewed)
 fare data is not left skewed, But To understand, Left skewed

titanic["Sqr-Fare"] = titanic["Fare"]^{**}(2)

it can be any value

titanic["Sqr-Fare"].skew()

(2)(3)(4) ...

titanic["Sqr-Fare"].hist()

* When The Data is left skewed.

Out:

* We apply n^{th} power To The Data (or) column.

* Exponential

`titanic["sqr-exp Fare"] = np.exp(titanic["Fare"])`

`titanic["sqr-exp Fare"].skew()`

`titanic["sqr-exp Fare"].hist()`

`plt.show`

another method
For Left Skewed
Distribution

* Evaluates e^x For Each Element in The given Input.

* Reciprocal Transformation

`titanic["Rec-Fare"] = 1 / [titanic["Fare"] + 0.01]`

`titanic["Rec-Fare"].skew()`

`titanic["Rec-Fare"].hist()`

`plt.show()`

"Reciprocal" is that we divide Every Value with $(1) / [\text{column}] + 0.01$ of The data to Eliminate "0" value In The Data.

Box Cox Transformation

Fare
7250
71283
-
-
-

$$\rightarrow \frac{7250 - 1}{71283 - 1}$$

$$T[x] = \frac{x^\lambda - 1}{\lambda}$$

$\therefore \lambda$ [lambda] Varies From -5 to +

\therefore Where x is The response Variable and

\Rightarrow In This Transformation, all Values of

" λ " are Considered and The optimal

Value for a given Variables are Selected.

```

from scipy import stats
# stats.boxcox(titanic['Fare'])

```

If we don't write This
Value
Error: Data must be two
+0.01 (or) +1 (User defined)

[Out]: (array([2.3844558, 6.43357655, 2.512739737,
2.607914, 5.7648427, 4.06768781...]),
[0.01])

From help

we see Two return values

1. array values (boxcox)
2. optimal Lambda parameter (max log)

which ever gives best result

Fitting data, Fitting Lambda

```

# titanic["Fare_boxcox"], param = stats.boxcox

```

Multiple Variable assignment
if a, b stores

```

# print ("λ", param) (or) param

```

[Out]: -0.09

```

# titanic["Fare_boxcox"].skew()

```

[Out]: -0.04

lowest values
while Compare with
other Transformations

Code

```

# titanic["Fare_boxcox"], param = stats.boxcox(titanic.Fare + 1)

```

```

# titanic["Fare_boxcox"].hist()

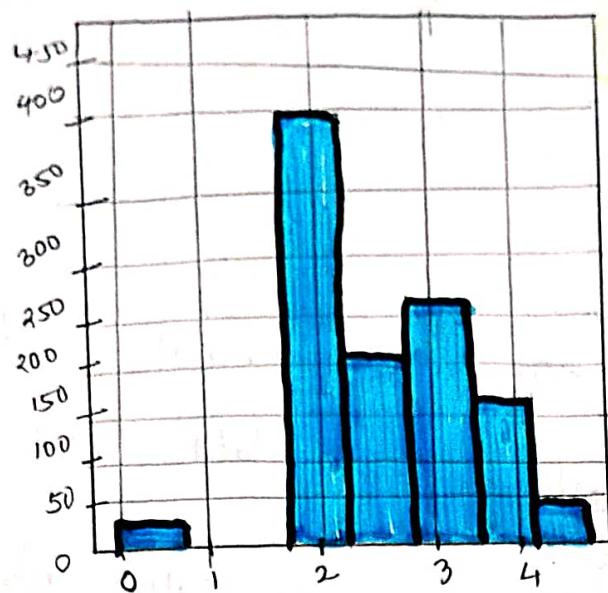
```

```

# titanic["Fare_boxcox"].skew()

```

[Out]: -0.04



~~inf/~~
4/4/22

5:30pm.

5/4/22
11:00 AM

2 \Rightarrow Feature Scaling \Rightarrow Applicable only for **Continuous Data**

Feature Scaling?

Ans:- it refers (or) Techniques used to **normalize** the ranges of **Independent Variables** in our data. (or) **input variables**

* The methods to set the features value range within a **similar scale**.

* Variables with bigger magnitude / larger value range dominate over those with smaller magnitude / value range.

Ex:- 10,000,000
if we take
10
10

\rightarrow Both are equal importance

* Scale of the features is an important consideration when building **Machine Learning models**.

* Feature scaling is generally the **last step** in the data **pre processing pipeline**, performed just before **Training the machine learning algorithms**.

5. Gradient descent **converges faster** when **Features** are on **similar scales**

Example :-

C.C	Mileage
1600	14
1800	15
2200	16.5
2100	18
2600	22.5
1800	19.4
1900	25.4

Continuous Data

* Feature Scaling *

of a car

a. When we ask Machine which is important In The both columns?

* Machine is Giving Importance To C.C

Because it is Having Large Values.

* But, We know both are important
So, we have Train's model in ^{Initially} Such that,

it should consider Both The columns Similar

So, we Reduce The Value

$$= \frac{800}{1000} = \frac{80}{100} = \frac{8}{10} = \frac{4}{5}$$

(Every One is same)

We are scaling down
Dividing Every value with

"10" In this case.

* To Reduce The Value.

* Feature Scaling

Importance in
Some ML Algorithms

1. Linear Regression & logistic Regression

* The regression coefficients of linear models are directly influenced by scale of the variable.

2. Support Vector Machines :-

* Feature scaling helps decrease the time to find support vectors for SVM's

3. K-means clustering

* Euclidean distances are sensitive to feature magnitude

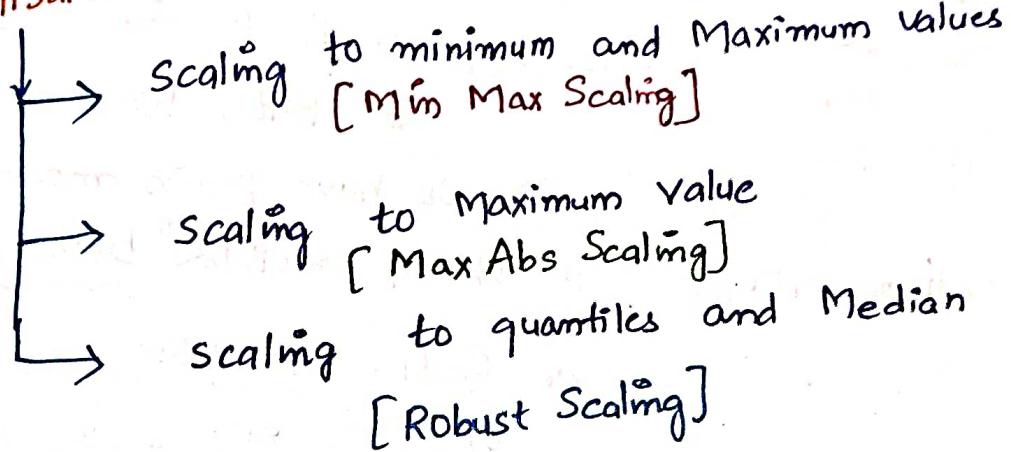
4. Principal Component analysis (PCA) :-

* PCA requires the features to be centered at "0".

* Various Feature Scaling Techniques :-

1. Standardisation

2. Normalisation Technique



import pandas as pd

import matplotlib.pyplot as plt

%matplotlib inline

→ Using (titanic Data Set)

code:
titanic = pd.read_csv("titanic.csv", usecols=["Age"])

titanic.head()

Out

Age
22.0
38.0
26.0
35.0
35.0

check Null values in "Age" column.

titanic["Age"].isnull().sum()

[Out] Age 177 → Null values

Replacing with "Median Value" for Null values

titanic["Age"] = titanic["Age"].fillna(titanic["Age"].median(),
inplace = True)

titanic["Age"].isnull().sum()

[Out] : Age 0 → zero Null values



* Standardisation

Converting Each Value To Zscore

Ex :-

X	$\frac{x-\mu}{\sigma} = Z \text{ score}$
23	$\frac{23-60}{2} = -18.5$
45	$\frac{45-60}{2} = -7.5$
67	$\frac{67-60}{2} = 3.5$
89	$\frac{89-60}{2} = 14.5$
78	$\frac{78-60}{2} = 9$

These five values converting To Zscores
is called as "standard" scaling

$$\text{Avg}/\text{mean} = 60.$$

$$(\text{std}) \quad \sigma = \sqrt{\frac{\sum(x_i - \bar{x})^2}{n}}$$

Code :-

```
# From sklearn.preprocessing import StandardScaler (sc)
# sc = StandardScaler () # call The Function (short cut)
# main (don't forget)
# titanic ["Age-SC"] = sc. fittransform (titanic [["Age"]])
# titanic ["Age-SC"] = # creating new column
```

titanic ["Age-SC"] = # fit-transform stores in

Out :-

	Age-SC
0	-0.5651
1	0.6638
2	-0.2583
3	0.43312
:	
890	0.20272

Here Every Value is converted To Zscore.

$$\left[\frac{x - \mu}{\sigma} \right] = Z \text{ score.}$$

EX:- Age [22] # First record $x = 22$

$$\# \text{titanic}.mean[\mu] = 29.361 \quad \# \text{titanic}.std[\sigma] = 13.01$$

Original Data "Age-SC" $\rightarrow [-0.5651]$

titanic

Out :-

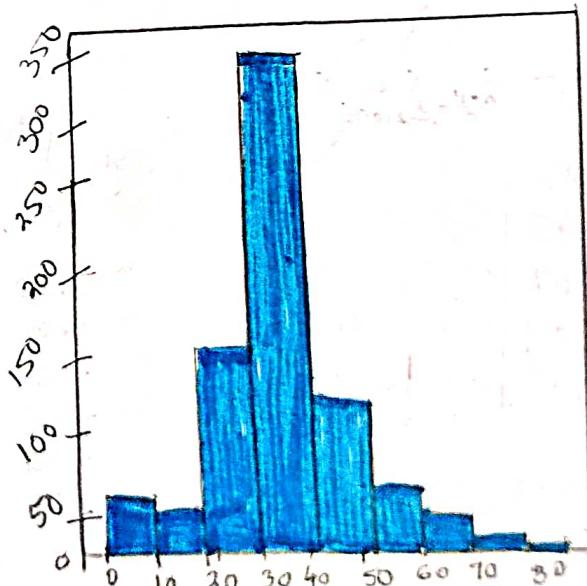
	Age	Age-SC
0	22.0	-0.5651
1	38.0	0.6638
2	25.0	-0.2583
:		
890	26.0	-0.2583
890	32.0	0.20272

Converted To Z-Score

histogram of Original ["Age"]

titanic ["Age"].hist()

plt. show

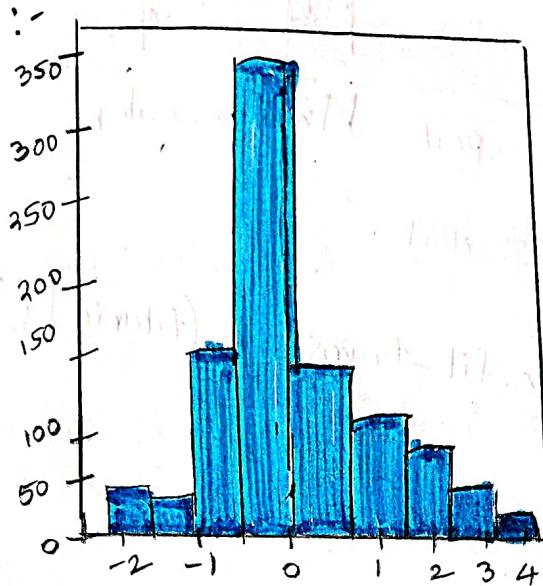


Histogram of Converted "Age"

titanic[["Age_sc"]].hist()

plt.show()

Out :-



Here, we see, Histogram is not going to change, But, Values does.

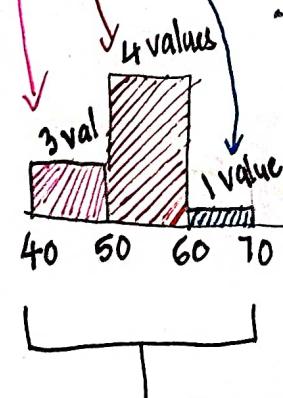
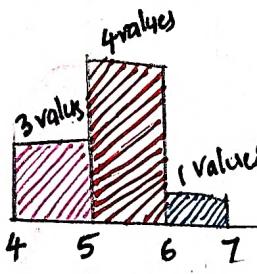
Example :-

* How ? Histogram Remains Same, But changed Values

Ans :-

Ex :-

X	Changed X
40	4
50	5
42	4.2
43	4.3
54	5.4
58	5.8
59	5.9
62	6.2



Here only 'X' values Only changes But Histogram Remains Same.

II # Normalisation

Ex:-

X	
65	$\frac{65-34}{53}$
45	$\frac{45-34}{53}$
34	$\frac{34-34}{53} = 0$
67	$\frac{67-34}{53}$
67	$\frac{67-34}{53}$
87	$\frac{87-34}{53} = 1$

$$\frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

* MinMax Scaling

→ MinMax Scaling Scales The Values between "0" to "1" (Every Thing is one value) No -ve value

Code: From sklearn.preprocessing import

min_max = MinMaxScaler() # Shortcut

titanic["Age_mm"] = min_max.fit_transform(titanic[["Age"]])

Creating new column

titanic (calling/show data)

Out :-

	Age	Age_mm
0	22.0	-0.563
1	38.0	0.448
2	26.0	0.259
3		0.32
4		0.32
5	26.0	-0.25
6	32.0	0.20

* Robust Scaling

$$\frac{x - \text{X}_{\text{median}}}{\text{IQR}}$$

$$\therefore \text{IQR} = Q_3 - Q_1$$

Code :

From sklearn.preprocessing import RobustScaler

rs = RobustScaler() → shortcut()

titanic["Age-rs"] = rs.fit_transform(titanic[["Age"]])

titanic["Age-rs"] (or) = RobustScaler().fit_transform(...)

If not using any shortcut like rs = RobustScaler()

titanic

Out

Age	Age-SC	Age-MM	Age-RS
22.0	0.565	0.461	0.4615
38.0	0.66	0.769	0.769
26.0	0.25	-0.153	-0.153
...			
26.0	-0.25	0.153	0.153
39.0	0.2	0.396	0.3076

* Max Abs Scaling

$$\frac{x}{x_{\max}}$$

Ex:-

x	
68	68/80
44	44/80
52	52/80
69	69/80
80	80/80 = 1

Max value

code

From sklearn.preprocessing import

MaxAbsScaler

mas = MaxAbsScaler()

titanic["Age-mas"] = mas.fit_transform(titanic[["Age"]])

titanic

Out :-

	Age	Age - sc	Age - mm	Age - rs	Age - mas
0	0.2750
1	0.4950
2	0.3250
...
889	0.3250
890	0.4000

891 x 5 col

Original Variable
Standard Scaler
Min max Scaler
Robust Scaler
Max Abs Scaler

✓
5/04/22
4:58 pm.

Example :- Discretization

Continuous Data

12.5
22.8
25.2
18.3
12.4
15.8
23.3

⇒ Mileage of a Car (Discrete, categorical) Data

$\leq 15 \rightarrow$ low. mileage
 $15 - 20 \rightarrow$ Avg. mileage
 $20 - 25 \rightarrow$ good. mileage
 $> 25 \rightarrow$ very good

By our own choice To divide the into parts