**2** → ==Feature Scaling== ⟹ Applicable only For ==Continous 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.
same scale

* Variables with bigger magnitude / Larger value range dominate over those with Smaller magnitude /value range.

EX:- 10,000,00 , x 10
So we Take
10 — 10 → Both are Equal Import
10

* 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 :-

of a car

| C.C | Milage |
|------|--------|
| 1600 | 14 |
| 1800 | 15 |
| 2200 | 16.5 |
| 2100 | 18 |
| 2600 | 22.5 |
| 1800 | 19.4 |
| 1900 | 25.4 |

Continous Data

# * Feature Scaling *

q. 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 importa So, we have Trains model in → Intially Such that, it should Consider" Both The columns Similar

So, we Reduce The Value with 200

$$= \frac{800}{1000} = \frac{80}{100} = \frac{8}{10} = \frac{4}{5}$$ (Every One i Same)

with 100

We are scaling down Dividing Every Value wit "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
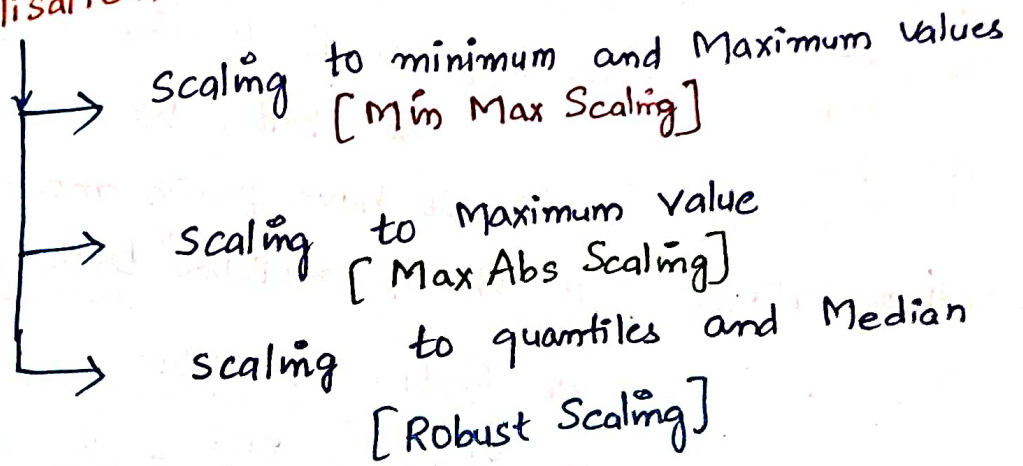
4. Principal Component analysis (PCA) :-

* PCA requires The Features to be Centered at "0".

* **Various Feature Scaling Techniques :-**

1. Standardisation
2. Normalisation Technique

→ Scaling to minimum and Maximum values [Min Max Scaling]

→ Scaling to Maximum value [Max Abs Scaling]

→ Scaling to quantiles and Median [Robust Scaling]

```
# import pandas as pd
# import matplotlib.pyplot as plt
%. matplotlib inlime.
```

→ 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 . isnull(). sum()

Out        Age      177 ⟶ Null values

# Replacing with "Median value" For Null values

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

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

Out :       Age    0 ⟶ Zero Null Values

---

**I** * Standardisation



Converting Each Value To Zscore

Ex :-

| X | $\frac{x-\mu}{\sigma}$ = Zscore |
|---|---|
| 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 |
| 18 | $\frac{78-60}{2}$ = 9 |

These Five values converting To Zscores
is called as "Standard" Scaling

Avg/mean = 60.

(Std) $\sigma$ = 2 EX?

**Code :-**

From sklearn.Preprocessing import Standard Scaler (sc)

\# sc = Standard Scaler ()       # call The Function (Short cut)

     main (don't Forget)

\# titanic ["Age_sc"] = sc. fit_tranform (titanic [["Age"]

       #creating new column      # fit_transform

           Calculation     Converting The value

\# titanic ["Age_sc"]      STores in

**Out :-**

| Age_sc |
|--------|
| 0   - 0.5657 |
| 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

\# titanic.mean $\left[\frac{X - \mu}{\sigma}\right] = \left[\frac{22 - 29.36}{13.01}\right]$

N = 29.361

\# titanic.std ()

$\sigma = 13.01$

original Data

"Age_sc" → [- 0.5657]

Converted To Z_Score

\# titanic

**Out**

| | Age | Age_sc |
|---|------|--------|
| 0 | 22.0 | - 0.5657 |
| 1 | 38.0 | - 0.6638 |
| 2 | 25.0 | - 0.2583 |
| . | . | . |
| 889 | 26.0 | -0.2583 |
| 890 | 32.0 | 0.20272 |

\# histogram of Original ["Age"]
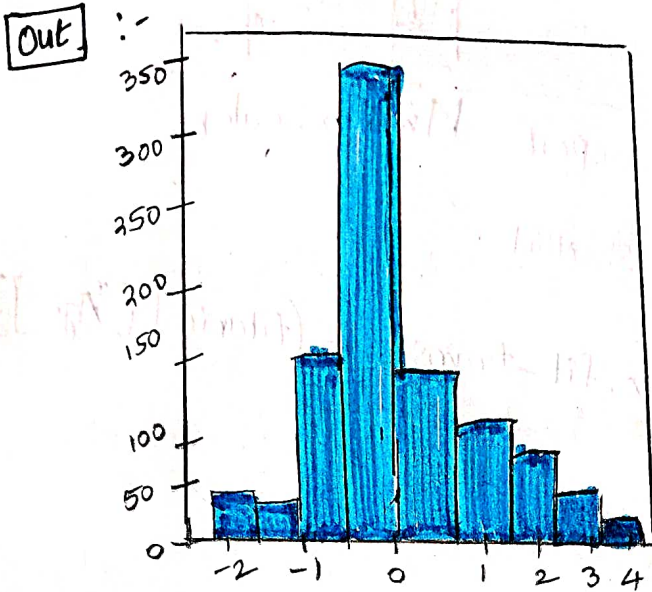
\# titanic ["Age"]. hist ()

    plt. show

# Histogram of Converted "Age"

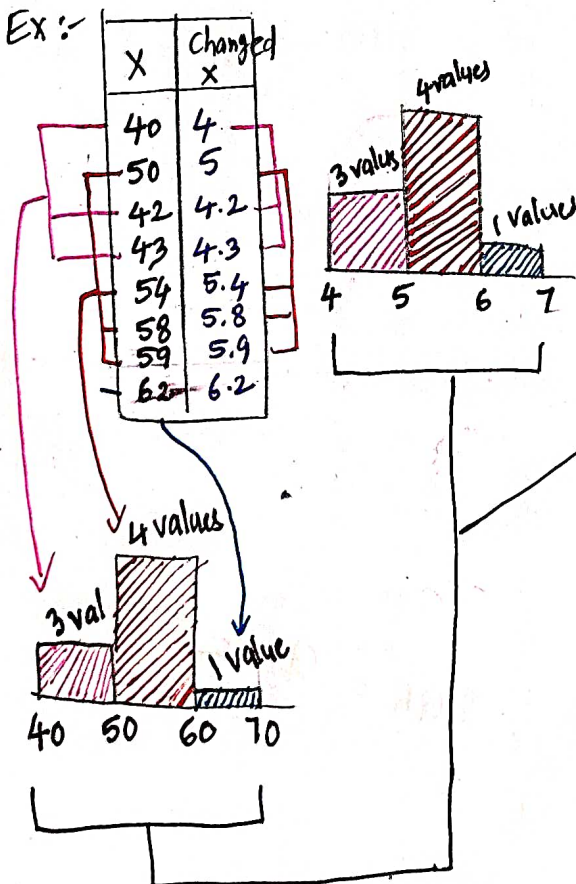# titanic ["Age_sc"] . hist ()

# 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 |



3 values   4 values   1 value
4    5    6    7

4 values
3 val        1 value
40  50  60  70

# Here only "x" values only changes But Histogram Remains Same.

## II # Normalisation

### * MinMax Scaling

⇒ MinMax Scaling Scales, The Values between "0" to "1" (Every Thing is (+ve) value) No -ve value

Ex:-

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

min. value ← 34

Max. value ← 87

$$\frac{X - X_{min}}{X_{max} - X_{min}}$$

Code: From sklearn. Preprocessing import MinMax Scaler

# min_max = MinMax Scaler () # Shortcut

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

↳ creating new column

fit_transform

# titanic (calling/show Data)

Out :-

| | Age | Age_sc | Age_mm |
|---|---|---|---|
| 0 | 22.0 | -0.565 | 0.271 |
| 1 | 38.0 | 0.457 | 0.472 |
| 2 | 26.0 | 0.258 | 0.32 |
| ... | | | |
| 889 | 26.0 | -0.25 | 0.32 |
| 890 | 32.0 | 0.20 | 0.39 |

### * Robust Scaling

$$\frac{X - X_{median}}{IQR}$$

∴ $IQR = Q_3 - Q_1$

## Code :

From sklearn.preproce import RobustScaler (using)

\# rs = RobustScaler() ⟶ shorcut()

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

titanic ["Age-rs] $\overset{(or)}{=}$ 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 |
| 32.0 | 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

Original Variable
StanDard Scaler
Min max scaler
Robust scaler
Max Abs Scaler

| | Age | Age -sc | Age -mm | Age -rs | Age -mas |
|---|---|---|---|---|---|
| 0 | . | . | . | . | 0.2750 |
| 1 | | | | . | 0.4150 |
| 2 | . | . | . | . | 0.3250 |
| : | | | | | : |
| 889 | . | . | . | . | 0.3250 |
| 890 | . | . | . | . | 0.4000 |

891 x 5 col

5/04/22
4:58 pm.

Example :- Discretization

Continous Data ⟹ Milage of a Car (Discrete . categorical) Data

< 15 → low . milage

15 - 20 → Avg . milage

20 - 25 → good. milage

> 25 → Very good

By our Own choice To Divide the into parts

| |
|---|
| 12.5 |
| 22.8 |
| 25.2 |
| 18.3 |
| 12.4 |
| 15.8 |
| 23.3 |