

Feature Scaling → Normalisation

Normalisation → technique often → As a part of → Data preparation

↓
goal is →

for Machine learning ↗

To change a value of numerical columns

↓
make it one scale

↓
without distorting
difference

↙ ↘
on less information range

Type of Normalisation

- Min Max Scaling
- Mean normalisation
- Max Absolute Scaling
- Robust Scaling

There are more Scaling function in Sciklearn

↳ Most important is → Min Max Scaling

Min Max Scaling \Rightarrow

Weight
130
67
81
61
32
54

Normalize
Min Max Scaling



$$x_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

$$x_i = 130$$

$$\textcircled{1} = \frac{130 - 32}{130 - 32} = 1$$

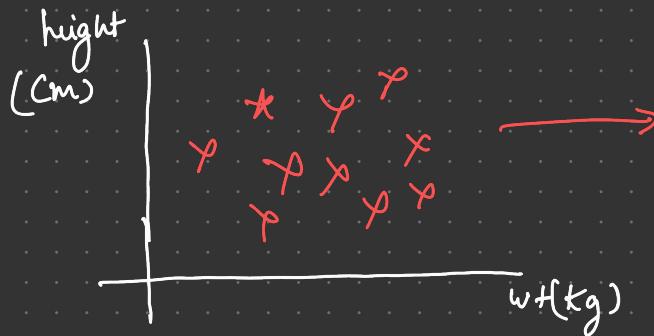
$$\textcircled{2} = \frac{67 - 32}{130 - 32} =$$

By this we
can see
that
Range will
be
 $[0, 1]$

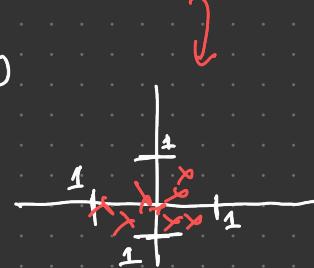
$$\left[\frac{130 - 32}{130 - 32} \right] \text{ max}$$

$$\left[\frac{32 - 32}{130 - 32} \right] = 0 \text{ min}$$

Biometrical Way



Applying min max scaling



We are putting all values

$[0, 1]$

Code form

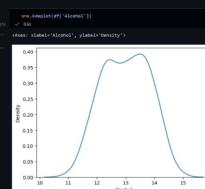
```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

df = pd.read_csv('wine_data.csv', header=None, usecols=[0,1,2])
df.columns=['ClassLabel', 'Alcohol', 'MalicAcid']

The pd.read_csv() function is a Pandas function that is used to read CSV files into DataFrames.
• The header=None parameter tells Pandas to ignore the first row of the CSV file. This is useful if the first row of the CSV file does not contain column names.
• The usecols=[0,1,2] parameter tells Pandas to only read the first three columns of the CSV file. This is useful if you only need to use a subset of the columns in the CSV file.
• The DataFrame is a Pandas data structure that is used to store tabular data. It is similar to a spreadsheet.
• The columns attribute of a DataFrame is a list of the names of the columns in the DataFrame.
• The rename() method of a Series or DataFrame is used to rename the columns or index labels.

Data Visualization

Kdeplot is a visualization technique used to estimate and plot the probability density function (PDF) of a data set



In Data Visualization \rightarrow KDE Plot (seaborn library)

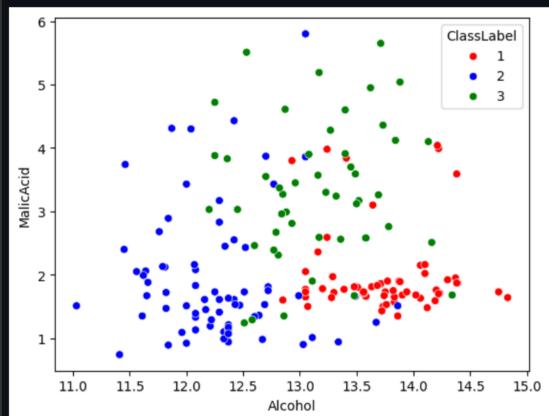
Help to plot \rightarrow PDF (Probability density function)

↳ This Doesn't Make any assumption about the
↳ Underlying distribution of the Data

By this Versatile tool for → Exploring Data

↳ Especially Data → is not
Normaly Distributed

```
color_dist={1:'red',2:'green',3:'blue'}  
sns.scatterplot(x=df['Alcohol'], y=df['MalicAcid'], hue = df['ClassLabel'], palette=color_dist)  
✓ 0.3s  
Axes: xlabel='Alcohol', ylabel='MalicAcid'
```



→ This is for Data
Preprocessing

```
from sklearn.model_selection import train_test_split  
X_train,X_test,Y_train,Y_test = train_test_split(df.drop('ClassLabel',axis=1),  
                                                df['ClassLabel'],test_size=0.3, random_state=0)  
✓ 0.0s
```

→ Splitting Data
for test and
train

```
X_train.shape,X_test.shape  
✓ 0.0s
```

```
((124, 2), (54, 2))  
  
from sklearn.preprocessing import MinMaxScaler  
scaler = MinMaxScaler()  
#fit the scaler to the train set, it will learn the parameters  
scaler.fit(X_train)  
  
#transform train and test sets  
X_train_Scaled = scaler.transform(X_train)  
X_test_Scaled = scaler.transform(X_test)
```

```
X_train_Scaled = pd.DataFrame(X_train_Scaled, columns=X_train.columns)  
X_test_Scaled = pd.DataFrame(X_test_Scaled, columns=X_test.columns)  
[19] ✓ 0.0s
```

From
Min Max Scaler ↗
we transformed

• When ever we fit
it will get fit
in training Data ↗

transform train
to bhi Karuge
but to bhi
Karuge

So when we use Scilearn library to mat
Change it into Numpy Array

To handle that we converted
as
Dataprame.

```
np.round(X_train.describe(),1)
```

0.1s

	Alcohol	MalicAcid
count	124.0	124.0
mean	13.0	2.4
std	0.8	1.1
min	11.0	0.9
25%	12.4	1.6
50%	13.0	1.9
75%	13.6	3.2
max	14.8	5.6


```
np.round(X_train_Scaled.describe(),1)
```

0.0s

	Alcohol	MalicAcid
count	124.0	124.0
mean	0.5	0.3
std	0.2	0.2
min	0.0	0.0
25%	0.4	0.2
50%	0.5	0.2
75%	0.7	0.5
max	1.0	1.0

After Doing min max Value

it will get zero

Mean → standard deviation Behave differently

Mostly min and max work like
that



May be in some distribution mai thoda sa change a sakte hai ↴

But its ignored

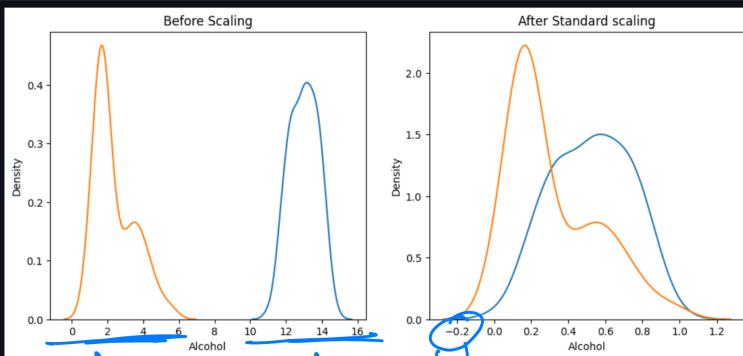
Because Data ka jo form hota hai usko retain kar sakte hai mean Max Scalar.

```
fig, (ax1,ax2) = plt.subplots(ncols= 2, figsize=(12,5))
#Before Scaling
ax1.set_title('Before Scaling')
sns.kdeplot(X_train['Alcohol'], ax=ax1)
sns.kdeplot(X_train['MalicAcid'], ax=ax1)

#After Scaling
ax2.set_title("After Standard scaling")
sns.kdeplot(X_train_Scaled['Alcohol'], ax= ax2)
sns.kdeplot(X_train_Scaled['MalicAcid'], ax=ax2)

plt.show()
```

Python



Mean Alcohol
was in the
Range $0-6$

Malic Acid $\rightarrow 10-16$

After Doing
Mean
Max Scaling It Came
in Same
Range.

There is Some Sort going
Negatively

As we are doing KDE plot

It take Inference
Statistically from

↳ There is a chance of changing → Shape after
Distribution
↳

After Distributions that isn't necessary
It should maintain
Same Shape

↳ we have issue in Outlier → Outlier bhi Squuz hojata
hai
In min max Scaling

Mean Normalization :-

↳ $X_i' = \frac{X_i - X_{\text{mean}}}{X_{\text{max}} - X_{\text{min}}}$ → Here we mean Centric
doing it

↳ -1 to 1 Range is Value

Mean \leq Value \rightarrow +ve number

Mean \geq Value \rightarrow -ve number

↳ get Rarely used \rightarrow There isn't any Class in Sciklearn

↳ get only usefull in \rightarrow Centred data

Even though People
use
Standardization

Max Abs Scaling \rightarrow

Max Absolute Scaling \rightarrow

$$x'_i = \frac{x_i}{|x_{\max}|}$$

We have a Class in Sciklearn \rightarrow MaxAbsScaler

↳

We can use directly

If get utilized where \rightarrow we have

+ Sparse data \rightarrow "0's" Bhaut hai

Robust Scaling \Rightarrow

$$x_i' = \frac{x_i - \text{median}}{\text{IQR}}$$

IQR $\left\{ \begin{matrix} 75^{\text{th}} \text{ Percentile} \\ \text{Value} \end{matrix} \right\}$

It is Robust \rightarrow bcoz \rightarrow if Data have lot of Outlier

\hookrightarrow we can't tell before which will do good work with which Data.

By Using \rightarrow tabular Data \rightarrow we can't say

Normalization

vs

Standardization

1) Is feature Scaling needed? \rightarrow Decision Tree not

needed
for this we need to understand Algo.

2) Mostly Problem will solved doing \rightarrow Standardization

that is why Standard Scalar Class host using

\downarrow

Most of the Time better result

④ Normalization \rightarrow Min Max Scaler

When we are aware of Min and Max for Example

$\xrightarrow{\text{CNN}}$ \leftarrow Image processing
Pixel have

$0 - 255$ Range Min Max Scaling

- ↳ Outlier have \rightarrow Robust Scaling
- ↳ Don't know \rightarrow Standard Scaling
- ↳ Min Max \rightarrow Min Max Scaling
- ↳ Spars Values? \rightarrow Max Abs Scaling