

## Simple List Imputer

```
[2]: from sklearn.impute import SimpleImputer
import numpy as np

data = np.array([
    [1, 2, np.nan],
    [5, 8, 9],
    [4, np.nan, 0],
    [0, np.nan, 0]
])

imputer = SimpleImputer(strategy='mean', fill_value=0)

imputed_data = imputer.fit_transform(data)

print(data)
print(imputed_data)

[[ 1.  2. nan]
 [ 5.  8.  9.]
 [ 4. nan  0.]
 [ 0. nan  0.]]
[[ 1.  2.  6.]
 [ 5.  8.  9.]
 [ 3.  4.  7.5]
 [ 0.  0.  0.]]
```

## KNN Imputer

```
[3]: from sklearn.impute import KNNImputer
import numpy as np

data = np.array([
    [1, 2, np.nan],
    [5, 8, 9],
    [0, np.nan, 0],
    [4, 0, 0]
])

knn_imputer = KNNImputer(n_neighbors=3)

imputed_data = knn_imputer.fit_transform(data)

print(data)
print(imputed_data)

[[ 1.  2. nan]
 [ 5.  8.  9.]
 [ 0. nan  0.]
 [ 9.  8.  9.]]
[[ 1.  2.  6.5]
 [ 5.  8.  9.]
 [ 5.5  0.  0.]
 [ 9.  8.  9.]]
```

## Mean Removal

```
[31]: import numpy as np
from sklearn import preprocessing

input_data = np.array([[1, 1.5, 3, 0.4], [0.5, 1.3, 4.1], [1.2, 1.4, 2.9, 0.7]])
data_standardized = preprocessing.scale(input_data)

print("\n Mean = ", data_standardized.mean(axis = 0))
print("\n Std deviation = ", data_standardized.std(axis = 0))

Mean = [ 0.38115126 17 -3.70074320 17  0.00000000e+00 -1.85872710e-17]
Std deviation = [ 1.  1.  1.  1.]
```

## Min-Max Scale

```
[33]: import numpy as np
from sklearn import preprocessing

input_data = np.array([[1, 1.5, 3, 0.4], [0.5, 1.3, 4.1], [1.2, 1.4, 2.9, 0.7]])

data_scaler = preprocessing.MinMaxScaler(feature_range = (0,1))
data_scaled = data_scaler.fit_transform(input_data)

print("\n Min Max Scaled Data = ",data_scaled)

Min Max Scaled Data = [[ 1.  0.  1.  0.  1.  0.  1.]
 [ 0.  1.  0.2718644 1.  1.]
 [ 0.33333333 0.66666667 0.  0.2  1]]
```

## Normalization 1

```
[38]: import numpy as np
from sklearn import preprocessing

input_data = np.array([[1, 1.5, 3, 0.4], [0.5, 1.3, 4.1], [1.2, 1.4, 2.9, 0.7]])

data_normalized = preprocessing.normalize(input_data, norm = 'l2')

print("\n l1 Normalized Data = ",data_normalized)

l1 Normalized Data = [[ 0.21302734  0.16703627  0.21302734  0.4661865]
 [ 0.  0.3574286 -0.3547619  0.48809526]
 [ 0.4952381  0.21304762 -0.2763968  0.4952381]]
```

## Normalization 2

```
[21]: import numpy as np
from sklearn import preprocessing

input_data = np.array([[1, 1.5, 3, 0.4], [0.5, 1.3, 4.1], [1.2, 1.4, 2.9, 0.7]])

data_normalized = preprocessing.normalize(input_data, norm = 'l2')

print("\n l2 Normalized Data = ",data_normalized)

l2 Normalized Data = [[ 0.38345117  0.16172558  0.38345117  0.81862016]
 [ 0.  0.57267755 -0.54790607  0.78181010]
 [ 0.87357808  0.39623866 -0.50376101 -0.76688891]]
```

## Binarizing

```
[34]: import numpy as np
from sklearn.preprocessing import Binarizer

ages = np.array([15], [22], [18], [39], [16], [45])

binarizer = Binarizer(threshold = 18)
Binary_Ages = Binarizer.fit_transform(ages)

print(" Original ages: \n",ages)
print(" Binarized ages: \n",Binary_Ages)

Original ages:
[[15]
 [22]
 [18]
 [39]
 [16]
 [45]]
Binarized ages:
[[0]
 [1]
 [0]
 [1]
 [0]
 [1]]

Gender Data Frame :
```

```
Gender      Name
0      F  Krishna
1      M   Radha
2      M  Krishna
3      F   Kiana
4      M   Dhoni
5      F    Sara
6      M   Rajesh
7      F    Lata
8      F   Anuska
9      M    Ratan
```

```
[27]: import pandas as pd
from sklearn import preprocessing

my_data = {
    "Gender": ['F', 'M', 'M', 'F', 'M', 'F', 'M', 'F', 'M', 'M'],
    "Name": ['Krishna', 'Radha', 'Krishna', 'Kiana', 'Dhoni', 'Sara', 'Rajesh', 'Lata', 'Anuska', 'Ratan']
}

bik = pd.DataFrame(my_data)
print("Gender Data Frame : \n")
print(bik)

my_label = preprocessing.LabelEncoder()
bik["Gender"] = my_label.fit_transform(bik["Gender"])

print(bik["Gender"].unique())
print("Data frame after Label encoding: \n")
print(bik)

Gender Data Frame :
```

```
Gender      Name
0      0  Krishna
1      1   Radha
2      1  Krishna
3      0   Kiana
4      1   Dhoni
5      0    Sara
6      1   Rajesh
7      1    Lata
8      0   Anuska
9      1    Ratan
```

## Label Encoding using Category Codes

```
[30]: import pandas as pd

my_data = {
    "Gender": ['F', 'M', 'M', 'F', 'M', 'F', 'M', 'F', 'M', 'M'],
    "Name": ['Krishna', 'Radha', 'Krishna', 'Kiana', 'Dhoni', 'Sara', 'Rajesh', 'Lata', 'Anuska', 'Ratan']
}

bik = pd.DataFrame(my_data)
print(bik.dtypes)

Gender      object
Name        object
dtype: object

[37]: bik["Gender"] = bik["Gender"].astype("category")
print(bik.dtypes)

Gender      category
Name        object
dtype: object

[34]: import pandas as pd

my_data = {
    "Gender": ['F', 'M', 'M', 'F', 'M', 'F', 'M', 'F', 'M', 'M'],
    "Name": ['Krishna', 'Radha', 'Krishna', 'Kiana', 'Dhoni', 'Sara', 'Rajesh', 'Lata', 'Anuska', 'Ratan']
}

bik = pd.DataFrame(my_data)
print(bik)
bik["Gender"] = bik["Gender"].astype("category")
print("\n Data frame after label encoding using category codes : \n")
bik["Gender"] = bik["Gender"].cat.codes
print(bik)

Gender Data Frame :
```

```
Gender      Name
0      0  Krishna
1      1   Radha
2      1  Krishna
3      0   Kiana
4      1   Dhoni
5      0    Sara
6      1   Rajesh
7      1    Lata
8      0   Anuska
9      1    Ratan
```

```
[43]: advertising = pd.read_csv("advertising.csv")
print(advertising.head())

TV      Radio      News Paper      Sales
0  230.1  37.8      69.2  22.1
1  44.5   39.3      89.1  16.4
2  17.2  45.5      69.3  12.0
3  151.5  41.3      38.5  16.5
4  186.8  18.8      38.4  17.9
```

```
[45]: print(advertising.info())
print("DataShape : ",advertising.shape)

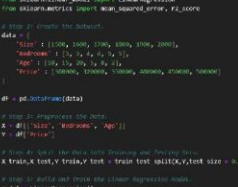
<class 'pandas.core.frame.DataFrame'>
Int64Index: 5 entries, 0 to 4
Data columns (total 4 columns):
 #   column      non-null count  dtype
---  ---
 0   TV          5 non-null      float64
 1   Radio       5 non-null      float64
 2   News Paper  5 non-null      float64
 3   Sales       5 non-null      float64
dtypes: float64(4)
memory usage: 292.0 bytes
None
DataShape: (5, 4)
```

```
[47]: # Check null values.
print("Checking for null values.")
print(advertising.isnull().sum())

Checking for null values.
TV      0
Radio   0
News Paper  0
Sales   0
dtype: float64
```

```
[49]: # Importing Library for Visualisation.
# Outlin Analysis.

fig,axs = plt.subplots(3,figsize = (7,5))
plt1 = sns.boxplot(advertising['TV'],ax = axs[0])
plt2 = sns.boxplot(advertising['Radio'],ax = axs[1])
plt3 = sns.boxplot(advertising['News Paper'],ax = axs[2])
plt.tight_layout()
```



```
[51]: # Loading the Dataset.
dataset = pd.read_csv("advertising.csv")
dataset.head()

x = dataset[['TV', 'Radio', 'News Paper']]
y = dataset[['Sales']]

[53]: # Splitting the dataset.
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size = 0.2, random_state = 100)

[57]: import pandas as pd
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score

dataset = pd.read_csv("advertising.csv")
dataset.head()
x = dataset[['TV', 'Radio', 'News Paper']]
y = dataset['Sales']
x_train,x_test,y_train, y_test = train_test_split(x,y,test_size = 0.2, random_state = 100)

print("TV Data: \n",x_train.head())
print(x_train.shape)
print(x_test.shape)
print("Sales Data: \n",y_train.head())
print(y_train.shape)
print(y_test.shape)

model = LinearRegression()
model.fit(x_train,y_train)

TV Data:
TV      Radio      News Paper
0  230.1  37.8      69.2
1  44.5   39.3      89.1
2  17.2  45.5      69.3
3  151.5  41.3      38.5
4  186.8  18.8      38.4
dtype: object
Sales Data:
TV      Radio      News Paper
0  22.1
1  16.4
2  12.0
3  16.5
4  17.9
dtype: float64
Name: Sales, dtype: float64
(4,)
(1,)
```

```
[59]: y_pred = model.predict(x_test)
mse = mean_squared_error(y_test,y_pred)
r2 = r2_score(y_test,y_pred)
print("Mean Squared Error: ",mse)
print("R2 Score : ",r2)
print("Coefficients : ",model.coef_)
print("Intercept : ",model.intercept_)

Mean Squared error: 1.7278076731812
R2 Score : nan
Coefficient : [ 0.65780954 -0.00005181  0.17621761]
Intercept : -0.652178494167467

C:\Users\raj\Desktop>python C:\Users\raj\Documents\sklearn\metrics_regression.py:123: UndefinedMetricWarning: R2 score is not well-defined with less than two samples.
  warnings.warn(msg, UndefinedMetricWarning)
```

## Case Study Implementation: Predicting House Prices.

```
[67]: # Step 1: Import libraries.
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
```