# Data Standardization

July 10, 2024

The process of standardizing the data to a common format and common range

```python
[4]: import numpy as np
     import pandas as pd
     import sklearn.datasets
     from sklearn.preprocessing import StandardScaler #function that will be use to
      ↪standardize our data set
     from sklearn.model_selection import train_test_split #it helsp to split the
      ↪data into train and test dataset
```

```python
[22]: #loading the dataset

      data = pd.read_csv('breast_cancer_data.csv')
      data.head()
```

```
[22]:          id diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean  \
      0    842302         M        17.99         10.38          122.80     1001.0
      1    842517         M        20.57         17.77          132.90     1326.0
      2  84300903         M        19.69         21.25          130.00     1203.0
      3  84348301         M        11.42         20.38           77.58      386.1
      4  84358402         M        20.29         14.34          135.10     1297.0

         smoothness_mean  compactness_mean  concavity_mean  concave points_mean  \
      0          0.11840           0.27760          0.3001              0.14710
      1          0.08474           0.07864          0.0869              0.07017
      2          0.10960           0.15990          0.1974              0.12790
      3          0.14250           0.28390          0.2414              0.10520
      4          0.10030           0.13280          0.1980              0.10430

         …  texture_worst  perimeter_worst  area_worst  smoothness_worst  \
      0  …          17.33           184.60      2019.0            0.1622
      1  …          23.41           158.80      1956.0            0.1238
      2  …          25.53           152.50      1709.0            0.1444
      3  …          26.50            98.87       567.7            0.2098
      4  …          16.67           152.20      1575.0            0.1374

         compactness_worst  concavity_worst  concave points_worst  symmetry_worst  \
      0             0.6656           0.7119                0.2654          0.4601
```

|   | | | | |
|---|---|---|---|---|
| 1 | 0.1866 | 0.2416 | 0.1860 | 0.2750 |
| 2 | 0.4245 | 0.4504 | 0.2430 | 0.3613 |
| 3 | 0.8663 | 0.6869 | 0.2575 | 0.6638 |
| 4 | 0.2050 | 0.4000 | 0.1625 | 0.2364 |

|   | fractal_dimension_worst | Unnamed: 32 |
|---|---|---|
| 0 | 0.11890 | NaN |
| 1 | 0.08902 | NaN |
| 2 | 0.08758 | NaN |
| 3 | 0.17300 | NaN |
| 4 | 0.07678 | NaN |

[5 rows x 33 columns]

```
[14]: dataset = sklearn.datasets.load_breast_cancer()
```

```
[15]: print(dataset)
```

```
{'data': array([[1.799e+01, 1.038e+01, 1.228e+02, …, 2.654e-01, 4.601e-01,
        1.189e-01],
       [2.057e+01, 1.777e+01, 1.329e+02, …, 1.860e-01, 2.750e-01,
        8.902e-02],
       [1.969e+01, 2.125e+01, 1.300e+02, …, 2.430e-01, 3.613e-01,
        8.758e-02],
       …,
       [1.660e+01, 2.808e+01, 1.083e+02, …, 1.418e-01, 2.218e-01,
        7.820e-02],
       [2.060e+01, 2.933e+01, 1.401e+02, …, 2.650e-01, 4.087e-01,
        1.240e-01],
       [7.760e+00, 2.454e+01, 4.792e+01, …, 0.000e+00, 2.871e-01,
        7.039e-02]]), 'target': array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 1, 1, 1,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
        0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0,
        1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0,
        1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1,
        1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0,
        0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
        1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1,
        1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0,
        0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0,
        1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1,
        1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
        0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1,
        1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0,
        0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0,
        0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0,
```

```
      1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1,
      1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 0,
      1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1,
      1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
      1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
      1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1,
      1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1,
      1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
      1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1]), 'frame': None,
'target_names': array(['malignant', 'benign'], dtype='<U9'), 'DESCR': '..
_breast_cancer_dataset:\n\nBreast cancer wisconsin (diagnostic)
dataset\n---------------------------------------\n\n**Data Set
Characteristics:**\n\n    :Number of Instances: 569\n\n    :Number of
Attributes: 30 numeric, predictive attributes and the class\n\n    :Attribute
Information:\n        - radius (mean of distances from center to points on the
perimeter)\n        - texture (standard deviation of gray-scale values)\n
- perimeter\n        - area\n        - smoothness (local variation in radius
lengths)\n        - compactness (perimeter^2 / area - 1.0)\n        - concavity
(severity of concave portions of the contour)\n        - concave points (number
of concave portions of the contour)\n        - symmetry\n        - fractal
dimension ("coastline approximation" - 1)\n\n        The mean, standard error,
and "worst" or largest (mean of the three\n        worst/largest values) of
these features were computed for each image,\n        resulting in 30 features.
For instance, field 0 is Mean Radius, field\n        10 is Radius SE, field 20
is Worst Radius.\n\n        - class:\n                - WDBC-Malignant\n
- WDBC-Benign\n\n    :Summary Statistics:\n\n
===================================== ====== ======\n
Min    Max\n    ===================================== ====== ======\n    radius
(mean):                           6.981  28.11\n    texture (mean):
9.71   39.28\n    perimeter (mean):                         43.79  188.5\n    area
(mean):                            143.5  2501.0\n    smoothness (mean):
0.053  0.163\n    compactness (mean):                       0.019  0.345\n
concavity (mean):                         0.0    0.427\n    concave points (mean):
0.0    0.201\n    symmetry (mean):                          0.106  0.304\n
fractal dimension (mean):                 0.05   0.097\n    radius (standard error):
0.112  2.873\n    texture (standard error):                0.36   4.885\n
perimeter (standard error):               0.757  21.98\n    area (standard error):
6.802  542.2\n    smoothness (standard error):              0.002  0.031\n
compactness (standard error):             0.002  0.135\n    concavity (standard
error):           0.0    0.396\n    concave points (standard error):      0.0
0.053\n    symmetry (standard error):                0.008  0.079\n    fractal
dimension (standard error):   0.001  0.03\n    radius (worst):
7.93   36.04\n    texture (worst):                          12.02  49.54\n
perimeter (worst):                        50.41  251.2\n    area (worst):
185.2  4254.0\n    smoothness (worst):                       0.071  0.223\n
compactness (worst):                      0.027  1.058\n    concavity (worst):
0.0    1.252\n    concave points (worst):                   0.0    0.291\n
symmetry (worst):                         0.156  0.664\n    fractal dimension
```

(worst):          0.055  0.208\n     ======================================
====== ======\n\n    :Missing Attribute Values: None\n\n     :Class Distribution:
212 - Malignant, 357 - Benign\n\n     :Creator:  Dr. William H. Wolberg, W. Nick
Street, Olvi L. Mangasarian\n\n    :Donor: Nick Street\n\n     :Date: November,
1995\n\nThis is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic)
datasets.\nhttps://goo.gl/U2Uwz2\n\nFeatures are computed from a digitized image
of a fine needle\naspirate (FNA) of a breast mass.  They
describe\ncharacteristics of the cell nuclei present in the image.\n\nSeparating
plane described above was obtained using\nMultisurface Method-Tree (MSM-T) [K.
P. Bennett, "Decision Tree\nConstruction Via Linear Programming." Proceedings of
the 4th\nMidwest Artificial Intelligence and Cognitive Science Society,\npp.
97-101, 1992], a classification method which uses linear\nprogramming to
construct a decision tree.  Relevant features\nwere selected using an exhaustive
search in the space of 1-4\nfeatures and 1-3 separating planes.\n\nThe actual
linear program used to obtain the separating plane\nin the 3-dimensional space
is that described in:\n[K. P. Bennett and O. L. Mangasarian: "Robust
Linear\nProgramming Discrimination of Two Linearly Inseparable
Sets",\nOptimization Methods and Software 1, 1992, 23-34].\n\nThis database is
also available through the UW CS ftp server:\n\nftp ftp.cs.wisc.edu\ncd math-
prog/cpo-dataset/machine-learn/WDBC/\n\n.. topic:: References\n\n   - W.N.
Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction \n     for
breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on \n
Electronic Imaging: Science and Technology, volume 1905, pages 861-870,\n
San Jose, CA, 1993.\n   - O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast
cancer diagnosis and \n     prognosis via linear programming. Operations
Research, 43(4), pages 570-577, \n     July-August 1995.\n   - W.H. Wolberg,
W.N. Street, and O.L. Mangasarian. Machine learning techniques\n     to diagnose
breast cancer from fine-needle aspirates. Cancer Letters 77 (1994) \n
163-171.', 'feature_names': array(['mean radius', 'mean texture', 'mean
perimeter', 'mean area',
       'mean smoothness', 'mean compactness', 'mean concavity',
       'mean concave points', 'mean symmetry', 'mean fractal dimension',
       'radius error', 'texture error', 'perimeter error', 'area error',
       'smoothness error', 'compactness error', 'concavity error',
       'concave points error', 'symmetry error',
       'fractal dimension error', 'worst radius', 'worst texture',
       'worst perimeter', 'worst area', 'worst smoothness',
       'worst compactness', 'worst concavity', 'worst concave points',
       'worst symmetry', 'worst fractal dimension'], dtype='<U23'), 'filename':
'breast_cancer.csv', 'data_module': 'sklearn.datasets.data'}

Loading the dataset to our pandas dataframe

```
[17]: df = pd.DataFrame(dataset.data, columns=dataset.feature_names)
```

```
[18]: df.head()
```

```
[18]:    mean radius  mean texture  mean perimeter  mean area  mean smoothness  \
    0        17.99         10.38          122.80     1001.0          0.11840
    1        20.57         17.77          132.90     1326.0          0.08474
    2        19.69         21.25          130.00     1203.0          0.10960
    3        11.42         20.38           77.58      386.1          0.14250
    4        20.29         14.34          135.10     1297.0          0.10030

        mean compactness  mean concavity  mean concave points  mean symmetry  \
    0            0.27760          0.3001              0.14710         0.2419
    1            0.07864          0.0869              0.07017         0.1812
    2            0.15990          0.1974              0.12790         0.2069
    3            0.28390          0.2414              0.10520         0.2597
    4            0.13280          0.1980              0.10430         0.1809

        mean fractal dimension  …  worst radius  worst texture  worst perimeter  \
    0                  0.07871  …         25.38          17.33           184.60
    1                  0.05667  …         24.99          23.41           158.80
    2                  0.05999  …         23.57          25.53           152.50
    3                  0.09744  …         14.91          26.50            98.87
    4                  0.05883  …         22.54          16.67           152.20

        worst area  worst smoothness  worst compactness  worst concavity  \
    0      2019.0            0.1622             0.6656           0.7119
    1      1956.0            0.1238             0.1866           0.2416
    2      1709.0            0.1444             0.4245           0.4504
    3       567.7            0.2098             0.8663           0.6869
    4      1575.0            0.1374             0.2050           0.4000

        worst concave points  worst symmetry  worst fractal dimension
    0                0.2654          0.4601                  0.11890
    1                0.1860          0.2750                  0.08902
    2                0.2430          0.3613                  0.08758
    3                0.2575          0.6638                  0.17300
    4                0.1625          0.2364                  0.07678

    [5 rows x 30 columns]
```

[19]: `df.shape`

[19]: (569, 30)

x= features which are the dataset values and analyzing the featureswe get the target. y = the target

[27]: 
```
x = df
y =dataset.target
```

[25]: 
```
print(x)
```

|     | mean radius | mean texture | mean perimeter | mean area | mean smoothness \ |
|-----|-------------|--------------|----------------|-----------|-------------------|
| 0   | 17.99       | 10.38        | 122.80         | 1001.0    | 0.11840           |
| 1   | 20.57       | 17.77        | 132.90         | 1326.0    | 0.08474           |
| 2   | 19.69       | 21.25        | 130.00         | 1203.0    | 0.10960           |
| 3   | 11.42       | 20.38        | 77.58          | 386.1     | 0.14250           |
| 4   | 20.29       | 14.34        | 135.10         | 1297.0    | 0.10030           |
| ..  | …           | …            | …              | …         | …                 |
| 564 | 21.56       | 22.39        | 142.00         | 1479.0    | 0.11100           |
| 565 | 20.13       | 28.25        | 131.20         | 1261.0    | 0.09780           |
| 566 | 16.60       | 28.08        | 108.30         | 858.1     | 0.08455           |
| 567 | 20.60       | 29.33        | 140.10         | 1265.0    | 0.11780           |
| 568 | 7.76        | 24.54        | 47.92          | 181.0     | 0.05263           |

|     | mean compactness | mean concavity | mean concave points | mean symmetry \ |
|-----|------------------|----------------|---------------------|-----------------|
| 0   | 0.27760          | 0.30010        | 0.14710             | 0.2419          |
| 1   | 0.07864          | 0.08690        | 0.07017             | 0.1812          |
| 2   | 0.15990          | 0.19740        | 0.12790             | 0.2069          |
| 3   | 0.28390          | 0.24140        | 0.10520             | 0.2597          |
| 4   | 0.13280          | 0.19800        | 0.10430             | 0.1809          |
| ..  | …                | …              | …                   | …               |
| 564 | 0.11590          | 0.24390        | 0.13890             | 0.1726          |
| 565 | 0.10340          | 0.14400        | 0.09791             | 0.1752          |
| 566 | 0.10230          | 0.09251        | 0.05302             | 0.1590          |
| 567 | 0.27700          | 0.35140        | 0.15200             | 0.2397          |
| 568 | 0.04362          | 0.00000        | 0.00000             | 0.1587          |

|     | mean fractal dimension | … | worst radius | worst texture \ |
|-----|------------------------|---|--------------|-----------------|
| 0   | 0.07871                | … | 25.380       | 17.33           |
| 1   | 0.05667                | … | 24.990       | 23.41           |
| 2   | 0.05999                | … | 23.570       | 25.53           |
| 3   | 0.09744                | … | 14.910       | 26.50           |
| 4   | 0.05883                | … | 22.540       | 16.67           |
| ..  | …                      | … | …            | …               |
| 564 | 0.05623                | … | 25.450       | 26.40           |
| 565 | 0.05533                | … | 23.690       | 38.25           |
| 566 | 0.05648                | … | 18.980       | 34.12           |
| 567 | 0.07016                | … | 25.740       | 39.42           |
| 568 | 0.05884                | … | 9.456        | 30.37           |

|     | worst perimeter | worst area | worst smoothness | worst compactness \ |
|-----|-----------------|------------|------------------|---------------------|
| 0   | 184.60          | 2019.0     | 0.16220          | 0.66560             |
| 1   | 158.80          | 1956.0     | 0.12380          | 0.18660             |
| 2   | 152.50          | 1709.0     | 0.14440          | 0.42450             |
| 3   | 98.87           | 567.7      | 0.20980          | 0.86630             |
| 4   | 152.20          | 1575.0     | 0.13740          | 0.20500             |
| ..  | …               | …          | …                | …                   |
| 564 | 166.10          | 2027.0     | 0.14100          | 0.21130             |
| 565 | 155.00          | 1731.0     | 0.11660          | 0.19220             |

|     |        |        |         |         |
|-----|--------|--------|---------|---------|
| 566 | 126.70 | 1124.0 | 0.11390 | 0.30940 |
| 567 | 184.60 | 1821.0 | 0.16500 | 0.86810 |
| 568 |  59.16 |  268.6 | 0.08996 | 0.06444 |

|     | worst concavity | worst concave points | worst symmetry \ |
|-----|-----------------|----------------------|------------------|
| 0   | 0.7119 | 0.2654 | 0.4601 |
| 1   | 0.2416 | 0.1860 | 0.2750 |
| 2   | 0.4504 | 0.2430 | 0.3613 |
| 3   | 0.6869 | 0.2575 | 0.6638 |
| 4   | 0.4000 | 0.1625 | 0.2364 |
| ..  | ...    | ...    | ...    |
| 564 | 0.4107 | 0.2216 | 0.2060 |
| 565 | 0.3215 | 0.1628 | 0.2572 |
| 566 | 0.3403 | 0.1418 | 0.2218 |
| 567 | 0.9387 | 0.2650 | 0.4087 |
| 568 | 0.0000 | 0.0000 | 0.2871 |

|     | worst fractal dimension |
|-----|-------------------------|
| 0   | 0.11890 |
| 1   | 0.08902 |
| 2   | 0.08758 |
| 3   | 0.17300 |
| 4   | 0.07678 |
| ..  | ...     |
| 564 | 0.07115 |
| 565 | 0.06637 |
| 566 | 0.07820 |
| 567 | 0.12400 |
| 568 | 0.07039 |

[569 rows x 30 columns]

```
[28]: print(y)
```

```
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1 0 0 0 0 0 0 0 0 1 0 1 1 1 1 1 0 0 1 0 0 1 1 1 1 0 1 0 0 1 1 1 1 0 1 0 0
 1 0 1 0 0 1 1 1 0 0 1 0 0 0 1 1 1 0 1 1 0 0 1 1 1 0 0 1 1 1 1 0 1 1 0 1 1
 1 1 1 1 1 1 0 0 0 1 0 0 1 1 1 0 0 1 0 1 0 0 1 0 0 1 1 0 1 1 0 1 1 1 1 0 1
 1 1 1 1 1 1 1 1 0 1 1 1 1 0 0 1 0 1 1 0 0 1 1 0 0 1 1 1 1 0 1 1 0 0 0 1 0
 1 0 1 1 1 0 1 1 0 0 1 0 0 0 0 1 0 0 0 1 0 1 0 1 1 0 1 0 0 0 0 1 1 0 0 1 1
 1 0 1 1 1 1 1 0 0 1 1 0 1 1 0 0 1 0 1 1 1 1 0 1 1 1 1 1 0 1 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 1 1 1 1 1 1 0 1 0 1 1 0 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1
 1 0 1 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 0 1 0 1 1 1 1 0 0 0 1 1
 1 1 0 1 0 1 0 1 1 1 0 1 1 1 1 1 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 0 0 1 0 0
 0 1 0 0 1 1 1 1 1 0 1 1 1 1 1 0 1 1 1 0 1 1 0 0 1 1 1 1 1 0 1 1 1 1 1 1
 1 0 1 1 1 1 1 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 0 1 0 0 1 0 1 1 1 1 1 0 1 1
 0 1 0 1 1 0 1 0 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 0 1
 1 1 1 1 1 1 0 1 0 1 1 0 1 0 1 1 1 1 1 0 0 1 0 1 0 1 1 1 1 1 0 1 1 0 1 0 1 0 0]
```

```
1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 0 0 0 0 0 0 1]
```

Splitting the data into traing data and testing data before standardizing. And train test split function was imported already and the function is used below. test_size is how much data we want in our test size, usually 10-20% for test data random state is to reproduce the code = 3 depending on what you want. Is an identity for spliting the data in a specific way.

In most cases it is better to standardize the data before splitting it.

[32]: 
```
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size =0.2,␣
  ↪random_state=3)
```

[33]: 
```
print(x.shape, x_train.shape, x_test.shape)
```

```
(569, 30) (455, 30) (114, 30)
```

Data Standardization; to do this I print the standard deviation of the whole dataset by calling the data 'dataset' because it contains all the data. The data I will call is from 'dataset = sklearn.datasets.load_breast_cancer()'

And find the standard deviation.

If our data contains all all the value in the same range the standard deviation should be 1.

And my dataset below is 228.29 so the dataset are not in the same range and they varies alot

[35]: 
```
print(dataset.data.std())
```

```
228.29740508276657
```

Using StandardScaler since it was imported already

[37]: 
```
scaler = StandardScaler()
```

[38]: 
```
scaler.fit(x_train)
```

[38]: 
```
StandardScaler()
```

[39]: 
```
x_train_standardized = scaler.transform(x_train)
```

[40]: 
```
print(x_train_standardized)
```

```
[[ 1.40381088  1.79283426  1.37960065 …  1.044121    0.52295995
    0.64990763]
 [ 1.16565505 -0.14461158  1.07121375 …  0.5940779   0.44153782
  -0.85281516]
 [-0.0307278  -0.77271123 -0.09822185 … -0.64047556 -0.31161687
  -0.69292805]
 …
 [ 1.06478904  0.20084323  0.89267396 …  0.01694621  3.06583565
  -1.29952679]
 [ 1.51308238  2.3170559   1.67987211 …  1.14728703 -0.16599653
```

```
 0.82816016]
 [-0.73678981 -1.02636686 -0.74380549 … -0.31826862 -0.40713129
 -0.38233653]]
```

For x test

[41]: `x_test_standardized = scaler.transform(x_test)`

Now looking at the standard deviation of our x_train data it shows that our std is 1 And for x test is is 0.87 and it is close to 1 This means that the data are in similar range

[42]: `print(x_train_standardized.std())`

```
1.0
```

[43]: `print(x_test_standardized.std())`

```
0.8654541077212674
```

[ ]: 

[ ]: 

[ ]: 

[ ]: 

[ ]: