

1. Task Description

Task : Data Normalization Techniques and Mean-Centering with Z-Score

(A)(Task for normalization and z-score) Find 10000 row dataset and apply all type of normalization in dataset . 20(B) Apply z-score using sklearn library and do mean-centering of Sales dataset(<https://www.kaggle.com/datasets/nishathakkar/100-sales>)

2. Task Output Screenshot

```
task5.py
34 df_standard.to_csv('zscore_normalized.csv', index=False)
35 df_robust.to_csv('robust_normalized.csv', index=False)
36 df_maxabs.to_csv('maxabs_normalized.csv', index=False)
37
38
39 df_sales = pd.read_csv('100_sales.csv')
40
41
42 numeric_columns = df_sales.select_dtypes(include=[np.number]).columns
43
```

PS C:\Users\krupa\Desktop\task 5> python task5.py

Min-Max Normalization:

	Feature1	Feature2	Feature3	Feature4	Feature5
0	0.978279	0.187570	0.379842	0.259583	0.242114
1	0.121459	0.950565	0.254214	0.117668	0.228934
2	0.592418	0.780821	0.271672	0.663364	0.793394
3	0.462414	0.913356	0.896967	0.445565	0.934287
4	0.349887	0.387641	0.318009	0.298045	0.223293

Z-Score Normalization:

	Feature1	Feature2	Feature3	Feature4	Feature5
0	1.676864	-1.359474	-0.413753	-0.815669	-0.885798
1	-1.382280	1.554737	-0.844982	-1.388497	-0.959118
2	0.334796	0.691379	-0.784987	0.587613	1.022580
3	-0.117105	1.426108	1.361000	-0.169168	1.510311
4	-0.508253	-0.667834	-0.625959	-0.681790	-0.950950

Robust Scaling:

	Feature1	Feature2	Feature3	Feature4	Feature5
0	0.985181	-0.779451	-0.229110	-0.466362	-0.584906
1	-0.743255	0.883857	-0.473428	-0.751452	-0.546944
2	0.286797	0.391089	-0.430476	0.345484	0.589276
3	-0.065457	0.810441	0.776587	-0.092276	0.868919
4	-0.282453	-0.384692	-0.349360	-0.388893	-0.542261

Max-Abs Scaling:

	Feature1	Feature2	Feature3	Feature4	Feature5
--	----------	----------	----------	----------	----------

```
File Edit Selection View Go Run Terminal Help task 5
EXPLORER
TASK 5
100_Sales.csv
maxabs_normalized.csv
minmax_normalized.csv
robust_normalized.csv
sales_data_mean_centered.csv
sales_data_zscore_normalized.csv
task5.py
zscore_normalized.csv

task5.py
34 df_standard.to_csv('zscore_normalized.csv', index=False)
35 df_robust.to_csv('robust_normalized.csv', index=False)
36 df_maxabs.to_csv('maxabs_normalized.csv', index=False)
37
38
39 df_sales = pd.read_csv('100_sales.csv')
40
41
42 numeric_columns = df_sales.select_dtypes(include=[np.number]).columns
43
44
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
4 -0.282453 -0.384692 -0.340360 -0.388893 -0.542261

Max-Abs Scaling:
Feature1 Feature2 Feature3 Feature4 Feature5
0 0.978281 0.107629 0.379888 0.259516 0.242192
1 0.121552 0.950568 0.254270 0.117684 0.221014
2 0.592461 0.708841 0.271727 0.663370 0.793415
3 0.462471 0.913362 0.896974 0.445574 0.934293
4 0.349956 0.307687 0.318861 0.298057 0.223373
C:\Users\krupa\AppData\Local\Programs\Python\Python312\Lib\site-packages\sklearn\utils\extmath.py:1108: RuntimeWarning: invalid value encountered in divide
updated_mean = (last_sum + new_sum) / updated_sample_count
C:\Users\krupa\AppData\Local\Programs\Python\Python312\Lib\site-packages\sklearn\utils\extmath.py:1113: RuntimeWarning: invalid value encountered in divide
T = new_sum / new_sample_count
C:\Users\krupa\AppData\Local\Programs\Python\Python312\Lib\site-packages\sklearn\utils\extmath.py:1133: RuntimeWarning: invalid value encountered in divide
new_unnormalized_variance -= correction**2 / new_sample_count

Z-Score Normalized Sales Data:
Unit_Cost Total_Revenue Total_Profit Unnamed: 0 Unnamed: 10
0 -0.168895 0.798622 1.168192 NaN NaN
1 -0.394831 -0.548427 -0.442948 NaN NaN
2 1.783101 -0.147989 -0.497510 NaN NaN
3 -0.983250 -0.893431 -0.967434 NaN NaN
4 1.783101 1.323690 0.452390 NaN NaN

77°F Mostly cloudy
```

```
File Edit Selection View Go Run Terminal Help task 5
EXPLORER
TASK 5
100_Sales.csv
maxabs_normalized.csv
minmax_normalized.csv
robust_normalized.csv
sales_data_mean_centered.csv
sales_data_zscore_normalized.csv
task5.py
zscore_normalized.csv

task5.py
34 df_standard.to_csv('zscore_normalized.csv', index=False)
35 df_robust.to_csv('robust_normalized.csv', index=False)
36 df_maxabs.to_csv('maxabs_normalized.csv', index=False)
37
38
39 df_sales = pd.read_csv('100_sales.csv')
40
41
42 numeric_columns = df_sales.select_dtypes(include=[np.number]).columns
43
44
PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS
4 0.349956 0.307687 0.318861 0.298057 0.223373
C:\Users\krupa\AppData\Local\Programs\Python\Python312\Lib\site-packages\sklearn\utils\extmath.py:1108: RuntimeWarning: invalid value encountered in divide
updated_mean = (last_sum + new_sum) / updated_sample_count
C:\Users\krupa\AppData\Local\Programs\Python\Python312\Lib\site-packages\sklearn\utils\extmath.py:1113: RuntimeWarning: invalid value encountered in divide
T = new_sum / new_sample_count
C:\Users\krupa\AppData\Local\Programs\Python\Python312\Lib\site-packages\sklearn\utils\extmath.py:1133: RuntimeWarning: invalid value encountered in divide
new_unnormalized_variance -= correction**2 / new_sample_count

Z-Score Normalized Sales Data:
Unit_Cost Total_Revenue Total_Profit Unnamed: 0 Unnamed: 10
0 -0.168895 0.798622 1.168192 NaN NaN
1 -0.394831 -0.548427 -0.442948 NaN NaN
2 1.783101 -0.147989 -0.497510 NaN NaN
3 -0.983250 -0.893431 -0.967434 NaN NaN
4 1.783101 1.323690 0.452390 NaN NaN

Mean-Centered Sales Data:
Unit_Cost Total_Revenue Total_Profit Unnamed: 0 Unnamed: 10
0 -31.628 1.160166e+06 509728.516 NaN NaN
1 -73.938 -7.967049e+05 -193275.624 NaN NaN
2 333.912 -2.149851e+05 -217063.234 NaN NaN
3 -184.120 -1.297896e+06 -422156.164 NaN NaN
4 333.912 1.922937e+06 197395.516 NaN NaN
PS C:\Users\krupa\Desktop\task 5>

77°F Mostly cloudy
```

3) Algorithm Used In Task

1) Min-Max Normalization:

Algorithm: This technique scales the data to a fixed range, typically between 0 and 1. It preserves the relationships between values while ensuring they all fit within a common scale.

2) Z-Score Normalization (Standardization):

Algorithm: Standardization (Z-Score normalization) scales data to have a mean of 0 and a standard deviation of 1, which helps remove the effects of scale from the data.

3) Robust Scaling:

Algorithm: This algorithm scales the data using the **median** and the **interquartile range (IQR)**, making it less sensitive to outliers compared to Min-Max scaling

4) Max-Abs Scaling:

Algorithm: This algorithm scales data based on the **absolute maximum value** in each feature, transforming the data such that the maximum absolute value is 1. This is useful when the data contains both positive and negative values.

5) Mean-Centering:

Algorithm: Mean-centering is the process of subtracting the mean of each feature from the values, shifting the dataset such that the mean of each feature becomes 0. This is not technically normalization but helps in some machine learning algorithms