# Walchand College of Engineering, Sangli.

Experiment No. 2

**Aim:** To perform normalization of data (min-max) and (z-score)

**Theory:** Data normalization is a basic element of data mining. It means transforming the data, namely converting the source data into another format that allows processing data effectively. The main purpose of data normalization is to minimize or even exclude duplicated data. This is a very essential and important issue because it is increasingly problematic to keep data in relational databases, which stores identical data in more than one place.

The use of data mining normalization has a number of advantages.
- Data mining algorithms get more effective and efficient
- Data is converted into the format that everyone can get their heads around
- Data can be extracted from database faster.
- It is possible to analyze the data in specific manner.

**Formula:**

A] Min-Max Normalization

$$D' = \left(\frac{D - minF}{maxF - minF}\right) \times \left(new\_maxF - new\_minF\right) + new\_minF$$

Where $V'$ - new value

$\phantom{xx}V$ - Respected value of the attribute

$\phantom{xx}$maxF - maximum value from given dataset

$\phantom{xx}$minF - minimum value from given dataset

$\phantom{xx}$new-minF : minimum value by user

$\phantom{xx}$new-maxF : maximum value by user

B) Z-Score Normalization.

$$V' = \frac{V - \bar{F}}{\sigma F}$$

Where $V$ - actual data value

$\phantom{xxx}F$ - mean value of data

$\phantom{xxx}\sigma F$ - standard deviation of data

Algorithm:

1) Take dataset and read data from csv file.

2) If we choose min-max normalization then find min-max value from given dataset and if it is z-score normalization then. Calculate mean and standard deviation.

3) For min-max normalization consider new-min, new-max by considering suitable range.

4) Calculate normalization value for all data in the given dataset with using respective normalized formula.

Example : A] Min-Max normalization

Data [marks]     min = 20     max = 75
    20           new-min = 0   new_man = 1
    75
    84
    55
    63

For marks 20:

$$V'_{20} = \frac{20-20}{75-20} \times (1-0) + 0 = 0$$

$$V'_{75} = \frac{75-20}{75-20} \times (1-0) + 0 = 1.0$$

$$V'_{84} = \frac{34-20}{75-20} \times (1-0) + 0 = 0.2545$$

$$V'_{55} = \frac{55-20}{75-20} \times (1-0) + 0 = 0.6363$$

$$V'_{63} = \frac{63-20}{75-20} \times (1-0) + 0 = 0.7818$$

Normalized Table

| Marks | Normalized Marks |
|-------|------------------|
| 20    | 0                |
| 75    | 1                |
| 34    | 0.25            |
| 55    | 0.64            |
| 63    | 0.78            |

B) Z-score Normalization

Marks : $[9, 12, 15, 20]$

Mean $= \dfrac{9+12+15+20}{4} = 14$ $\qquad p = 14$

Standard Deviation.

$$\sigma = \sqrt{\dfrac{(9-14)^2 + (12-14)^2 + (15-14)^2 + (20-14)^2}{4}}$$

$$= \sqrt{\dfrac{25+4+1+36}{4}} = \sqrt{16.5} = 4.062$$

For marks 9

$$Z_9 = \frac{9-14}{4.062} = -1.251$$

For mark 12:

$$Z_{12} = \frac{12-14}{4.062} = -0.4923$$

$$Z_{15} = \frac{15-14}{4.062} = 0.2462$$

$$Z_{20} = \frac{20-14}{4.062} = 1.4771$$

| Marks | Z-score Normalized Marks |
|-------|--------------------------|
| 9 | -1.23 |
| 12 | -0.49 |
| 25 | 0.25 |
| 20 | 1.48 |

Conclusions:
The transformed values obtained by calculation and by knime and weka are identical successfully performed transformation of given data using min-max and x-score normalization

Knime Result :

Min-Max

| Input Marks | Knime min. man Normalized |
|---|---|
| 20 | 0 |
| 25 | |
| 34 | 0.255 |
| 55 | 0.636 |
| 63 | 0.783 |

Z - score

| I/P Mark | Z-score Normalized |
|---|---|
| 9 | -1.066 |
| 12 | -0.426 |
| 15 | 0.213 |
| 20 | 1.279 |

Conclusion :

Normalization is an important processor task in data preprocessing. It is used to ensure consistency d in data records. In order to bring all attribute on the same scale min-man normalization is used. Standardiaing scor on some scale by dividing a scores deviation by standard deviation for that z-score is used It scales a dote to particular small scale which helps to allow processing data efficiency.

Program:

```cpp
#include <bits/stdc++.h>
#include <fstream>
using namespace std;

void calculateMinMax(ifstream &inputFile, double &minValue, double &maxValue)
{
    double currentValue;
    inputFile >> currentValue;
    while (inputFile)
    {
        if (currentValue > maxValue)
            maxValue = currentValue;
        if (currentValue < minValue)
            minValue = currentValue;
        inputFile >> currentValue;
    }
}

void performMinMaxNormalization(ifstream &inputFile, ofstream &outputFile, double oldMin,
double oldMax, double newMin, double newMax)
{
    outputFile << "Original Data,"
            << "Normalized Data"
            << "\n";

    double currentValue;
    inputFile >> currentValue;
    while (inputFile)
    {
        double previousValue = currentValue;
        currentValue = (((currentValue - oldMin) / (oldMax - oldMin)) * (newMax - newMin)) +
newMin;
        outputFile << previousValue << "," << currentValue << "\n";
        inputFile >> currentValue;
    }
}

void performZScoreNormalization(ifstream &inputFile, ofstream &outputFile)
```

```cpp
{
    double sum = 0.0, count = 0.0, squareSum = 0.0, mean, standardDeviation;
    double currentValue;

    // Calculate mean
    while (inputFile)
    {
        sum += currentValue;
        count++;
        inputFile >> currentValue;
    }
    mean = sum / count;

    // Calculate standard deviation
    inputFile.clear();
    inputFile.seekg(0, ios::beg);
    while (inputFile)
    {
        squareSum += (currentValue - mean) * (currentValue - mean);
        inputFile >> currentValue;
    }
    inputFile.clear();
    inputFile.seekg(0, ios::beg);

    standardDeviation = sqrt(squareSum / count);

    // Perform z-score normalization
    outputFile << "Original Data,"
            << "Normalized Data"
            << "\n";

    while (inputFile)
    {
        double prev = currentValue;
        currentValue = (currentValue - mean) / standardDeviation;
        outputFile << prev << "," << currentValue << endl;
        inputFile >> currentValue;
    }
}
```

```cpp
int main()
{
    double currentValue, minValue, maxValue, newMinValue, newMaxValue;
    double sum, count, squareSum, mean, standardDeviation;

    ifstream inputFileMinMax("exp2_input_MinMax.csv");
    ifstream inputFileMinMax_2("exp2_input_MinMax.csv");
    ifstream inputFileZScore("exp2_input_Zscore.csv");

    int option;
    cout << "\nEnter an option: \n1. Min-Max Normalization \n2. Z-Score Normalization\nOption: ";
    cin >> option;

    ofstream outputFileMinMax("exp2_output_MinMax.csv", ios::app);
    ofstream outputFileZScore("exp2_output_ZScore.csv", ios::app);

    switch (option)
    {
    case 1: // Min-Max Normalization
        if (!inputFileMinMax)
        {
            cout << "Error opening file, please try again.";
            exit(0);
        }
        calculateMinMax(inputFileMinMax, minValue, maxValue);
        cout << "Enter new minimum value: ";
        cin >> newMinValue;
        cout << "\nEnter new maximum value: ";
        cin >> newMaxValue;
        performMinMaxNormalization(inputFileMinMax_2, outputFileMinMax, minValue,
maxValue, newMinValue, newMaxValue);
        break;

    case 2: // Z-Score Normalization
        if (!inputFileZScore)
        {
            cout << "Error opening file, please try again.";
            exit(0);
        }
```
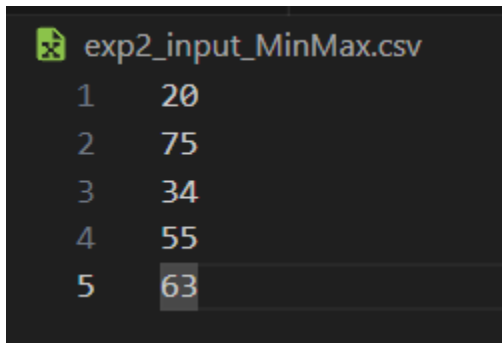
```
        performZScoreNormalization(inputFileZScore, outputFileZScore);
        break;

    default:
        cout << "Invalid option";
    }

    return 0;
}
```
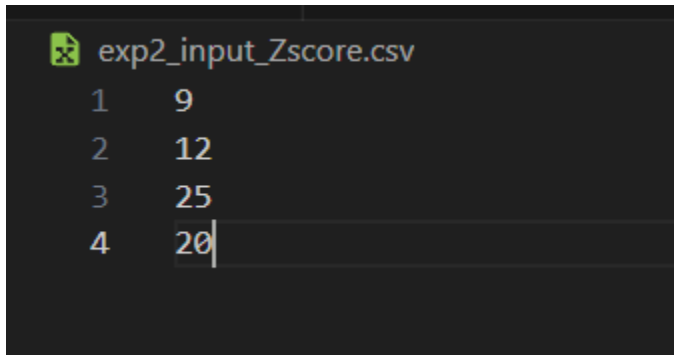
I/P

Min-Max



Z-Score



,

O/P

MIN_MAX

```
exp2_output_MinMax.csv
1    Original Data,Normalized Data
2    20,0.266667
3    75,1
4    34,0.453333
5    55,0.733333
6    63,0.84
7    |
```

Z-Score

```
exp2_output_Zscore.csv
1    Original Data,Normalized Data
2    20,0.96013
3    9,-0.593022
4    12,-0.169435
5    25,1.66611
6    20,0.96013
7
```