# Analysis and Systems of

# Big Data Practise

# Lab - 3

Sreepathy Jayanand
CED17IO38

# Q1

1. Suppose that the data for analysis includes the attribute *age*. The *age* values for the data tuples are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.

(a) Use *min-max normalization* to transform the values of age to the range [0:1].
(b) Use *z-score normalization* to transform the values of *age*.
(c) Use normalization by *decimal scaling* to transform the values of *age* such that the transformed value is less than 1.

2. Use the given dataset and perform the operations listed below.

**Logic :** Using the formulas mentioned below, the new set of values is found out.

**Min-Max Value:**

age[i] = new_min + (age[i] - old_min) * (new_max - new_min) / (old_max - old_min)

**Z-Score Value:**

age[i] = (age[i] - Mean)/ Standard_Deviation

**Decimal Scaling:**

age[i] = age[i] / 100

## Libraries used:

- Statistics - For mean and standard deviation

## Code:

```python
import statistics

age = [13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25,
25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70]
output_1 = []
output_2 = []
output_3 = []
mean = statistics.mean(age)
stdev = statistics.stdev(age)
for i in range(len(age)):
    temp1 = ((age[i] - min(age))) / (max(age) - min(age))
    output_1.append(round(temp1, 2))
    temp2 = (age[i] - mean) / stdev
    output_2.append(round(temp2, 2))
    temp3 = age[i] / 100
    output_3.append(temp3)

print("MIN-MAX representation : ")
print(output_1)
print("Z-SCORE representation : ")
print(output_2)
print("DECIMAL SCALED representation : ")
print(output_3)
```

**Output:**

```
[SJ-$ python3 1.py
MIN-MAX representation :
[0.0, 0.04, 0.05, 0.05, 0.11, 0.12, 0.12, 0.14, 0.16, 0.16, 0.21, 0.21, 0.21, 0.
21, 0.3, 0.35, 0.35, 0.39, 0.39, 0.39, 0.39, 0.4, 0.47, 0.56, 0.58, 0.68, 1.0]
Z-SCORE representation :
[-1.31, -1.16, -1.08, -1.08, -0.85, -0.77, -0.77, -0.69, -0.62, -0.62, -0.38, -0
.38, -0.38, -0.38, 0.0, 0.23, 0.23, 0.39, 0.39, 0.39, 0.39, 0.47, 0.78, 1.16, 1.
24, 1.7, 3.09]
DECIMAL SCALED representation :
[0.13, 0.15, 0.16, 0.16, 0.19, 0.2, 0.2, 0.21, 0.22, 0.22, 0.25, 0.25, 0.25, 0.2
5, 0.3, 0.33, 0.33, 0.35, 0.35, 0.35, 0.35, 0.36, 0.4, 0.45, 0.46, 0.52, 0.7]
SJ-$ ▊
```

# Q2

1. Suppose that the data for analysis includes the attribute *age*. The *age* values for the data tuples are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.

From the Avocado Dataset,

Sort the attribute "Total Volume" in the given dataset and distribute the data into equal sized/frequency bins. Let the number of bins be 250. **Smooth** the sorted data by

(i) bin-means

(ii) bin-medians

(iii) bin-boundaries


b. The dataset represents weekly retail scan data for National retail volume (units) and price. However, the company is interested in knowing the *monthly (total per month) and annual sales (total per year)*, rather than the total per week. So, **reduce** the data accordingly.

c. Summarise the number of missing values for each attribute

d. Populate data for the **missing values** of the attribute= "Average Price" by averaging all the values of
the "Avg Price" attribute that fall under the same **"REGION"** attribute value.

e. **Discretise** the attribute= "Date" using *concept hierarchy* into {Old, New, Recent}


**Logic :** a) Different bins with size 250 is created. The data is smoothened by bin-means, bin medians, bin boundaries

b) The data is reduced by considering the total per month and year instead of week.

c) The columns with datatype float containing null value is found by converting the value to string datatype and removing the '.' , if the resultant string is not a digit then it is considered as null value. For columns with string datatype we are directly checking whether it is a missing value or no.

d) The missing values are replaced by average in that particular region.

e) The date in a data frame is replaced by old if the corresponding year is either 2015 or 2016, the date is replaced by new if the corresponding year is 2017 and the date is replaced by recent if the corresponding year is 2018.

## Libraries used:
- Csv - To import csv files
- Math - For ceil function to compute number of bins
- Numpy - For array manipulation

## Code:

```
import csv

import math
import numpy as np


with open("Lab_Data.csv", 'r') as f:
    data = list(csv.reader(f, delimiter = ","))

total_volume = []

for i in range(1, len(data)):
        total_volume.append(float(data[i][2]))

total_volume = np.array(total_volume)
total_volume = np.sort(total_volume)

no_of_bins = 250
bin_size = math.ceil(len(total_volume) / no_of_bins)

bin1 = np.zeros((no_of_bins, bin_size))
bin2 = np.zeros((no_of_bins, bin_size))
```

```python
bin3 = np.zeros((no_of_bins, bin_size))

for i in range (0, no_of_bins*bin_size, bin_size):
        k = int(i/bin_size)
        mean = sum(total_volume[i:i+bin_size])/bin_size
        for j in range(bin_size):
                bin1[k,j] = round(mean, 6)
print("Bin Mean: \n",bin1)

for i in range (0, no_of_bins*bin_size, bin_size):
        k = int(i/bin_size)
        for j in range (bin_size):
                bin2[k,j] = total_volume[i+math.floor(bin_size/2)]
print("Bin Median: \n",bin2)

for i in range (0, no_of_bins*bin_size, bin_size):
        k = int(i/bin_size)
        for j in range (bin_size):
                if (total_volume[i+j]-total_volume[i]) < (total_volume[i+bin_size-1]-
total_volume[i+j]):
                        bin3[k,j] = total_volume[i]
                else:
                        bin3[k,j] = total_volume[i+bin_size-1]
print("Bin Boundaries: \n",bin3)




date = data[1][0][-7:]
region = data[1][12]
c = 0



list_date = []
list_avg_price = []
list_total_volume = []
list_4046 = []
```

```python
list_4225 = []
list_4770 = []
list_total_bags = []
list_small_bags = []
list_large_bags = []
list_xlarge_bags = []
list_region = []

sum_avg_price = 0
sum_total_volume = 0
sum_4046 = 0
sum_4225 = 0
sum_4770 = 0
sum_total_bags = 0
sum_small_bags = 0
sum_large_bags = 0
sum_xlarge_bags = 0

for i in range(1, len(data)):
        if(date == data[i][0][-7:] and region == data[i][12]):
                d = data[i][1]

                if(d.replace('.', '', 1).isdigit() == True):
                        sum_avg_price += float(data[i][1])
                        c += 1

                sum_total_volume += float(data[i][2])
                sum_4046 += float(data[i][3])
                sum_4225 += float(data[i][4])
                sum_4770 += float(data[i][5])
                sum_total_bags += float(data[i][6])
                sum_small_bags += float(data[i][7])
                sum_large_bags += float(data[i][8])
                sum_xlarge_bags += float(data[i][9])
        else:
```

```python
                list_date.append(date)

                if(c != 0):

                        list_avg_price.append(sum_avg_price/c)
                else:
                        list_avg_price.append(0)
                list_total_volume.append(sum_total_volume)
                list_4046.append(sum_4046)
                list_4225.append(sum_4225)
                list_4770.append(sum_4770)
                list_total_bags.append(sum_total_bags)
                list_small_bags.append(sum_small_bags)
                list_large_bags.append(sum_large_bags)
                list_xlarge_bags.append(sum_xlarge_bags)
                list_region.append(region)
                date = data[i][0][-7:]
                region = data[i][12]
                d = data[i][1]
                if(d.replace('.', '', 1).isdigit() == True):
                        sum_avg_price = float(data[i][1])
                        c = 1
                else:
                        sum_avg_price = 0
                        c = 0
                sum_total_volume = float(data[i][2])
                sum_4046 = float(data[i][3])
                sum_4225 = float(data[i][4])
                sum_4770 = float(data[i][5])
                sum_total_bags = float(data[i][6])
                sum_small_bags = float(data[i][7])
                sum_large_bags = float(data[i][8])
                sum_xlarge_bags = float(data[i][9])

print("Date   avg_price   total_volume      4046  4225  4770   total_bags
small_bags   large_bags    xlarge_bags   list_region")
```

```python
for i in range(20):
        print(list_date[i], "   ", round(list_avg_price[i], 2), "       ",
round(list_total_volume[i], 2), "    ", round(list_4046[i], 2), "    ",
                round(list_4225[i], 2), "        ", round(list_4770[i], 2), "    ",
round(list_total_bags[i], 2), "        ", round(list_small_bags[i], 2), "    ",
                round(list_large_bags[i], 2), "        ", round(list_xlarge_bags[i], 2),
"        ", list_region[i])




temp_lists0 = 0
temp_lists1 = 0
temp_lists2 = 0
temp_lists3 = 0
temp_lists4 = 0
temp_lists5 = 0
temp_lists6 = 0
temp_lists7 = 0
temp_lists8 = 0
temp_lists9 = 0
temp_lists10 = 0
temp_lists11 = 0
temp_lists12 = 0
for i in range(1, len(data)):
        d0 = data[i][0]

        d1 = data[i][1]
        d2 = data[i][2]
        d3 = data[i][3]
        d4 = data[i][4]
        d5 = data[i][5]
        d6 = data[i][6]
        d7 = data[i][7]
        d8 = data[i][8]
        d9 = data[i][9]
        if(d0 == ''):
```

```python
            print(d0)
            temp_lists0+=1
    elif(d1 == '' or d1.replace('.', '', 1).isdigit() == False):
            temp_lists1+=1
    elif(d2 == '' or d2.replace('.', '', 1).isdigit() == False):
            temp_lists2+=1
    elif(d3 == '' or d3.replace('.', '', 1).isdigit() == False):
            temp_lists3+=1
    elif(d4 == '' or d4.replace('.', '', 1).isdigit() == False):
            temp_lists4+=1
    elif(d5 == '' or d5.replace('.', '', 1).isdigit() == False):
            temp_lists5+=1
    elif(d6 == '' or d6.replace('.', '', 1).isdigit() == False):
            temp_lists6+=1
    elif(d7 == '' or d7.replace('.', '', 1).isdigit() == False):
            temp_lists7+=1
    elif(d8 == '' or d8.replace('.', '', 1).isdigit() == False):
            temp_lists8+=1
    elif(d9 == '' or d9.replace('.', '', 1).isdigit() == False):
            temp_lists9+=1
    elif(data[i][10] == ''):
            temp_lists10+=1
    elif(data[i][11].isdigit() == False):
            temp_lists11+=1
    elif(data[i][12] == ''):
            temp_lists12+=1


print("The count of the missing values are-")
print("Date - ", temp_lists0)
print("Averge price - ", temp_lists1)
print("Total volume - ", temp_lists2)
print("4046 - ", temp_lists3)
print("4225 - ", temp_lists4)
print("4770 - ", temp_lists5)
print("Total bags - ", temp_lists6)
print("Small bags - ", temp_lists7)
```

```python
print("Large bags - ", temp_lists8)
print("XLarge bags - ", temp_lists9)
print("Type - ", temp_lists10)
print("Year - ", temp_lists11)
print("Region - ", temp_lists12)


can = 0
for i in range(1, len(data)):
        d = data[i][1]
        if(data[i][1] == '' or d.replace('.', '', 1).isdigit() == False):
                s = 0
                count = 0
                for j in range(1, len(data)):
                        d = data[j][1]
                        if(data[j][12] == data[i][12] and d.replace('.', '', 1).isdigit() ==
True):
                                s = s + float(data[j][1])
                                count+=1

                data[i][1] = str(round(s/count, 6))
print("Earlier missing entries : ")
print(data[5])
print(data[6])
print(data[7])
print(data[8])
print(data[9])
print(data[10])
print(data[11])




for i in range(1, len(data)):
        if((data[i][11] == '2015') or (data[i][11] == '2016')):
                data[i][0] = "Old"
```

```
        elif(data[i][11] == '2017'):
                data[i][0] = "New"
        elif(data[i][11] == '2018'):
                data[i][0] = "Recent"
print(data[10])
print(data[3000])
print(data[6000])
print(data[9000])
print(data[12000])

print(data[15000])
print(data[18000])
```

## Output:

(a)Bin mean, Bin median, Bin boundaries (Total 250 bins)

```
Bin Mean:
 [[7.35713014e+02 7.35713014e+02 7.35713014e+02 ... 7.35713014e+02
  7.35713014e+02 7.35713014e+02]
 [1.03670205e+03 1.03670205e+03 1.03670205e+03 ... 1.03670205e+03
  1.03670205e+03 1.03670205e+03]
 [1.17840082e+03 1.17840082e+03 1.17840082e+03 ... 1.17840082e+03
  1.17840082e+03 1.17840082e+03]
 ...
 [1.35563292e+07 1.35563292e+07 1.35563292e+07 ... 1.35563292e+07
  1.35563292e+07 1.35563292e+07]
 [3.11970276e+07 3.11970276e+07 3.11970276e+07 ... 3.11970276e+07
  3.11970276e+07 3.11970276e+07]
 [3.89735224e+07 3.89735224e+07 3.89735224e+07 ... 3.89735224e+07
  3.89735224e+07 3.89735224e+07]]
Bin Median:
 [[7.74200000e+02 7.74200000e+02 7.74200000e+02 ... 7.74200000e+02
  7.74200000e+02 7.74200000e+02]
 [1.03500000e+03 1.03500000e+03 1.03500000e+03 ... 1.03500000e+03
  1.03500000e+03 1.03500000e+03]
 [1.17595000e+03 1.17595000e+03 1.17595000e+03 ... 1.17595000e+03
  1.17595000e+03 1.17595000e+03]
 ...
 [8.38991804e+06 8.38991804e+06 8.38991804e+06 ... 8.38991804e+06
  8.38991804e+06 8.38991804e+06]
 [3.13460915e+07 3.13460915e+07 3.13460915e+07 ... 3.13460915e+07
  3.13460915e+07 3.13460915e+07]
 [3.73523606e+07 3.73523606e+07 3.73523606e+07 ... 3.73523606e+07
  3.73523606e+07 3.73523606e+07]]
Bin Boundaries:
 [[8.45600000e+01 8.45600000e+01 8.45600000e+01 ... 9.34950000e+02
  9.34950000e+02 9.34950000e+02]
 [9.36690000e+02 9.36690000e+02 9.36690000e+02 ... 1.11744000e+03
  1.11744000e+03 1.11744000e+03]
 [1.11847000e+03 1.11847000e+03 1.11847000e+03 ... 1.23327000e+03
  1.23327000e+03 1.23327000e+03]
 ...
 [7.36092584e+06 7.36092584e+06 7.36092584e+06 ... 2.80125209e+07
  2.80125209e+07 2.80125209e+07]
 [2.80413354e+07 2.80413354e+07 2.80413354e+07 ... 3.39939313e+07
  3.39939313e+07 3.39939313e+07]
 [3.41267310e+07 3.41267310e+07 3.41267310e+07 ... 6.25056465e+07
  6.25056465e+07 6.25056465e+07]]
```

(b) Reduced format

(c) Average price missing 48 entries

(d) Replaced avg price with mean 1.570755

(e)According to year, they have been categorised as Old, New, recent. Different examples have been taken from 2015, 2016, 2017, 2018 to illustrate.