

CSC 522 - Fall 2018

HOMEWORK 1

Homework Group 43

Name	UnityID
Aman Chauhan	achauha3
Khantil Choksi	khchoksi

Problem 1: Classifying attributes:

	Attribute Name	binary / discrete / continuous	nominal / ordinal / interval or ratio
(a)	Hair color (Black, Blonde, Red)	discrete	nominal
(b)	Level of agreement (yes, maybe, no)	discrete	ordinal (Yes - can be considered as maximum level of agreement; maybe - can be considered as neutral level of agreement and no - can be considered as minimum level of agreement)
(c)	Income earned in a week	continuous	ratio
(d)	Celsius temperature	continuous	interval
(e)	Genotype (Bb, bb, BB, bB)	discrete	nominal
(f)	ISBN numbers for books	discrete	nominal
(g)	Time in terms of AM or PM	discrete (like 6AM, 7AM, 2PM, hence there are only finite values in a day)	ordinal (You can order these times as 8AM comes before 11PM for the same day)
(h)	Waiting number for restaurant	discrete	ordinal
(i)	Years of work experience	continuous (1.5, 2.5)	ratio
(j)	Categorization of clothing (hat, shirt, pants, shoes)	discrete	nominal
(k)	Angles as measured in degrees between 0 and 360	continuous	ratio
(l)	Ratings of movies (G, PG, R)	discrete	ordinal (based on level of violence and profanity in language)

(m)	Coat check number	discrete	nominal (follows only property of distinctness e.g. if the number is "equal to" or "not equal to")
-----	-------------------	----------	---

Problem 2: Term Frequency - Inverse Document Frequency

(A) Min and Max values for the functions

For equation (1), tf_{ij}' :

- Maximum value = $p * \log m$
If out of 'm' documents, word occurred only in 1 document with p frequency then $df_i = 1$.
- Minimum value = 0
If a word occurs at least 1 time in each and every document then $df_i = m$. Therefore, $\log \frac{m}{df_i} = 0$.

For equation (2), tf_{ij}''

- Maximum value = $p * \log \frac{\sum_{k=1}^m d_k}{p}$
If i^{th} word occurs only in j^{th} document and it is the only word occurring p times in j^{th} document then we will get maximum inverse document frequency.
- Minimum value = 0
If all the documents contains the i^{th} word. Therefore,

$$\sum_{k=1}^m d_k = \sum_{k=1}^{df_i} d_k \Rightarrow p * \log 1 \Rightarrow 0$$

(B) Importance of TF-IDF

Both tf_{ij}' and tf_{ij}'' are used in natural language processing for classification of documents, information retrieval and as a tool to compare similarity between documents.

TF_{ij}' is used in scenarios where the length of all documents in the corpus are comparable. TF_{ij}'' is used in scenarios where the length of documents in the corpus vary a lot.

TF_{ij}' does not take into consideration the length of document in which that word occurs in the IDF portion. Thus it gives equal weightage to words which are unique in both small and long documents. TF_{ij}'' considers the length of the document in which that word is occurring. Thus, the IDF portion of this equation gives more weightage to terms that are unique in shorter documents.

Problem 3:

(A) Answer: Stratified Sampling

Reason: The sampling method depends on the overall population of the given data. Here the patients having albinism and normal skin is proportional of 5% to 95%.

- Random sampling involves the random selection of datapoint from the entire population; so each possible sample is equally likely to occur. If we choose to get a random sample of size 100 over the entire population then there is chance that the random sample turns out to be not well balanced across the patients having albinotic vs. normal skin and hence it is biased causing a significant error in estimation.
- In contrast, stratified sampling divides members of the population into homogeneous subgroups before sampling. A random sample is taken from each subgroup in direct proportion to the size of the stratum compared to the population. The sample subsets are then combined to create a final random sample.
- Here the population density varies greatly within a data region, stratified sampling will ensure that estimates can be made with equal accuracy in different parts of the region.
- So, after stratified sampling,
 number of albinotic patients sampled = $\frac{26150}{523000} * 1000 = 50$
 number of normal skin patients sampled = $\frac{496850}{523000} * 1000 = 950$

(B) The recorded systolic blood pressure (SBP) of 250 can be considered an **outlier** because this is a *true recorded value* which is significantly away from the critical value.

Reason 1 - If we consider a distribution of SBP, then the question suggests that an SBP of 180 would be critical, and thus can be considered as an upper extreme value of the distribution. The value of 250 is even beyond that upper extreme. Hence this is either noise or an error.

Reason 2 - This value was recorded by a person using an instrument, which means that the observation was supervised, ruling out the possibility of error generated by the machine.

Q4

Import the libraries

```
1 import random
2 import math
3 import sys
4 import os
```

Helper Functions

Function to print the matrix

```
1 def pprint(mat):
2     pad = '{:>}' + str(len(str(max([max(mat[i]) for i in range(len(mat))]))))
3     + '}'
4     print('[' + '\n '.join(['[' + ', '.join([pad.format(mat[i][j]) for j in
5     range(len(mat[0]))]) + ']' for i in range(len(mat))] + '']')
```

Function to create identity matrix of size 'dim'

```
1 def identity_matrix(dim):
2     return [[1 if j==i else 0 for j in range(dim)] for i in range(dim)]
```

Function to set the column to a particular value

```
1 def column_manipulate(mat, col, val):
2     for row in mat:
3         row[col] = val
4     return mat
```

Function to sum all values in the matrix

```

1 def matrix_sum(mat):
2     cnt = 0
3     for row in mat:
4         for col in row:
5             cnt += col
6     return cnt

```

Transpose any given MxN matrix

```

1 def transpose(mat):
2     result = [[None for j in range(len(mat))] for i in range(len(mat[0]))]
3     for i in range(len(mat)):
4         for j in range(len(mat[i])):
5             result[j][i] = mat[i][j]
6     return result

```

Calculate sum for a given row and sum of diagonal elements

```

1 def row_diagonal_sum(mat, row):
2     rowcnt = 0
3     dcnt = 0
4     for i in range(len(mat)):
5         for j in range(len(mat[i])):
6             if i==j:
7                 dcnt += mat[i][j]
8             if i==row:
9                 rowcnt += mat[i][j]
10    return rowcnt, dcnt

```

Generate NxN Gaussian Matrix for a given mean and variance

```

1 def gaussian_matrix(dim, mean, variance, dtype):
2     if dtype=="int":
3         return [[math.floor(random.gauss(mean, math.sqrt(variance))) for j
4 in range(dim)] for i in range(dim)]
5     else:
6         return [[random.gauss(mean, math.sqrt(variance)) for j in
7 range(dim)] for i in range(dim)]

```

Function to multiply 2 matrices

```
1 def matrix_multiply(a, b):
2     result = [[0 for j in range(len(b[0]))] for i in range(len(a))]
3     for i in range(len(a)):
4         for j in range(len(b[0])):
5             for k in range(len(b)):
6                 result[i][j] += a[i][k] * b[k][j]
7     return result
```

Function to multiply individual elements of given 2 matrices

```
1 def multiply(a, b):
2     result = [[0 for j in range(len(a[0]))] for i in range(len(a))]
3     for i in range(len(a)):
4         for j in range(len(a[0])):
5             result[i][j] = a[i][j] * b[i][j]
6     return result
```

Function to add 2 matrices

```
1 def add(a, b):
2     result = [[0 for j in range(len(a[0]))] for i in range(len(a))]
3     for i in range(len(a)):
4         for j in range(len(a[0])):
5             result[i][j] = a[i][j] + b[i][j]
6     return result
```

Function to shift all rows up by 1

```
1 def row_shift(mat):
2     return mat[1:] + [mat[0]]
```

Function to find covariance between 2 vectors

```
1 def covariance(x, y):
2     mean_x = sum(x)/len(x)
3     mean_y = sum(y)/len(y)
4     sumo = 0.0
5     for i in range(len(x)):
6         sumo += (x[i] - mean_x) * (y[i] - mean_y)
7     sumo /= (len(x) - 1)
8     return sumo
```

(A) Generate 5x5 identity matrix

```
1 A = identity_matrix(5)
2 pprint(A)
```

```
1 [[1, 0, 0, 0, 0]
2  [0, 1, 0, 0, 0]
3  [0, 0, 1, 0, 0]
4  [0, 0, 0, 1, 0]
5  [0, 0, 0, 0, 1]]
```

(B) Change all elements in the 2nd column of A to 3.

```
1 A = column_manipulate(A, 1, 3)
2 pprint(A)
```

```
1 [[1, 3, 0, 0, 0]
2  [0, 3, 0, 0, 0]
3  [0, 3, 1, 0, 0]
4  [0, 3, 0, 1, 0]
5  [0, 3, 0, 0, 1]]
```

(C) Sum of all elements in the matrix

```
1 print(matrix_sum(A))
```

```
1 19
```

(D) Transpose the matrix A

```
1 A = transpose(A)
2 pprint(A)
```

```
1 [[1, 3, 0, 0, 0]
2  [0, 3, 0, 0, 0]
3  [0, 3, 1, 0, 0]
4  [0, 3, 0, 1, 0]
5  [0, 3, 0, 0, 1]]
```

(E) Calculate sum of the 3rd row, and the diagonal in the matrix A.

```
1 rowcnt, dcnt = row_diagonal_sum(A, 2)
2 print('Sum of 3rd row - ' + str(rowcnt))
3 print('Sum of major diagonal - ' + str(dcnt))
```

```
1 Sum of 3rd row - 4
2 Sum of major diagonal - 7
```

(F) Generate a 5*5 matrix B following Gaussian Distribution with mean 5 and variance 3.

```
1 B = gaussian_matrix(5, 5, 3, "float")
2 pprint(B)
```



```

1 [[6.215699316621983, 5.3314324181394275, 3.160592195218215,
   5.731340893558323, 3.8761641182554536]
2  [2.726979912172165, 3.363145178336113, 5.168241217648324,
   3.3272058142040364, 4.53812449216301]
3  [5.8952555626790035, 5.008129339178955, 3.3278857168552056,
   4.2163819866696235, 0.2846806253144418]
4  [5.760377878211342, 5.911725654130385, 3.335467243819485,
   5.183137167675729, 4.836748684684982]
5  [4.616059063383704, 6.42657606799423, 4.242358931555765, 4.656873001603596,
   2.9403643600492413]]

```

(G)

```

1 C1 = multiply(matrix_multiply([[1,0,0,0,0]]*5, row_shift(B)),
   identity_matrix(5))
2 C2 = matrix_multiply(matrix_multiply([[1,0,0,0,0], [0,0,0,0,0]], B), C1)
3 C3 = matrix_multiply([[0,0,0,0,0], [0,0,1,1,-1]], B)
4 C = add(C2, C3)
5 pprint(C)

```

```

1 [[ 16.9500871765304, 17.930381230690458, 16.33470285550438,
   19.069350744232608, 17.59051532069851]
2  [ 7.03957437750664, 4.49327892531511, 2.4209940291189254,
   4.742646152741756, 2.181064949950182]]

```

(H)

```

1 D1 = multiply([[2,3,4,5,6]]*5, identity_matrix(5))
2 D = matrix_multiply(C, D1)
3 pprint(D)

```

```

1 [[ 33.9001743530608, 53.79114369207137, 65.33881142201751,
   95.34675372116304, 105.54309192419106]
2  [ 14.07914875501328, 13.47983677594533, 9.683976116475701,
   23.71323076370878, 13.086389699701094]]

```

(I) Covariance Matrix

```

1 print('X\tY\tZ')
2 print('-\t-\t-')
3 print('{}\t{}\t{}'.format(2, 6, 1))
4 print('{}\t{}\t{}'.format(4, 5, 3))
5 print('{}\t{}\t{}'.format(6, 4, 5))
6 print('{}\t{}\t{}'.format(8, 3, 7))
7 X = [2, 4, 6, 8]
8 Y = [6, 5, 4, 3]
9 Z = [1, 3, 5, 7]
10 print('\nCovariance Matrix')
11 print('\tX\tY\tZ')
12 print('X\t{:0.2f}\t{:0.2f}\t{:0.2f}'.format(covariance(X, X), covariance(X,
Y), covariance(X, Z)))
13 print('Y\t{:0.2f}\t{:0.2f}\t{:0.2f}'.format(covariance(Y, X), covariance(Y,
Y), covariance(Y, Z)))
14 print('Z\t{:0.2f}\t{:0.2f}\t{:0.2f}'.format(covariance(Z, X), covariance(Z,
Y), covariance(Z, Z)))

```

```

1 X    Y    Z
2 -    -    -
3 2    6    1
4 4    5    3
5 6    4    5
6 8    3    7
7
8 Covariance Matrix
9      X    Y    Z
10 X    6.67   -3.33   6.67
11 Y   -3.33    1.67  -3.33
12 Z    6.67   -3.33   6.67

```

(J) Verify the equation

```

1 x = [2, 4, 6, 8, 10, 12, 14, 16, 18, 20]
2 mean = sum(x)/len(x)
3 std = math.sqrt(sum([(y-mean)**2 for y in x])/len(x))
4 print('Mean - ' + str(mean))
5 print('Standard Deviation (sd) - ' + '{}'.format(std))
6 print('Mean of Squares - ' + '{}'.format(sum([y**2 for y in x])/len(x)))
7 print('Sum of square of mean ( ' + str(mean**2) + ') and square of standard
deviation ( ' + str(std**2) + ') - ' + str(mean**2 + std**2))

```

```
1 Mean - 11.0
2 Standard Deviation (sd) - 5.744562646538029
3 Mean of Squares - 154.0
4 Sum of square of mean (121.0) and square of standard deviation (33.0) -
  154.0
```

NOTE: This is considering the sample is the population itself. Hence we use the 'Uncorrected sample standard deviation'

Q5

Import the libraries

```
1 import scipy.spatial.distance as dist
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import os
5 %matplotlib inline
```

Load the data

```
1 data = None
2 with open('Data' + os.sep + 'seeds.csv') as fp:
3     data = [x.strip().split(',') for x in fp.readlines()]
4
5 headers = data[0]
6 class_field = len(headers)-1
7 data = [[int(x[i]) if i==class_field else float(x[i]) for i in
8         range(len(x))] for x in data[1:]]
9 data = np.asarray(data)
10
11 print('Attributes - ')
12 print('\t'.join([x[:6] for x in headers]))
13 for i in range(len(data[:10])):
14     print('\t'.join(['{0:.3f}'.format(x) for x in data[i]]))
15 print('...')
16 print(str(len(data)-10) + ' more rows.')
```

```
1 Attributes -
2 area_A  perime  compac  kernel  kernel  asymme  kernel  Class
3 15.260  14.840  0.871  5.763  3.312  2.221  5.220  1.000
4 14.880  14.570  0.881  5.554  3.333  1.018  4.956  1.000
5 14.290  14.090  0.905  5.291  3.337  2.699  4.825  1.000
6 13.840  13.940  0.895  5.324  3.379  2.259  4.805  1.000
7 16.140  14.990  0.903  5.658  3.562  1.355  5.175  1.000
8 14.380  14.210  0.895  5.386  3.312  2.462  4.956  1.000
9 14.690  14.490  0.880  5.563  3.259  3.586  5.219  1.000
10 14.110  14.100  0.891  5.420  3.302  2.700  5.000  1.000
```

```

11 16.630 15.460 0.875 6.053 3.465 2.040 5.877 1.000
12 16.440 15.250 0.888 5.884 3.505 1.969 5.533 1.000
13 ...
14 200 more rows.

```

```

1 types = {'[ORIGINAL]': None, '[NORMALIZED]': None, '[STANDARDIZED]': None}
2 type_names = ['[ORIGINAL]', '[NORMALIZED]', '[STANDARDIZED]']
3 dists = {'[EUCLIDEAN]': dist.euclidean,
4         '[MAHALANOBIS]': dist.mahalanobis,
5         '[CITY BLOCK]': dist.cityblock,
6         '[MINKOWSKI (R=3)]': dist.minkowski,
7         '[CHEBYSHEV]': dist.chebyshev,
8         '[COSINE]': dist.cosine,
9         '[CANBERRA]': dist.canberra}
10 dist_funcs = (dist.euclidean,
11              dist.mahalanobis,
12              dist.cityblock,
13              dist.minkowski,
14              dist.chebyshev,
15              dist.cosine,
16              dist.canberra)
17 dist_names = ('[EUCLIDEAN]',
18              '[MAHALANOBIS]',
19              '[CITY BLOCK]',
20              '[MINKOWSKI (R=3)]',
21              '[CHEBYSHEV]',
22              '[COSINE]',
23              '[CANBERRA]')

```

(A) Select the attributes and normalize/standardize

```

1 select_headers = [headers[0], headers[4]]
2 select = data[:,[0,4]]
3
4 print('Total Observations - ' + str(len(select)) + '\n')
5 print('Attributes - ')
6 print('\t'.join([x[:6] for x in select_headers]))
7 for i in range(len(select[:10])):
8     print('\t'.join(['{0:.3f}'.format(x) for x in select[i]]))
9 print('...')
10 print(str(len(data)-10) + ' more rows.')

```

```

1 Total Observations - 210
2
3 Attributes -
4 area_A  kernel
5 15.260  3.312
6 14.880  3.333
7 14.290  3.337
8 13.840  3.379
9 16.140  3.562
10 14.380  3.312
11 14.690  3.259
12 14.110  3.302
13 16.630  3.465
14 16.440  3.505
15 ...
16 200 more rows.

```

```

1 def printrange(val, headers, dtype):
2     print(dtype + ' Range of values - ')
3     amin = np.amin(val, axis=0)
4     amax = np.amax(val, axis=0)
5     print('\t' + '\t'.join(headers))
6     print('min\t' + '\t'.join(['{0:.3f}'.format(x) for x in amin]))
7     print('max\t' + '\t'.join(['{0:.3f}'.format(x) for x in amax]))
8     print('range\t' + '\t'.join(['{0:.3f}'.format(x) for x in amax-amin]))
9     print('\n')
10    return amin, amax

```

```

1 amin, amax = printrange(select, select_headers, '[ORIGINAL]')
2 types['[ORIGINAL]'] = select
3
4 normal = np.copy(select)
5 normal = (normal - amin)/(amax - amin)
6 amin, amax = printrange(normal, select_headers, '[NORMALIZED]')
7 types['[NORMALIZED]'] = normal
8
9 amean = np.mean(select, axis=0)
10 astd = np.std(select, axis=0)
11 standard = np.copy(select)
12 standard = (standard - amean)/astd
13 amin, amax = printrange(standard, select_headers, '[STANDARDIZED]')

```

```

14 types['[STANDARDIZED]'] = standard
15

```

```

1 [ORIGINAL] Range of values -
2     area_A  kernel_width
3 min 10.590  2.630
4 max 21.180  4.033
5 range  10.590  1.403
6
7 [NORMALIZED] Range of values -
8     area_A  kernel_width
9 min 0.000   0.000
10 max 1.000   1.000
11 range  1.000   1.000
12
13 [STANDARDIZED] Range of values -
14     area_A  kernel_width
15 min -1.467  -1.668
16 max  2.182   2.055
17 range  3.648   3.723

```

```

1 print('\t'.join(type_names))
2 print('\t'.join(['\t'.join([x[:6] for x in select_headers]])*3))
3 for i in range(len(select[:20])):
4     t = []
5     for y in type_names:
6         t.append('\t'.join(['{0:.3f}'.format(x) for x in types[y][i]]))
7     print('\t'.join(t))
8 print('...')
9 print(str(len(select)-20) + ' more rows.')

```

	[ORIGINAL]		[NORMALIZED]		[STANDARDIZED]	
	area_A	kernel	area_A	kernel	area_A	kernel
3	15.260	3.312	0.441	0.486	0.142	0.142
4	14.880	3.333	0.405	0.501	0.011	0.197
5	14.290	3.337	0.349	0.504	-0.192	0.208
6	13.840	3.379	0.307	0.534	-0.347	0.320
7	16.140	3.562	0.524	0.664	0.445	0.805
8	14.380	3.312	0.358	0.486	-0.161	0.142
9	14.690	3.259	0.387	0.448	-0.054	0.001
10	14.110	3.302	0.332	0.479	-0.254	0.115
11	16.630	3.465	0.570	0.595	0.614	0.548
12	16.440	3.505	0.552	0.624	0.549	0.654

```

13 15.260  3.242  0.441  0.436  0.142  -0.044
14 14.030  3.201  0.325  0.407  -0.282 -0.153
15 13.890  3.199  0.312  0.406  -0.330 -0.158
16 13.780  3.156  0.301  0.375  -0.368 -0.272
17 13.740  3.114  0.297  0.345  -0.382 -0.384
18 14.590  3.333  0.378  0.501  -0.089  0.197
19 13.990  3.383  0.321  0.537  -0.295  0.330
20 15.690  3.514  0.482  0.630  0.290  0.678
21 14.700  3.466  0.388  0.596  -0.051  0.550
22 12.720  3.049  0.201  0.299  -0.733 -0.556
23 ...
24 190 more rows.

```

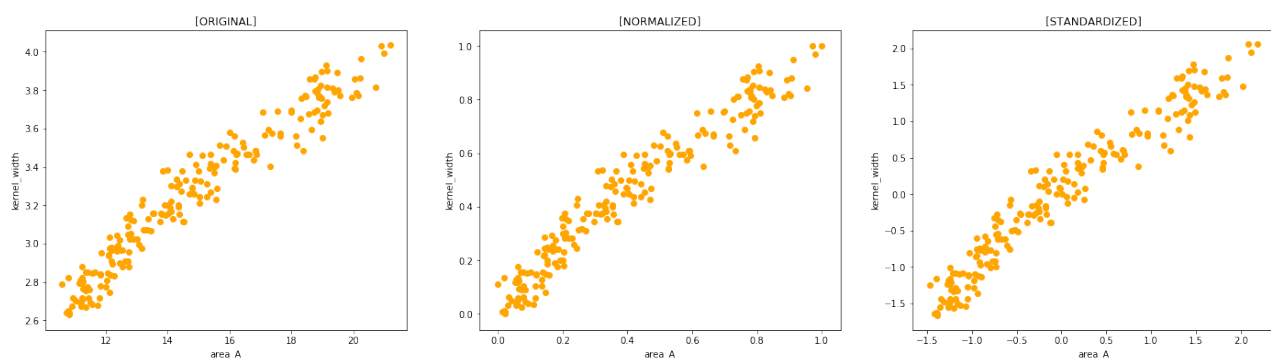
(B) Distance Metrics

(i) Scatter Plots

```

1 fig, axes = plt.subplots(nrows=1, ncols=3)
2 fig.set_figheight(6)
3 fig.set_figwidth(24)
4
5 for i, data in enumerate(type_names):
6     axes[i].set_title(data)
7     axes[i].set_xlabel(select_headers[0])
8     axes[i].set_ylabel(select_headers[1])
9     axes[i].scatter(types[data][:,0], types[data][:,1], c='orange')
10
11 plt.show()

```



Analysis

1. All the plots show that there is high colinearity between the 2 attributes.
2. Normalization and Standardization do not change the relative positioning of the elements, but scale the values.

3. Normalization only scales the values so that the range of values in the dataset is 1 and the minimum and maximum are 1 and 0 respectively. This effectively reduces the spread of the dataset to 1 unit if we plot its histogram.
4. Standardization scales the data so that the mean of the dataset is 0 and has a unit standard deviation. This helps preserve the spread of the data, allowing us to compare the distribution of this data with the standard normal distribution.

(ii) Mean values for all categories

```
1 def getmeanstd(val, headers, dtype):
2     print(dtype + ' Mean and Standard Deviation - ')
3     mean = np.mean(val, axis=0)
4     std = np.std(val, axis=0)
5     print('\t' + '\t'.join(headers))
6     print('mean\t' + '\t'.join(['{0:.3f}'.format(x) for x in mean]))
7     print('st.dev\t' + '\t'.join(['{0:.3f}'.format(x) for x in std]))
8     print()
9     return mean, std
```

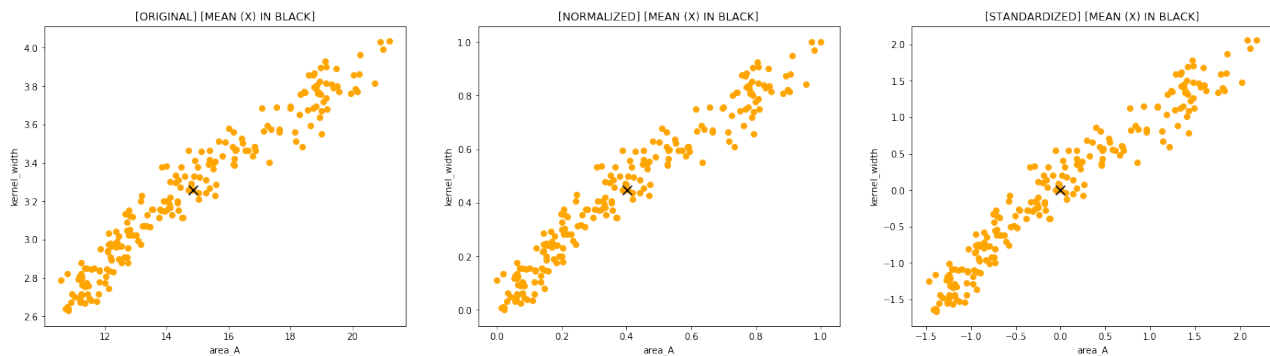
```
1 mean = {x:None for x in types.keys()}
2 std = {x:None for x in types.keys()}
3
4 for data in types.keys():
5     mean[data], std[data] = getmeanstd(types[data], select_headers, data)
6
7 fig, axes = plt.subplots(nrows=1, ncols=3)
8 fig.set_figheight(6)
9 fig.set_figwidth(24)
10
11 for i, data in enumerate(type_names):
12     axes[i].set_title(data + ' [MEAN (X) IN BLACK]')
13     axes[i].set_xlabel(select_headers[0])
14     axes[i].set_ylabel(select_headers[1])
15     axes[i].scatter(types[data][:,0], types[data][:,1], c='orange')
16     axes[i].scatter(mean[data][0], mean[data][1], s=100, c='black',
17                     marker='x')
18 plt.show()
```

```
1 [ORIGINAL] Mean and Standard Deviation -
2     area_A  kernel_width
3 mean      14.848  3.259
4 st.dev     2.903  0.377
```

```

5
6 [NORMALIZED] Mean and Standard Deviation -
7   area_A  kernel_width
8 mean    0.402  0.448
9 st.dev  0.274  0.269
10
11 [STANDARDIZED] Mean and Standard Deviation -
12   area_A  kernel_width
13 mean    -0.000 -0.000
14 st.dev   1.000  1.000

```



(iii) Distance from mean for all categories over all distance metrics

```

1 def getdists(types, mean, dist_funcs, dist_names, type_names):
2     t = []
3     for i in range(len(dist_funcs)):
4         ret = []
5         for x in type_names:
6             if dist_names[i]=='[MAHALANOBIS]':
7                 ret.append([dist_funcs[i](mean[x], y,
8 np.linalg.inv(np.cov(types[x].T)) for y in types[x]))
9             elif dist_names[i]=='[MINKOWSKI (R=3)]':
10                 ret.append([dist_funcs[i](mean[x], y, 3) for y in
11 types[x]])
12             else:
13                 ret.append([dist_funcs[i](mean[x], y) for y in types[x]])
14         ret = np.asarray(ret)
15         ret = np.transpose(ret)
16         t.append(ret)
17     return t

```

```
1 alldists = getdists(types, mean, dist_funcs, dist_names, type_names)
```

(iv) Top 10 nearest points for each distance metric

```
1 def gettop10s(alldists, type_names):
2     t = []
3     for x in alldists:
4         d = {key:None for key in type_names}
5         c = np.copy(x)
6         c.sort(axis=0)
7         c = c[9]
8         for i in range(len(type_names)):
9             d[type_names[i]] = np.where(x[:,i]<=c[i])[0]
10        t.append(d)
11    return t
```

```
1 top10s = gettop10s(alldists, type_names)
```

```
1 for i,x in enumerate(dist_names):
2     print(x)
3     for y in type_names:
4         print('\t' + y + ' : ' + ', '.join([str(n) for n in top10s[i][y]]))
5     print()
```

```
1 [EUCLIDEAN]
2     [ORIGINAL] : 1, 6, 24, 34, 38, 47, 48, 49, 55, 57
3     [NORMALIZED] : 6, 10, 24, 34, 38, 48, 49, 50, 55, 132
4     [STANDARDIZED] : 6, 10, 24, 34, 38, 48, 49, 50, 55, 132
5
6 [MAHALANOBIS]
7     [ORIGINAL] : 0, 6, 21, 24, 32, 33, 49, 53, 66, 67
8     [NORMALIZED] : 0, 6, 21, 24, 32, 33, 49, 53, 66, 67
9     [STANDARDIZED] : 0, 6, 21, 24, 32, 33, 49, 53, 66, 67
10
11 [CITY BLOCK]
12     [ORIGINAL] : 1, 6, 24, 34, 38, 47, 48, 49, 55, 57
13     [NORMALIZED] : 1, 6, 10, 24, 38, 48, 49, 50, 55, 132
14     [STANDARDIZED] : 1, 6, 10, 24, 38, 48, 49, 50, 55, 132
15
16 [MINKOWSKI (R=3)]
17     [ORIGINAL] : 1, 6, 24, 34, 38, 47, 48, 49, 55, 57
```

```

18     [NORMALIZED] : 0, 6, 10, 24, 38, 48, 49, 50, 55, 132
19     [STANDARDIZED] : 0, 6, 10, 24, 38, 48, 49, 50, 55, 132
20
21     [CHEBYSHEV]
22         [ORIGINAL] : 1, 6, 24, 34, 38, 47, 48, 49, 55, 57
23         [NORMALIZED] : 0, 5, 6, 10, 24, 38, 48, 49, 50, 55
24         [STANDARDIZED] : 0, 5, 6, 10, 24, 38, 48, 49, 50, 55
25
26     [COSINE]
27         [ORIGINAL] : 4, 20, 22, 31, 34, 43, 46, 49, 68, 134
28         [NORMALIZED] : 0, 9, 36, 49, 66, 112, 130, 139, 140, 162
29         [STANDARDIZED] : 14, 26, 66, 163, 164, 171, 182, 183, 191, 206
30
31     [CANBERRA]
32         [ORIGINAL] : 1, 6, 10, 24, 34, 38, 48, 49, 50, 55
33         [NORMALIZED] : 1, 6, 10, 24, 38, 48, 49, 50, 55, 132
34         [STANDARDIZED] : 6, 10, 18, 24, 28, 38, 48, 49, 137, 207

```

(v) Plot top 10 nearest points

```

1 def gplot(axes, i, j, title, dist, mean, headers, original, labels):
2     axes[i,j].set_title(title)
3     axes[i,j].scatter(dist[:,0], dist[:,1], s=50, c='orange')
4     axes[i,j].scatter(mean[0], mean[1], s=150, c='black', marker='x')
5     axes[i,j].set_xlabel(headers[0])
6     axes[i,j].set_ylabel(headers[1])
7     for k,label in enumerate(labels):
8         axes[i,j].annotate(label, (dist[k,0], dist[k,1]))

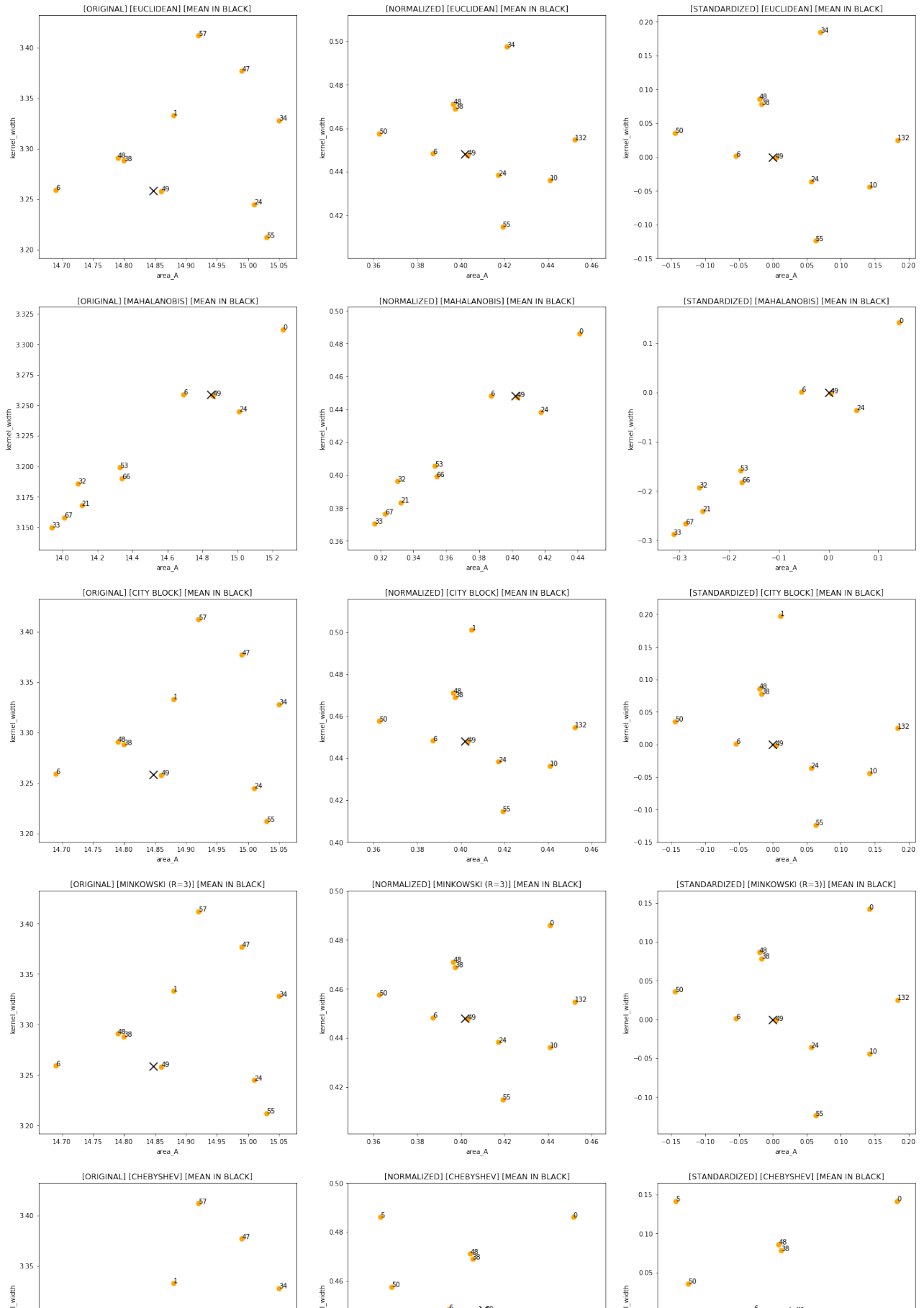
```

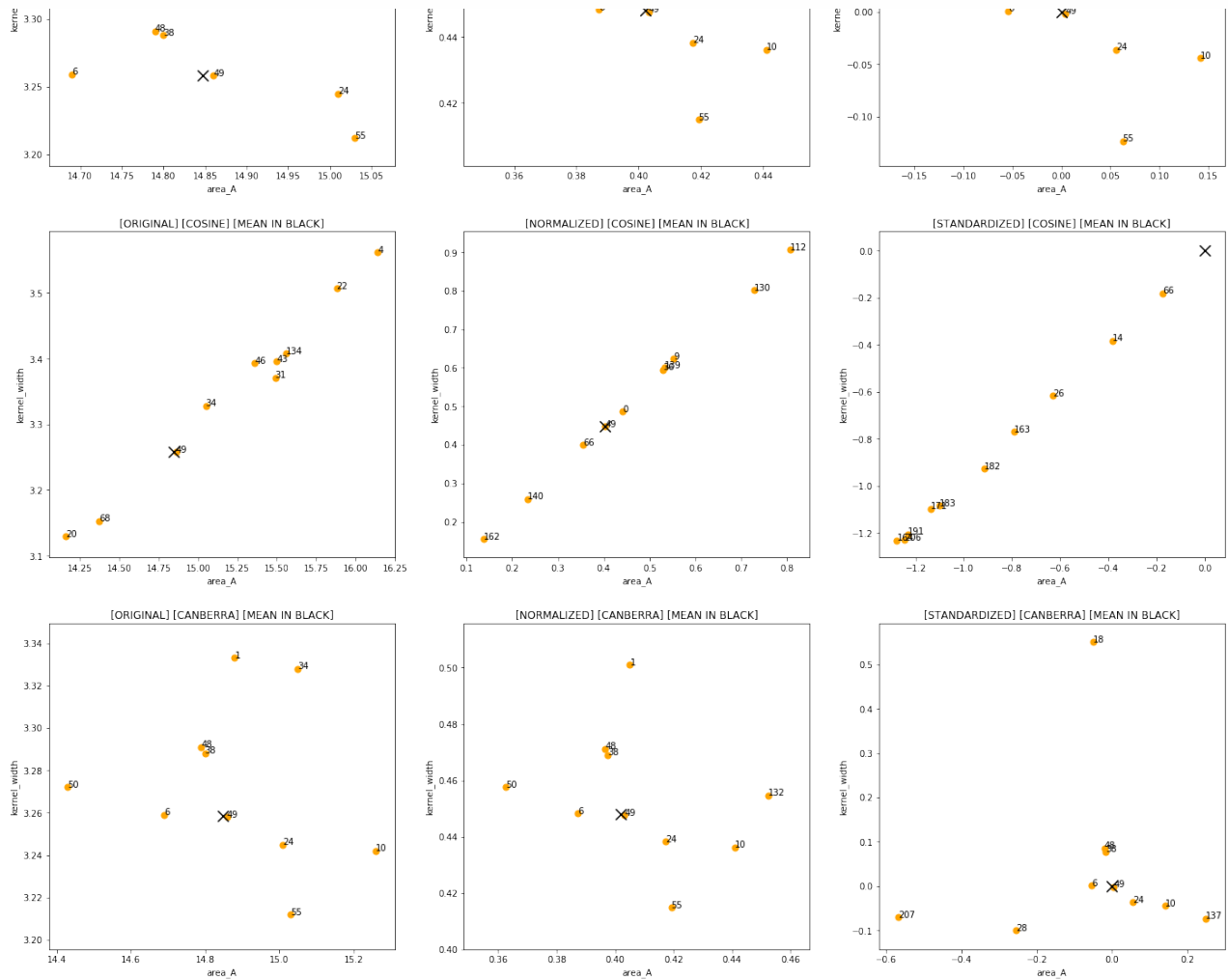
```

1 fig, axes = plt.subplots(nrows=len(alldists), ncols=len(type_names))
2 fig.set_figheight(56)
3 fig.set_figwidth(24)
4
5 for i in range(len(alldists)):
6     for j in range(len(type_names)):
7         gplot(axes, i, j, type_names[j] + ' ' + dist_names[i] + ' [MEAN IN
BLACK]', types[type_names[j]][top10s[i][type_names[j]]],
mean[type_names[j]], select_headers, types[type_names[j]], top10s[i]
[type_names[j]])
8
9 plt.show()

```

Please note that we have attached this image separately in our submission folder as Q5B5.png





(vi) Verification if the nearest points are similar across all distance metrics

```

1  s = {}
2  for y in type_names:
3      print(y)
4      t = set()
5      for i,x in enumerate(dist_names):
6          if x not in ['[MAHALANOBIS]', '[COSINE]', '[CANBERRA]']:
7              if len(t)==0:
8                  t = set(top10s[i][y])
9              else:
10                 t = t.intersection(top10s[i][y])
11                 print('\t' + x + ' : ' + ', '.join([str(n) for n in top10s[i][y]]))
12 s[y] = sorted(list(t))
13 print()

```

```

1  [ORIGINAL]
2      [EUCLIDEAN] : 1, 6, 24, 34, 38, 47, 48, 49, 55, 57
3      [MAHALANOBIS] : 0, 6, 21, 24, 32, 33, 49, 53, 66, 67
4      [CITY BLOCK] : 1, 6, 24, 34, 38, 47, 48, 49, 55, 57
5      [MINKOWSKI (R=3)] : 1, 6, 24, 34, 38, 47, 48, 49, 55, 57
6      [CHEBYSHEV] : 1, 6, 24, 34, 38, 47, 48, 49, 55, 57
7      [COSINE] : 4, 20, 22, 31, 34, 43, 46, 49, 68, 134
8      [CANBERRA] : 1, 6, 10, 24, 34, 38, 48, 49, 50, 55
9
10 [NORMALIZED]
11      [EUCLIDEAN] : 6, 10, 24, 34, 38, 48, 49, 50, 55, 132
12      [MAHALANOBIS] : 0, 6, 21, 24, 32, 33, 49, 53, 66, 67
13      [CITY BLOCK] : 1, 6, 10, 24, 38, 48, 49, 50, 55, 132
14      [MINKOWSKI (R=3)] : 0, 6, 10, 24, 38, 48, 49, 50, 55, 132
15      [CHEBYSHEV] : 0, 5, 6, 10, 24, 38, 48, 49, 50, 55
16      [COSINE] : 0, 9, 36, 49, 66, 112, 130, 139, 140, 162
17      [CANBERRA] : 1, 6, 10, 24, 38, 48, 49, 50, 55, 132
18
19 [STANDARDIZED]
20      [EUCLIDEAN] : 6, 10, 24, 34, 38, 48, 49, 50, 55, 132
21      [MAHALANOBIS] : 0, 6, 21, 24, 32, 33, 49, 53, 66, 67
22      [CITY BLOCK] : 1, 6, 10, 24, 38, 48, 49, 50, 55, 132
23      [MINKOWSKI (R=3)] : 0, 6, 10, 24, 38, 48, 49, 50, 55, 132
24      [CHEBYSHEV] : 0, 5, 6, 10, 24, 38, 48, 49, 50, 55
25      [COSINE] : 14, 26, 66, 163, 164, 171, 182, 183, 191, 206
26      [CANBERRA] : 6, 10, 18, 24, 28, 38, 48, 49, 137, 207

```

Analysis

- The top 10 points are not same for all distance metrics. Apart from **MAHALANOBIS, COSINE and CANBERRA**, almost all other metrics share more than 70% of the points.
 - For **original data**, apart from the above 3, all other distance metrics give the same points - 1, 6, 24, 34, 38, 47, 48, 49, 55, 57
 - For **normalized data**, apart from the above 3, 80% of the points are same across all metrics - 6, 10, 24, 38, 48, 49, 50, 55
 - For **standardized data**, apart from the above 3, 80% of the points are same across all metrics - 6, 10, 24, 38, 48, 49, 50, 55
- MAHALANOBIS** gives different points as it lets the data decide the coordinate system, ie - it transforms the coordinates along the 1st and 2nd PCA components(in case of 2 dimensional data) and then calculates the Euclidean distance between points in this tranformed space.
- COSINE** distance selects points that are radially close to each other from the center. In our case, one of those points is the mean, hence this metric tends to select points which are lying

on or near the line connecting the mean and origin. If the angle between 2 points from the center is near to 0, cosine reaches its maximum. Cosine Distance is $1 - \cos(\theta)$, hence this distance is minimum when the cosine is maximum.

4. **CANBERRA** distance for standardized distance is just random selection of 10 points. Canberra distance from mean (0,0) is 1 for all the points. This distance metric tends to select the points whose abscissa or ordinate is the same as the other point (in our case the mean). Hence most of the points are near the imaginary vertical and horizontal line intersecting at the mean.

(vii) Results

We can clearly see from the above scenario that most of the distance metrics are heavily affected by the range of values along each dimension of the data. If the range of one of the dimensions is much greater than the other, it tends to dominate the distance value. This might lead to unequal contribution of each dimension in calculating distance between 2 points.

Some distance metrics like Mahalanobis, Canberra and Cosine overcome this issue and give different points. Cosine is very favorable for calculating similarity between vectors. Mahalanobis is completely independent of the Euclidean co-ordinate axis, and calculates the distance along completely different co-ordinate system dictated by the variance in the data.

Thus it is important to transform the data to avoid effect of dimension with large range of values.

Normalization scales all the dimensions in the range of [0,1] Thus all dimensions have equal weightage. But this also incorporates the outliers in the actual data, thus it is not useful in scenarios which are highly sensitive to outliers.

Standardization on the other hand preserves the spread of the distribution along each dimension and also centers the mean of each dimension at 0 with a standard deviation of 1. This allows us to make comparisons with the Standard Normal Distribution and also helps eliminate outliers, while also reducing the range of data.

Q6

Import Libraries

```
1 from scipy import stats
2
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import os
6 %matplotlib inline
```

Load Data

```
1 data = None
2 with open('Data' + os.sep + 'hw1q6_data.csv') as fp:
3     data = [x.strip().split(',') for x in fp.readlines()]
4
5 headers = data[0]
6 class_field = len(headers)-1
7 data = [[int(x[i]) if i==class_field else float(x[i]) for i in
8         range(len(x))] for x in data[1:]]
9
10 print('Total Observations - ' + str(len(data)) + '\n')
11 print('Attributes - ')
12 print('\t'.join([x[:6] for x in headers]))
13 for i in range(len(data[:10])):
14     print('\t'.join(['{}'.format(x) for x in data[i]]))
15 print('...')
16 print(str(len(data)-10) + ' more rows.')
```

```
1 Total Observations - 768
2
3 Attributes -
4 Glucos  BloodP  SkinTh  BMI Diabet  Age Class
5 148.0    72.0    35.0    33.6    0.627    50.0    1
6 85.0     66.0    29.0    26.6    0.351    31.0    0
7 183.0    64.0    0.0 23.3    0.672    32.0    1
8 89.0     66.0    23.0    28.1    0.167    21.0    0
9 137.0    40.0    35.0    43.1    2.288    33.0    1
```

```

10 116.0  74.0    0.0 25.6    0.201  30.0    0
11  78.0  50.0   32.0   31.0    0.248  26.0    1
12 115.0   0.0 0.0 35.3    0.134  29.0    0
13 197.0  70.0   45.0   30.5    0.158  53.0    1
14 125.0  96.0   0.0 0.0 0.232  54.0    1
15 ...
16 758 more rows.

```

Helper Functions

```

1 def countPatients(data):
2     diabetic = 0
3     nondiabetic = 0
4
5     for i in data[:,6]:
6         if i == 1:
7             diabetic += 1
8         else:
9             nondiabetic += 1
10
11     return (diabetic, nondiabetic)

```

(A) Number of diabetic and nondiabetic patients

```

1 data = np.asarray(data)
2 diabetic, nondiabetic = countPatients(data)
3 print(f'Number of diabetic patients in the dataset: {diabetic}')
4 print(f'Number of nondiabetic patients in the dataset: {nondiabetic}')

```

```

1 Number of diabetic patients in the dataset: 268
2 Number of nondiabetic patients in the dataset: 500

```

(B) Missing rate for each feature

```

1 def countMissingValue(data):
2     missing_values = 0
3     for i in data:
4         if i == 0:
5             missing_values += 1
6     return '{:.2f}'.format(missing_values/len(data) * 100)
7
8 print('Missing Rates:')
9 for i in range(len(headers)-1):
10     percent = countMissingValue(data[:,i])
11     print(f'{headers[i]} : {percent}% ')

```

```

1 Missing Rates:
2 Glucose : 0.65%
3 BloodPressure : 4.56%
4 SkinThickness : 29.56%
5 BMI : 1.43%
6 DiabetesPedigreeFunction : 0.00%
7 Age : 0.00%

```

(C) Methods to handle missing data

Listwise Deletion

The most simple method for handling missing data is to simply remove/neglect the rows (objects) having at least one missing attribute and analyze the rest of the data.

Advantages:

- The remaining dataset is complete and accurate.
- Most frequently used method in many fields. Therefore it is also called the default option in many statistical procedures in statistical libraries.

Disadvantages:

- Reduced sample size caused by removing incomplete data
If there is not a large sample or the assumption of missing completely at random is not satisfied, this is not the optimal strategy.
- It increases chance of having a biased dataset to be substantial if data is not most completely at random.

Mean Substitution

As the name suggests, the missing data value for the attribute of a row(object) is replaced by mean value of that attribute in place of that.

Advantages:

- Allows the whole data to be used for analysis in a given incomplete dataset.
- As mean seems a reasonable estimate for a randomly selected observation from given normal distribution.

**Disadvantages*:*

- For missing values that are not strictly random, especially in presence of a great inequality in number of missing values for different variables, the mean substitution method may lead to inconsistent bias.
- May underestimate the error.

Removing missing data, rows of patient

```
1 remaining_data = data[(data[:,0:6]!=0).all(1)]
2 print('Total Observations - ' + str(len(remaining_data)) + '\n')
3 print('Attributes - ')
4 print('\t'.join([x[:6] for x in headers]))
5 for i in range(len(remaining_data[:10])):
6     print('\t'.join(['{}'.format(x) for x in remaining_data[i]]))
7 print('...')
8 print(str(len(remaining_data)-10) + ' more rows.')
```

```
1 Total Observations - 532
2
3 Attributes -
4 Glucos  BloodP  SkinTh  BMI  Diabet  Age  Class
5 148.0    72.0    35.0    33.6    0.627    50.0    1.0
6 85.0     66.0    29.0    26.6    0.351    31.0    0.0
7 89.0     66.0    23.0    28.1    0.167    21.0    0.0
8 137.0    40.0    35.0    43.1    2.288    33.0    1.0
9 78.0     50.0    32.0    31.0    0.248    26.0    1.0
10 197.0    70.0    45.0    30.5    0.158    53.0    1.0
11 189.0    60.0    23.0    30.1    0.398    59.0    1.0
12 166.0    72.0    19.0    25.8    0.587    51.0    1.0
13 118.0    84.0    47.0    45.8    0.551    31.0    1.0
14 103.0    30.0    38.0    43.3    0.183    33.0    0.0
15 ...
16 522 more rows.
```

```
1 print('After removing missing data, missing Rates:')
2 for i in range(len(headers)-1):
3     percent = countMissingValue(remaining_data[:,i])
4     print(f'{headers[i]} : {percent}%')
```

```
1 After removing missing data, missing Rates:
2 Glucose : 0.00%
3 BloodPressure : 0.00%
4 SkinThickness : 0.00%
5 BMI : 0.00%
6 DiabetesPedigreeFunction : 0.00%
7 Age : 0.00%
```

(D) Number of diabetic and nondiabetic patients in remaining data

```
1 diabetic, nondiabetic = countPatients(remaining_data)
2 print(f'Number of diabetic patients in the dataset: {diabetic}')
3 print(f'Number of nondiabetic patients in the dataset: {nondiabetic}')
```

```
1 Number of diabetic patients in the dataset: 177
2 Number of nondiabetic patients in the dataset: 355
```

(E) Summary Statistics

```

1 print('After removing missing data, Statistics:')
2 for i in range(len(headers)-1):
3     mean = np.average(remaining_data[:,i])
4     median = np.median(remaining_data[:,i])
5     std_deviation = np.std(remaining_data[:,i])
6     data_range = np.ptp(remaining_data[:,i])
7     percentile_25 = np.percentile(remaining_data[:,i],25)
8     percentile_50 = np.percentile(remaining_data[:,i],50)
9     percentile_75 = np.percentile(remaining_data[:,i],75)
10    print(f'\n{headers[i]} :-> Mean: {mean:.3f}, \n\t Median: {median:.3f},
\n\t Standard Deviation: {std_deviation:.3f}, \n\t Range: {data_range:.3f},
'
11          f'\n\t 25th percentile: {percentile_25:.3f}, \n\t 50th
percentile: {percentile_50:.3f}, \n\t 75th percentile:
{percentile_75:.3f}')

```

```

1 After removing missing data, Statistics:
2
3 Glucose :-> Mean: 121.030,
4     Median: 115.000,
5     Standard Deviation: 30.970,
6     Range: 143.000,
7     25th percentile: 98.750,
8     50th percentile: 115.000,
9     75th percentile: 141.250
10
11 BloodPressure :-> Mean: 71.506,
12     Median: 72.000,
13     Standard Deviation: 12.299,
14     Range: 86.000,
15     25th percentile: 64.000,
16     50th percentile: 72.000,
17     75th percentile: 80.000
18
19 SkinThickness :-> Mean: 29.182,
20     Median: 29.000,
21     Standard Deviation: 10.514,
22     Range: 92.000,
23     25th percentile: 22.000,
24     50th percentile: 29.000,
25     75th percentile: 36.000
26
27 BMI :-> Mean: 32.890,
28     Median: 32.800,

```

```

29     Standard Deviation: 6.875,
30     Range: 48.900,
31     25th percentile: 27.875,
32     50th percentile: 32.800,
33     75th percentile: 36.900
34
35 DiabetesPedigreeFunction :-> Mean: 0.503,
36     Median: 0.416,
37     Standard Deviation: 0.344,
38     Range: 2.335,
39     25th percentile: 0.259,
40     50th percentile: 0.416,
41     75th percentile: 0.659
42
43 Age :-> Mean: 31.615,
44     Median: 28.000,
45     Standard Deviation: 10.751,
46     Range: 60.000,
47     25th percentile: 23.000,
48     50th percentile: 28.000,
49     75th percentile: 38.000

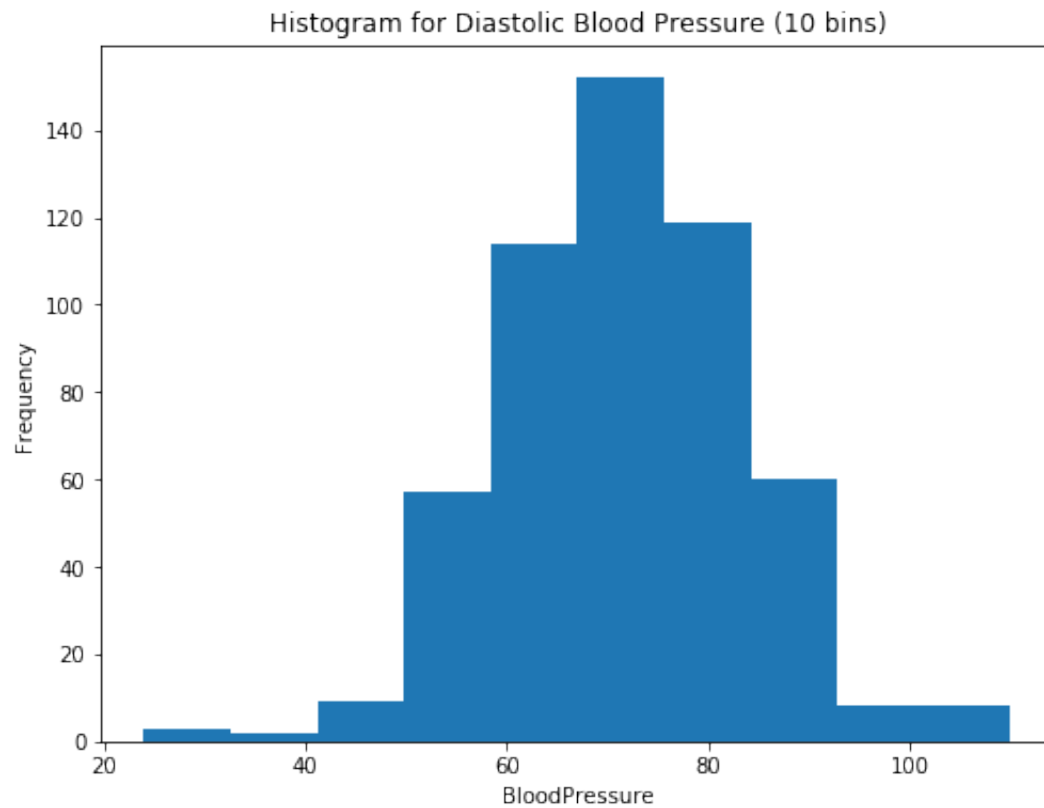
```

(F) Histograms for Blood Pressure and Diabetes Pedigree Function

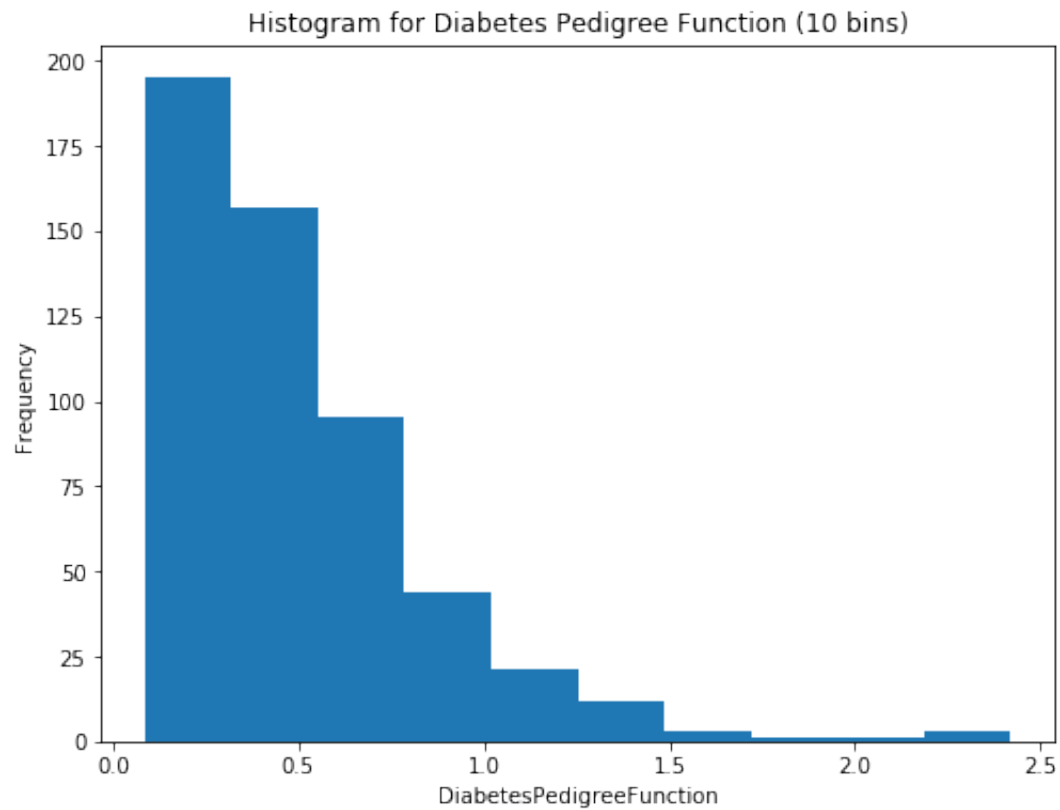
```

1  #Histogram for Blood Pressure
2  bins = np.histogram(remaining_data[:,1])
3  fig, axes = plt.subplots(nrows=1, ncols=1)
4  fig.set_figheight(6)
5  fig.set_figwidth(8)
6  axes.set_title('Histogram for Diastolic Blood Pressure (10 bins)')
7  axes.set_xlabel(headers[1])
8  axes.set_ylabel('Frequency')
9  axes.hist(remaining_data[:,1], bins[1])
10 plt.show()

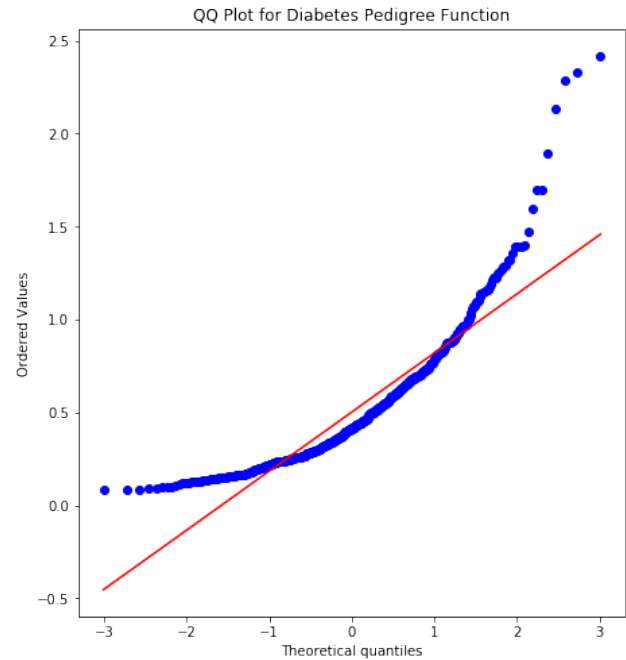
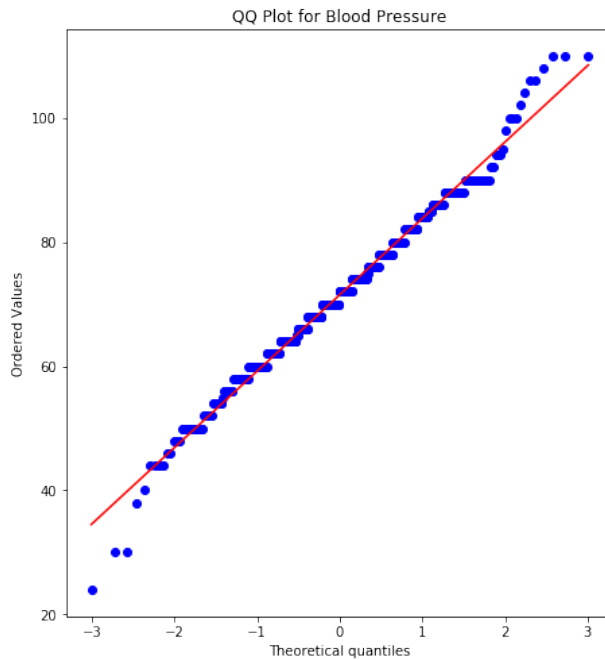
```



```
1 #Histogram for DiabetesPedigreeFunction
2 bins = np.histogram(remaining_data[:,4])
3 fig, axes = plt.subplots(nrows=1, ncols=1)
4 fig.set_figheight(6)
5 fig.set_figwidth(8)
6 axes.set_title('Histogram for Diabetes Pedigree Function (10 bins)')
7 axes.set_xlabel(headers[4])
8 axes.set_ylabel('Frequency')
9 axes.hist(remaining_data[:,4], bins[1])
10 plt.show()
```

```
1 fig, axes = plt.subplots(nrows=1, ncols=2)
2 fig.set_figheight(8)
3 fig.set_figwidth(16)
4 stats.probplot(remaining_data[:,1], plot=axes[0])
5 stats.probplot(remaining_data[:,4], plot=axes[1])
6 axes[0].set_title('QQ Plot for Blood Pressure')
7 axes[1].set_title('QQ Plot for Diabetes Pedigree Function')
8 plt.show()
```



(G) QQ Plot Analysis

QQ Plot is often used to compare the distribution of data with the standard normal distribution. This helps us in predicting the shape of the distribution and figure out how to select outliers.

If we extend the zero point on x-axis to the red line, the y co-ordinate of that intersection will be the mean of our distribution. We compare the other points with a horizontal line parallel to x-axis that passes through that intersection. The red line then dictates how far the points should be from the mean if our distribution is similar to the standard normal distribution.

Long tails mean that the values at extremes of the distribution have a higher chance of occurring when compared with the standard normal distribution. Short tails mean that there are limited values with significant probability of occurring compared to the standard normal distribution.

- QQ Plot for Blood Pressure - We can see that on the left hand of 0 on the x-axis, the values are deviating from the red line and **away** from mean. This means that the distribution of Blood Pressure has a **long left tail**. Similarly, the values on to the right of 0 are also deviating from the red line and **away** from the mean, which means that the Blood Pressure distribution has a **long right tail**.
- QQ Plot for Diabetes Pedigree Function - We can see that on the left hand of 0 on the x-axis, the values are deviating largely from the red line and **towards** the mean. This means that the distribution of Diabetes Pedigree Function has a **very short left tail**. On the contrary, the values to the right of 0 are deviating largely from the red line but **away** from the mean, which mean that Diabetes Pedigree Function distribution has a **long right tail**.

These conclusions can be backed by looking at the histograms we plotted above. They show the same structure as suggested here.