

math381HW7

Siyue Zhu

March 2022

1 Introduction

In this paper, I'm going to calculate and analyze the distance between different countries using mortality rates from Leukemia. I'm going to create one dimensional, two dimensional and three dimensional models using MDS. The data I'm going to use is mortality rates from Leukemia per million children between the ages of 0 to 14. My dataset can be find using this link: <https://people.sc.fsu.edu/~jburkardt/datasets/hartigan/file51.txt>. The dataset has 18 rows and 12 columns, representing 18 different countries through 1956 to 1967, which are 12 years in total.

mortality rates from Leukemia for each country												
Country	1956	1957	1958	1959	1960	1961	1962	1963	1964	1965	1966	1976
Australia	34	33	43	44	38	38	39	34	36	35	37	34
Austria	27	35	40	45	39	34	37	39	35	39	37	29
Belgium	41	29	31	40	39	39	47	40	34	29	34	30
Canada	36	35	32	42	35	35	35	38	36	37	33	33
Denmark	33	44	46	45	35	48	43	49	56	39	39	38
Finland	43	41	34	32	37	28	28	42	43	38	32	41
France	44	41	41	42	40	36	40	36	36	34	34	31
West Germany	35	33	35	36	35	35	34	38	36	38	34	39
Hungary	28	26	31	31	29	27	27	34	36	35	29	35
Isreal	38	44	45	28	61	33	32	28	29	36	30	30
Japan	26	25	29	28	30	29	32	32	31	33	33	31
Netherlands	32	37	39	39	39	31	41	37	40	42	36	33
Northern Ireland	17	25	27	39	29	27	22	31	33	26	30	32
Norway	45	44	36	31	44	44	56	31	36	37	50	40
Portugal	20	28	21	34	30	30	30	29	35	37	40	36
Scotland	27	31	32	30	35	30	31	26	29	28	29	25
Sweden	33	45	46	45	44	37	44	31	39	35	37	38
Switzerland	42	46	44	47	38	35	45	42	33	31	36	31

The following matrix A is our original matrix without any normalization, and I'm going to include this original matrix for our comparison.

$$A = \begin{pmatrix} 34 & 33 & 43 & 44 & 38 & 38 & 39 & 34 & 36 & 35 & 37 & 34 \\ 27 & 35 & 40 & 45 & 39 & 34 & 37 & 39 & 35 & 39 & 37 & 29 \\ 41 & 29 & 31 & 40 & 39 & 39 & 47 & 40 & 34 & 29 & 34 & 30 \\ 36 & 35 & 32 & 42 & 35 & 35 & 35 & 38 & 36 & 37 & 33 & 33 \\ 33 & 44 & 46 & 45 & 35 & 48 & 43 & 49 & 56 & 39 & 39 & 38 \\ 43 & 41 & 34 & 32 & 37 & 28 & 28 & 42 & 43 & 38 & 32 & 41 \\ 44 & 41 & 41 & 42 & 40 & 36 & 40 & 36 & 36 & 34 & 34 & 31 \\ 35 & 33 & 35 & 36 & 35 & 35 & 34 & 38 & 36 & 38 & 34 & 39 \\ 28 & 26 & 31 & 31 & 29 & 27 & 27 & 34 & 36 & 35 & 29 & 35 \\ 38 & 44 & 45 & 28 & 61 & 33 & 32 & 28 & 29 & 36 & 30 & 30 \\ 26 & 25 & 29 & 28 & 30 & 29 & 32 & 32 & 31 & 33 & 33 & 31 \\ 32 & 37 & 39 & 39 & 39 & 31 & 41 & 37 & 40 & 42 & 36 & 33 \\ 17 & 25 & 27 & 39 & 29 & 27 & 22 & 31 & 33 & 26 & 30 & 32 \\ 45 & 44 & 36 & 31 & 44 & 44 & 56 & 31 & 36 & 37 & 50 & 40 \\ 20 & 28 & 21 & 34 & 30 & 30 & 30 & 29 & 35 & 37 & 40 & 36 \\ 27 & 31 & 32 & 30 & 35 & 30 & 31 & 26 & 29 & 28 & 29 & 25 \\ 33 & 45 & 46 & 45 & 44 & 37 & 44 & 31 & 39 & 35 & 37 & 38 \\ 42 & 46 & 44 & 47 & 38 & 35 & 45 & 42 & 33 & 31 & 36 & 31 \end{pmatrix}$$

2 Calculating normalized matrices and distance matrices

Now, we are going to normalize our matrix A with two different methods. The first method is to find the mean and standard deviation of each column, and then subtract the mean and divide by the standard deviation in each column. This will transform our data into a new data set in which each column has mean zero and standard deviation 1. We can get our matrix B from the following equation:

$$B_{ij} = (A_{ij} - \mu_j) / \sigma_j$$

A second way to normalize our matrix A, is to find the min and max value in each column, and then map each column entry x to (x-min)/(max-min). This will convert the data set into a new data set in which each column as a minimum of 0 and a maximum of 1. In this way, the dimensions will be (in a certain sense) comparable. We can get our matrix c from the following equation:

$$C_{ij} = (A_{ij} - \min_j) / (\max_j - \min_j)$$

Also, we are going to calculate distance matrix X from matrix A, distance matrix Y from matrix B and distance matrix Z from matrix C. Each entry in X_{ij} , Y_{ij} and Z_{ij} representing the distance between the country at row i and the country at column j from matrices A, B and C.

Matrix X can be calculated from the following equation:

$$X_{ij} = \sqrt{\sum_{m=1}^{12} (A_{im} - A_{jm})^2}$$

Matrices Y and Z are both calculated in the same way with matrix X.

3 Results and Discussion

Let’s define model 1 as the model of MDS using matrix A, model 2 as the model of MDS using matrix B, and model 3 as the model of MDS using matrix C. The following are visualization for dimension 1 and 2 of model 1, 2 and 3. plots for dimension 3 is not showing because it is hard to visualize.

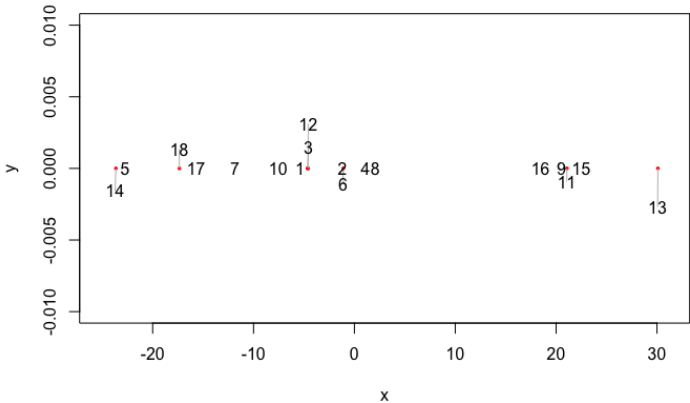


Figure 1: model 1 dimension 1

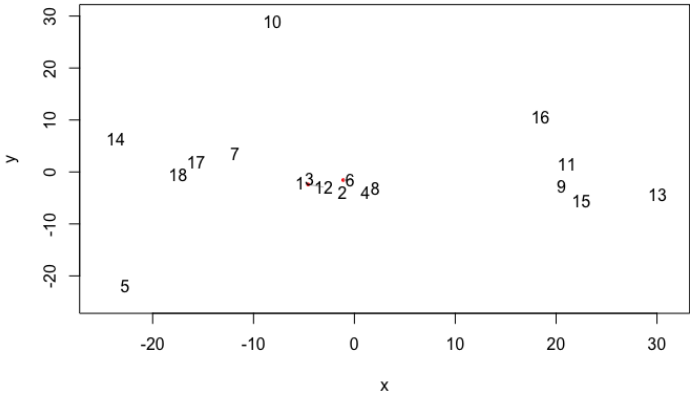


Figure 2: model 1 dimension 2

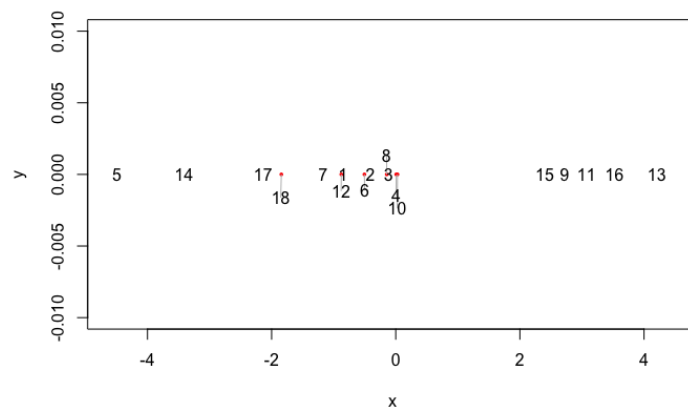


Figure 3: model 2 dimension 1

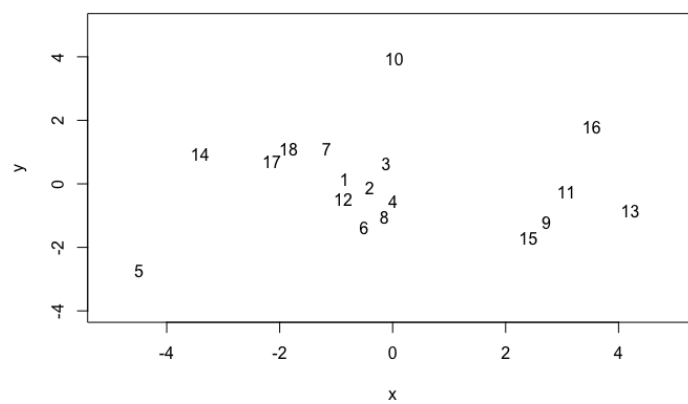


Figure 4: model 2 dimension 2

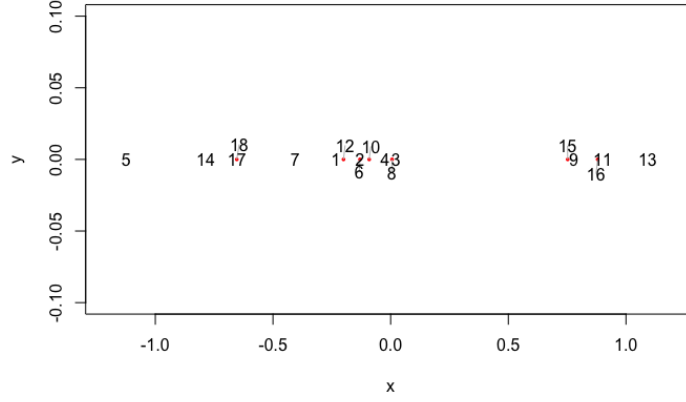


Figure 5: model 3 dimension 1

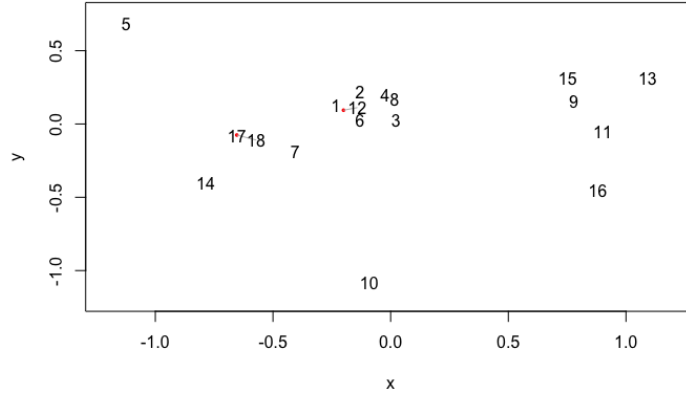
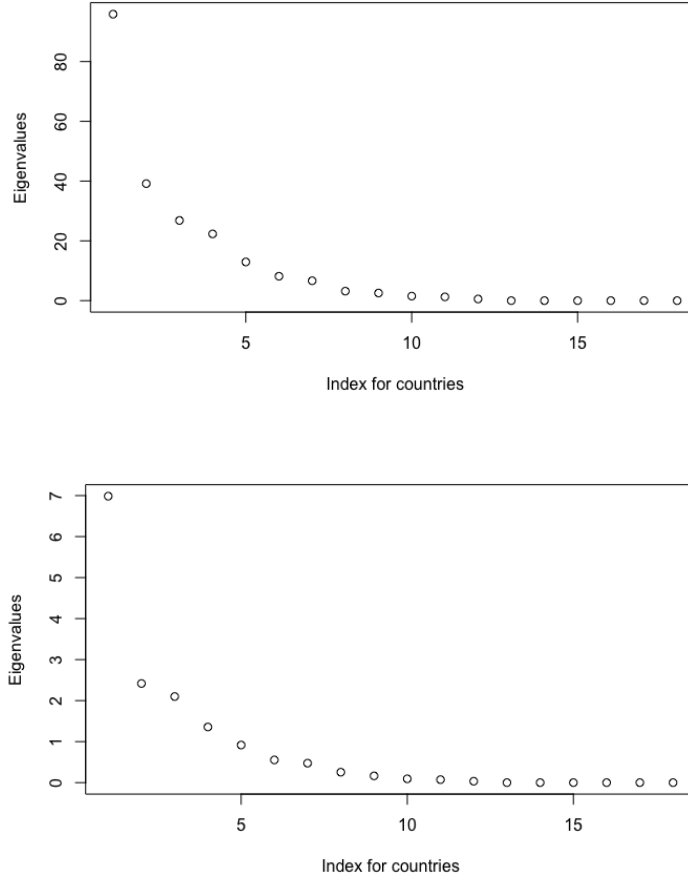


Figure 6: model 3 dimension 2

I'm not going to include model 1 in the following discussion, since model 1 is generated from a non-normalized matrix. Now, I'm going to discuss my results about model 2 and model 3, and make some comparison between model 2 and model 3.

3.1 Eigenvalues

Here is the plot of the eigenvalues for model 2 and model 3.



We can see from the plot of eigenvalues for matrix B and C, the first eigenvalue is the largest in both graphs, which means our 1 dimensional space contains the most information from the original data set. Also, for both plots of eigenvalues for matrix B and C, the first 12 indexes are greater than zero, and all other indexes are about equal to zero. This indicates that my model fits well for 12 dimensional space. We can also see that eigenvalues become very small after the first 12 values, which indicates that the seven-dimensional model is almost perfect for my data set.

3.2 Goodness of fit

There are two GOF values given by the `cmdscale` command in R, because there are different ways to handle negative eigenvalues. The first values is cal-

culated by

$$(\sum_{a=1}^{12} \lambda_a) / (\sum_{a=1}^{12} |\lambda_a|)$$

, and the second value is calculated by

$$(\sum_{a=1}^{12} \lambda_a) / (\sum_{a=1}^{12} \max(0, \lambda_a))$$

However, the two GOF values for each of our model and dimension are all equal to each other, thus we can conclude that there is no negative eigenvalue exist in our model. In the following table we are going to show one GOF for each model and each dimension.

Values of GOF			
Dimension	1	2	3
model 2	0.434	0.611	0.732
model 3	0.453	0.61	0.746

We can see that the GOF is getting larger as the dimension get larger, since larger dimensions include more information. GOF is a number between 0 and 1, the larger the GOF is, the more perfect our model is. Thus dimension 3 is better than dimension 2, and dimension 2 is better than dimension 1 for both model 1 and model 2.

3.3 Distancing Plot

The following are 6 plots of the distances in model 2 and model 3 versus the actual distances for dimension 1, 2 and 3.

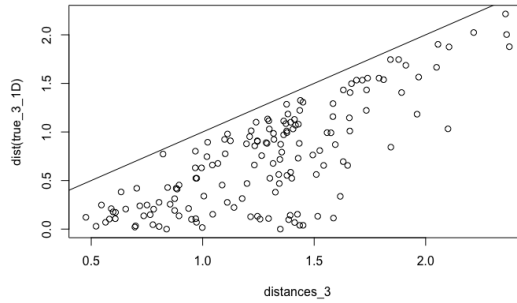


Figure 7: distances plot for 1 dimensional model for model 2

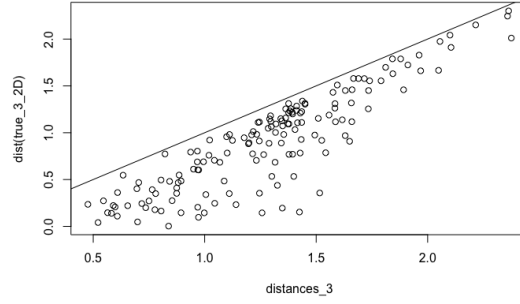


Figure 8: distances plot for 2 dimensional model for model 2

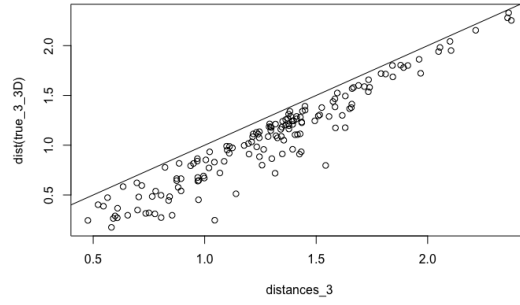


Figure 9: distances plot for 3 dimensional model for model 2

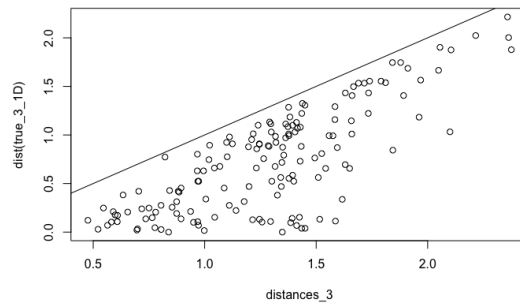


Figure 10: distances plot for 1 dimensional model for model 2

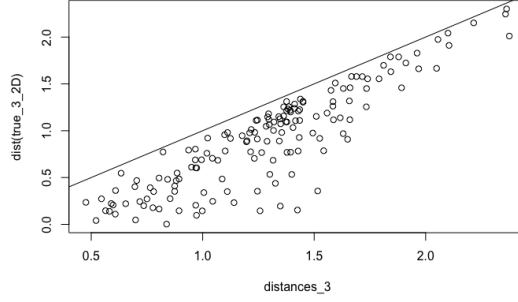


Figure 11: distances plot for 2 dimensional model for model 2

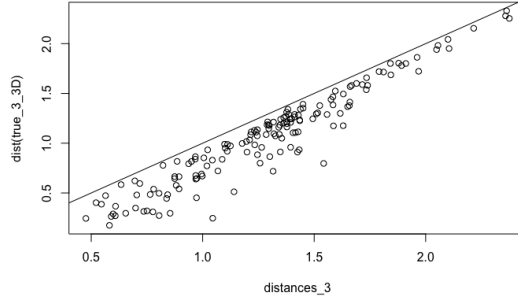


Figure 12: distances plot for 3 dimensional model for model 2

We can see from the 6 plots above, most points are close to $y=x$, which means our models make sense. Also, with the dimension gets higher, the points are getting closer to $y=x$. This tells us that our higher dimensional model are more accurate than lower dimensional model.

3.4 Mean and max absolute difference

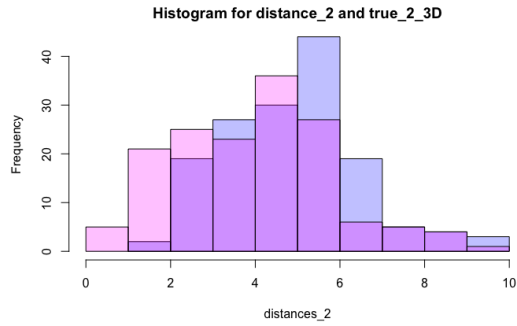
Now I'm going to calculate the mean and max absolute difference to help me better understand my model.

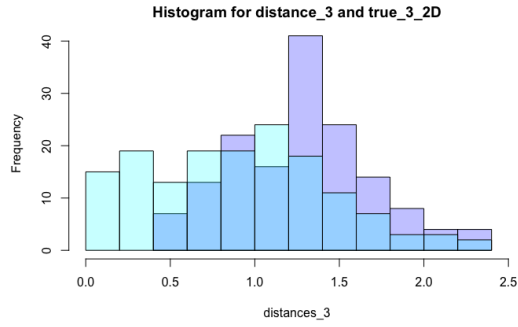
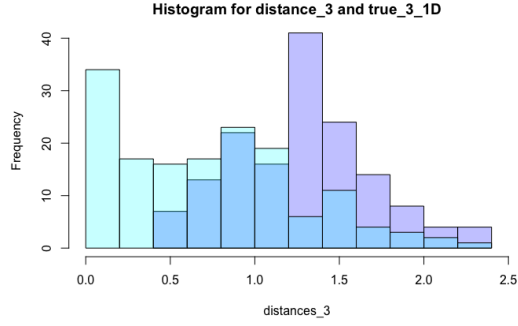
Mean and max absolute difference for model 2			
Dimension	1	2	3
Mean absolute difference	11.68936	7.611184	5.267824
Maximum absolute difference	40.20361	29.19582	26.44597

Mean and max absolute difference for model 3			
Dimension	1	2	3
Mean absolute difference	0.5197833	0.3523398	0.2084869
Maximum absolute difference	1.47299	1.272255	0.798645

We can see from the table above that both mean and max absolute difference is getting smaller with the dimension getting larger, which means that the difference between the distance calculated from our model is getting closer to the actual distance. Thus, the model with higher dimension is better than the model with lower dimension.

3.5 histograms of the original and model distances





We can see from all the histograms above with the dimension getting larger, the two histograms have more overlapping area. We can conclude that dimension 3 is the best model for both matrix B and matrix C, since the histograms for dimension 3 has the most area overlapping with the true distance.

3.6 Model dimension

The following is the table showing the correlation coefficient for x, y and z coordinates. We want to see what are the most related years corresponding to

x, y and z coordinates.

Column	X coordinate with matrix B	Y coordinate with matrix B	Z coordinate with matrix B	X coordinate with matrix C	Y coordinate with matrix C	Z coordinate with matrix C
1956	-0.699	-0.717	0.358	-0.378	-0.061	0.079
1957	-0.841	-0.882	0.319	-0.314	-0.026	0.016
1958	-0.750	-0.804	0.351	-0.227	0.302	-0.309
1959	-0.508	-0.547	-0.204	0.512	0.688	-0.601
1960	-0.464	-0.508	0.765	-0.757	-0.170	-0.042
1961	-0.841	-0.812	0.001	0.109	0.061	-0.019
1962	-0.799	-0.781	0.221	-0.105	-0.041	-0.034
1963	-0.581	-0.567	-0.509	0.578	0.441	-0.091
1964	-0.625	-0.580	0.662	0.610	0.071	0.267
1965	-0.527	-0.480	-0.319	0.177	-0.415	0.553
1966	-0.639	-0.577	-0.133	0.125	-0.441	0.385
1967	-0.464	-0.402	-0.529	0.325	-0.537	0.734

The absolute values for x-coordinate for year 1963 and 1964 are greater than 0.55 for both model B and model C. So we can assume the model in x dimension is most correlated to 1963 and 1964. And for y-coordinate, the absolute values for year 1959 are greater than 0.51 for both model B and model C. So we can assume the model in y dimension is most correlated to 1959. And for z-coordinate, the absolute values for year 1967 are greater than 0.52 for both model B and model C. So we can assume the model in y dimension is most correlated to 1967.

All my code can be found on github:

<https://github.com/zhusiyue1999/math381/tree/main>