

Homework 1

Tianyi Li

PID: tianyili

February 8, 2018

Q 1.a. (4 points) Two vectors x and y have zero mean (given $\bar{x} = \bar{y} = 0$). The cosine measure and correlation between them will be equal.

$$\cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| \cdot |\vec{y}|} \quad (1)$$

$$\begin{aligned} \text{corr}(\vec{x}, \vec{y}) &= \frac{\text{cov}(\vec{x}, \vec{y})}{\text{std}(\vec{x}) \cdot \text{std}(\vec{y})} \\ &= \frac{\frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sqrt{\frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\frac{1}{n-1} \cdot \sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{\sum_{i=1}^n (x_i \cdot y_i)}{|\vec{x}| \cdot |\vec{y}|} \\ &= \cos(\vec{x}, \vec{y}) \end{aligned} \quad (2)$$

Q 1.b. the mathematical relationship between cosine similarity and Euclidean distance when each data object vector has an L2 length (magnitude) of 1 (given $|\vec{x}| = |\vec{y}| = 1$):

$$\cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| \cdot |\vec{y}|} = \vec{x} \cdot \vec{y} \quad (3)$$

$$\begin{aligned} \text{dist}(\vec{x}, \vec{y}) &= |\vec{x} - \vec{y}| \\ &= \sqrt{(\vec{x} - \vec{y})^2} \\ &= \sqrt{(|\vec{x}|)^2 + (|\vec{y}|)^2 - 2 \cdot \vec{x} \cdot \vec{y}} \\ &= \sqrt{2 - 2 \cdot \cos(\vec{x}, \vec{y})} \end{aligned} \quad (4)$$

Such being the case, $\text{dist}(\vec{x}, \vec{y})^2 + 2 \cdot \cos(\vec{x}, \vec{y}) = 2$

Q 2. To prevent a bias towards longer documents, we first compute the *relative word frequency* $rf(w, A)$ of word w in an article set A :

$$rf(w, A) = \frac{f(w, A)}{\max\{f(w_i, A) : w_i \in A\}} \quad (5)$$

The domain weight $C_{w_i,D}$, in order to quantify the ability of word w_i in representing a targeted domain, should consider the word frequency in articles related to all events in a domain D , and compare its frequency in all articles.

Word frequency in articles related to all events in a domain D :

$$f_D(w_i) = \sum_{A \in A_D} rf(w_i, A) \quad (6)$$

Word frequency in all articles:

$$f_A(w_i) = \sum_{A \in A} rf(w_i, A) \quad (7)$$

$$\begin{aligned} C_{w_i,D} &= \frac{f_D(w_i)}{f_A(w_i)} \\ &= \frac{\sum_{A \in A_D} rf(w_i, A)}{\sum_{A \in A} rf(w_i, A)} \\ &= \frac{\sum_{A \in A_D} \frac{f(w_i, A)}{\max\{f(w_j, A): w_j \in A\}}}{\sum_{A \in A} \frac{f(w_i, A)}{\max\{f(w_j, A): w_j \in A\}}} \end{aligned} \quad (8)$$

The event weight $E_{w_i,e}$ quantifies the ability of word w_i in distinguishing event e from other events in the same domain.

Word frequency in articles related to an event e :

$$f_e(w_i) = \sum_{A \in A_e} rf(w_i, A) \quad (9)$$

Word frequency in all articles in the same domain D as event e :

$$f_{D_e}(w_i) = \sum_{A \in \{A_D: e \in D\}} rf(w_i, A) \quad (10)$$

$$\begin{aligned} E_{w_i,e} &= \frac{f_e(w_i)}{f_{D_e}(w_i)} \\ &= \frac{\sum_{A \in A_e} rf(w_i, A)}{\sum_{A \in \{A_D: e \in D\}} rf(w_i, A)} \\ &= \frac{\sum_{A \in A_e} \frac{f(w_i, A)}{\max\{f(w_j, A): w_j \in A\}}}{\sum_{A \in \{A_D: e \in D\}} \frac{f(w_i, A)}{\max\{f(w_j, A): w_j \in A\}}} \end{aligned} \quad (11)$$

Q 3. Set up a Cartesian coordinate plane on the ground with John being the origin. John's position is always $(0, 0)$, and Mike's position can be represented by (x, y) , $x, y \in [0, 1]$

Manhattan Distance between Mike and John:

$$\begin{aligned} distance &= |x| + |y| \\ min &= 1 \\ max &= \sqrt{2} \end{aligned} \quad (12)$$

Euclidean Distance between Mike and John:

$$\begin{aligned} distance &= \sqrt{x^2 + y^2} = 1 \\ min &= 1 \\ max &= 1 \end{aligned} \tag{13}$$

Chebyshev Distance between Mike and John:

$$\begin{aligned} distance &= \max(|x|, |y|) \\ min &= \frac{\sqrt{2}}{2} \\ max &= 1 \end{aligned} \tag{14}$$

Q 4.

1. Calculate Euclidean distance. Line:4 and 6. Formula for Euclidean distance:

$$\begin{aligned} EuclideanDist(\vec{x}, \vec{y}) &= |\vec{x} - \vec{y}| \\ &= \sqrt{(\vec{x} - \vec{y})^2} \\ &= \sqrt{\sum_i^n (x_i - y_i)^2} \end{aligned} \tag{15}$$

Neighborhood graphs:

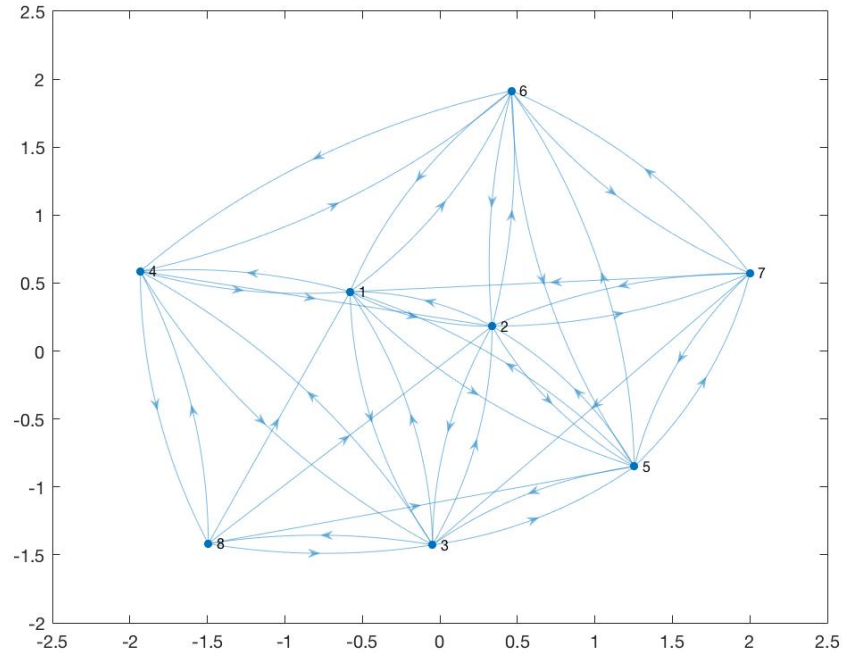


Figure 3: Q4.1a: Neighborhood graph connecting points x_i and x_j if x_i is one of the 5 nearest neighbors of x_j (8 data points)

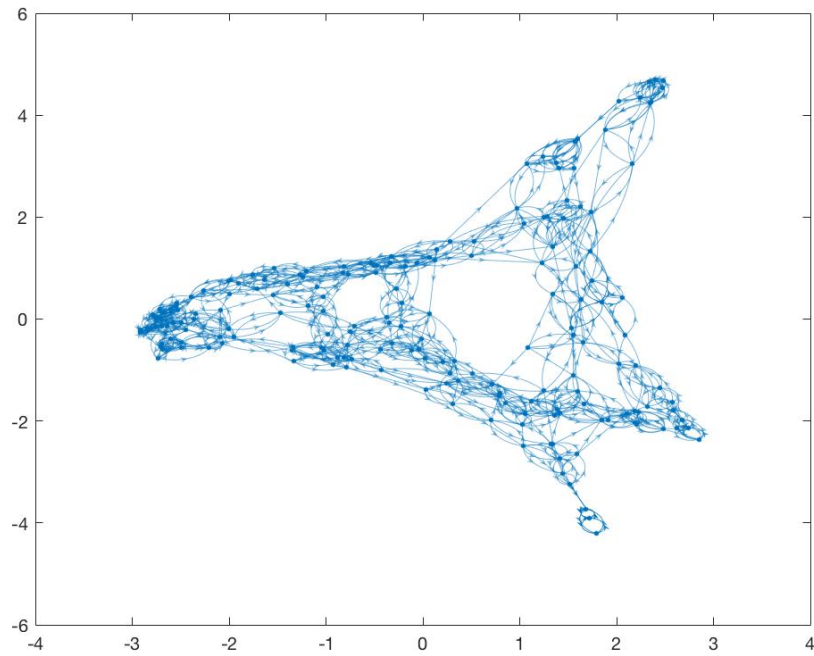


Figure 1: Q4.1a: Neighborhood graph connecting points x_i and x_j if x_i is one of the 5 nearest neighbors of x_j (200 data points)

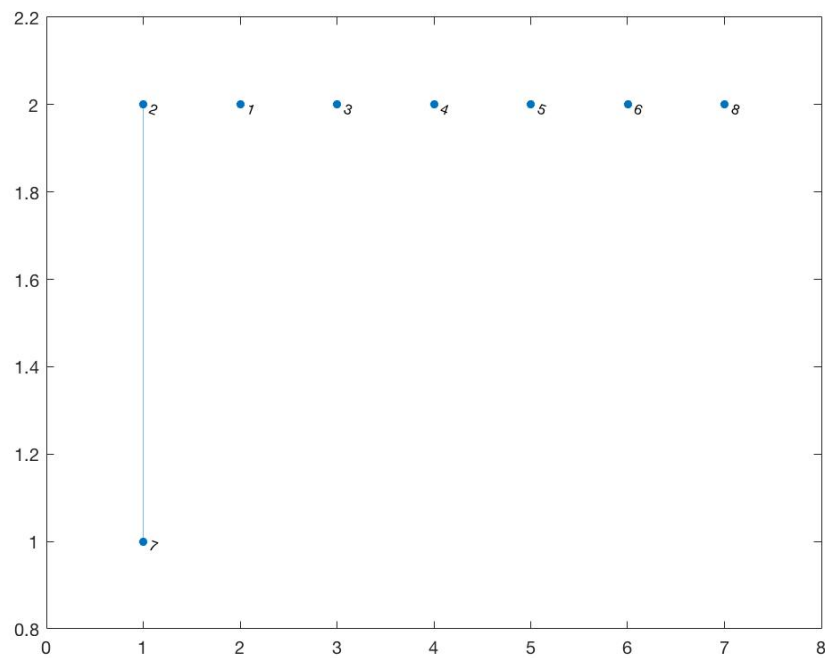


Figure 4: Q4.1b: Neighborhood graph connecting points x_i and x_j if their distance is less than 6 (8 data points)

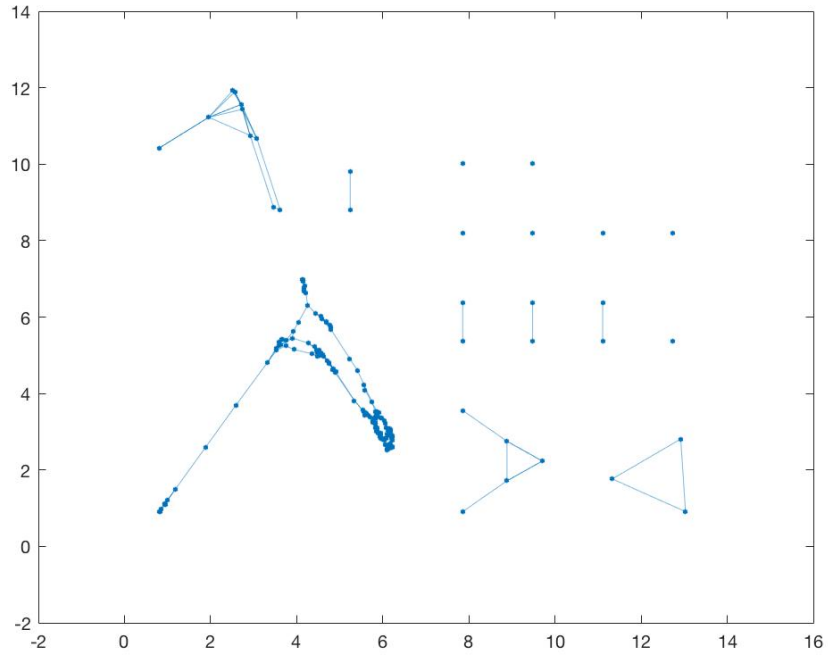


Figure 2: Q4.1b: Neighborhood graph connecting points x_i and x_j if their distance is less than 6 (200 data points)

geodesic(): line 52-63
8 by 8 distance matrices:

0	1	1	1	1	1	2	2
1	0	1	2	1	1	1	2
1	1	0	1	1	2	2	1
1	1	1	0	2	1	2	1
1	1	1	2	0	1	1	2
1	1	2	1	1	0	1	2
1	1	1	2	1	1	0	2
1	1	1	1	1	2	2	0

Table 1: Q4.2: Distance Matrix of 5 nearest neighbors graph

	0	Inf	Inf	Inf	Inf	Inf	Inf	Inf
	Inf	0	Inf	Inf	Inf	Inf	1	Inf
	Inf	Inf	0	Inf	Inf	Inf	Inf	Inf
	Inf	Inf	Inf	0	Inf	Inf	Inf	Inf
	Inf	Inf	Inf	Inf	0	Inf	Inf	Inf
	Inf	Inf	Inf	Inf	Inf	0	Inf	Inf
	Inf	1	Inf	Inf	Inf	Inf	0	Inf
	Inf	Inf	Inf	Inf	Inf	Inf	Inf	0

Table 2: Q4.2: Distance Matrix of 5 nearest neighbors graph

Code for Q4:

```

1 load data.mat
2 % the Euclidean distances between each pair of the data points (xi, xj)
3 % (a 200x200 distance matrix)
4 EuclideanDistMat = dist(X_data', X_data);
5 % distances among the first 8 data points (a 8x8 distance matrix)
6 First8 = EuclideanDistMat(1:8,1:8);
7
8 % construct neighborhood graphs
9 [adjMatLess6_200, adjMat5nn_200, distGraphLess6_200, distGraph5nn_200] =
    ↪ plotGraph(EuclideanDistMat);
10 [adjMatLess6_8, adjMat5nn_8, distGraphLess6_8, distGraph5nn_8] = plotGraph(
    ↪ First8);
11
12 % part 2: geodist
13 distMatLess6 = geodesic(distGraphLess6_8);
14 distMat5nn = geodesic(distGraph5nn_8);
15
16 function [adjMatLess6, adjMat5nn, distGraphLess6, distGraph5nn] = plotGraph
    ↪ ( graphMat )
17     size = length(graphMat);
18     sorted = zeros(200,5);
19     for i=1:size
20         tmp = sort(graphMat(i,:));
21         % minimal is oneself, counting from 2 to 6
22         sorted(i,:) = tmp(2:6);
23     end
24     adjMatLess6 = zeros(size);
25     adjMat5nn = zeros(size);
26
27     for i=1:size
28         adjMat5nn(i,:) = ismember(graphMat(i,:),sorted(i,:));
29         adjMat5nn(i,i) = 0;
30         for j=1:size

```

```

31         if graphMat(i,j) < 6
32             if i == j
33                 adjMatless6(i,j) = 0;
34             else
35                 adjMatless6(i,j) = 1;
36             end
37         end
38     end
39 end
40
41 distGraphLess6 = graph(adjMatless6);
42 distGraph5nn = digraph(adjMat5nn);
43 f1 = figure;
44 f2 = figure;
45 figure(f1);
46 plot(distGraphLess6);
47 figure(f2);
48 plot(distGraph5nn);
49
50 end
51
52 function geoDist = geodesic(distGraph)
53     size = numnodes(distGraph);
54     geoDist = zeros(size);
55     for i=1:size
56         for j=1:size
57             geoDist(i,j) = length(shortestpath(distGraph,i,j,'Method','
                    ↪ positive')) - 1;
58             if geoDist(i,j) == -1
59                 geoDist(i,j) = Inf;
60             end
61         end
62     end
63 end

```

Q 5.

(1) Boxplot of the first 100 rows of data:

Code:

```

1 thickness = csvread('thick.csv',1,1);
2 figure();
3 boxplot(thickness(1:100,:));

```

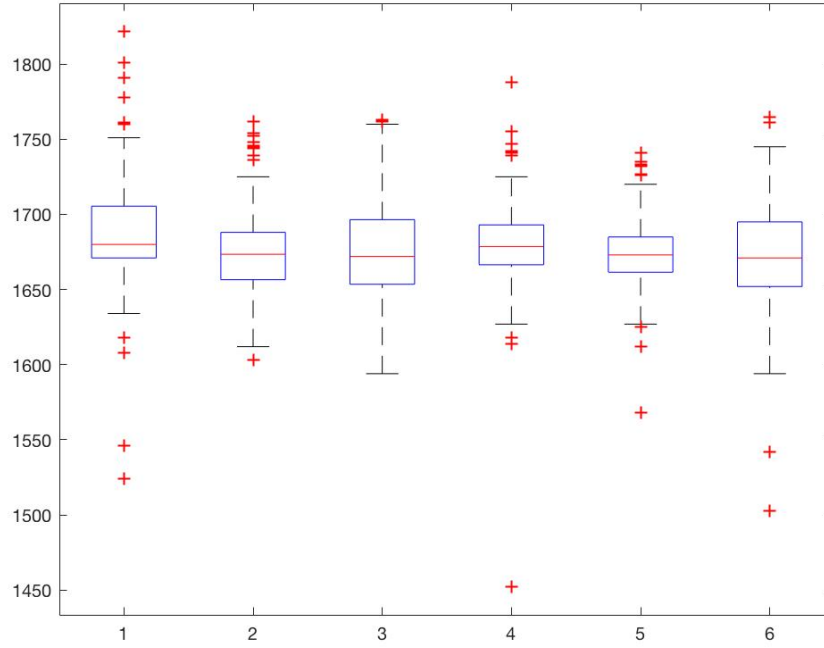


Figure 5: Q5.1: Boxplot of the first 100 rows of data

(2) Because the distribution of data values are not symmetrical. Only if there are same number of data values that are bigger than the lowest 25% of data and smaller than the highest 50%, and that bigger than the lowest 50% of data and smaller than the highest 75%.

Q 6.

1. Import: line 2
2. Suppose the whole data set $D = \{d\}$, where d is each data point. Formula for centering is $d' = d - \mu_D$, then scaling is $d'' = \frac{d'}{\sigma_D}$. Combined formula is

$$d'' = \frac{d - \mu_D}{\sigma_d}$$

. We apply this formula to each attribute (column) of the food data. (line 4-13)

3. The absolute value of normalized data points (z-score) represents the distance between the raw data value and the population mean in units of the standard deviation. z-score is negative when the raw data value is below the mean, positive when above.

4. Formula for correlation matrix in z-score (z_i^x is normalized x_i for each data point x_i ,

\vec{z}_x is normalized vector (column) \vec{x} (implemented in line 15-24):

$$\begin{aligned}
corr(\vec{x}, \vec{y}) &= \frac{cov(\vec{x}, \vec{y})}{std(\vec{x}) \cdot std(\vec{y})} \\
&= \frac{\frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \mu_x) \cdot (y_i - \mu_y)}{\sigma_{\vec{x}} \cdot \sigma_{\vec{y}}} \\
&= \frac{1}{n-1} \cdot \sum_{i=1}^n \left(\frac{(x_i - \mu_x)}{\sigma_{\vec{x}}} \cdot \frac{(y_i - \mu_y)}{\sigma_{\vec{y}}} \right) \\
&= \frac{1}{n-1} \cdot \sum_{i=1}^n (z_i^x \cdot z_i^y) \\
&= \frac{1}{n-1} \cdot \vec{z}_x \cdot \vec{z}_y
\end{aligned} \tag{16}$$

5. line 27

6. Plot the percentage of variance captured by the individual components in decreasing order: line 29-43

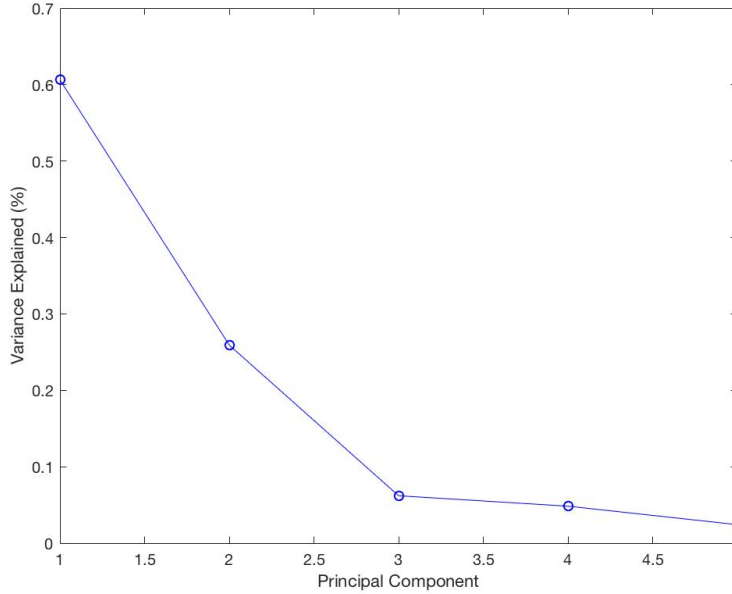


Figure 6: scree plot of the eigenvalues energies

7. I would use 2 component because the inflection point happens at the third component, and the first two components combined capture a reasonably high (close to 90%) percentage of variance. These steps are PCA algorithm, retaining $(\lambda_1^2 + \lambda_2^2) / (\lambda_1^2 + \lambda_2^2 + \lambda_3^2 + \lambda_4^2 + \lambda_5^2) = 10.8671 / 11.0364 = 98.47\%$ energy. ($\lambda_{1,2,3,4,5}$ are eigenvalues sorted in descending order)

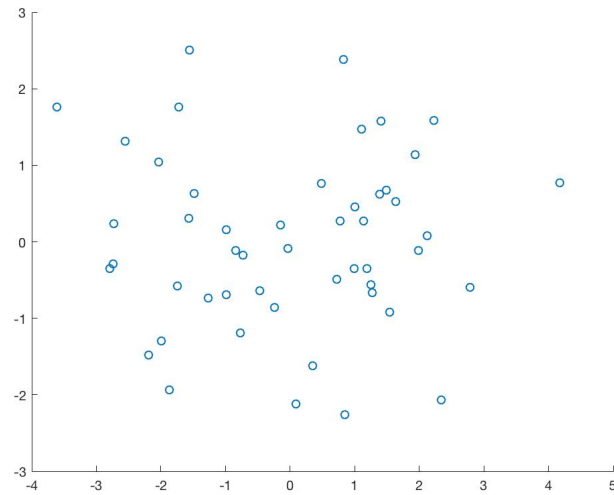


Figure 7: Projecting matrix to 2-D project with PCA

9. PCA is not recommended for high dimensional data, because it will be too expensive to compute too many eigenvalues and eigenvectors for each dimension.

10. SVD

11.2-D Scatter Plot for SVD

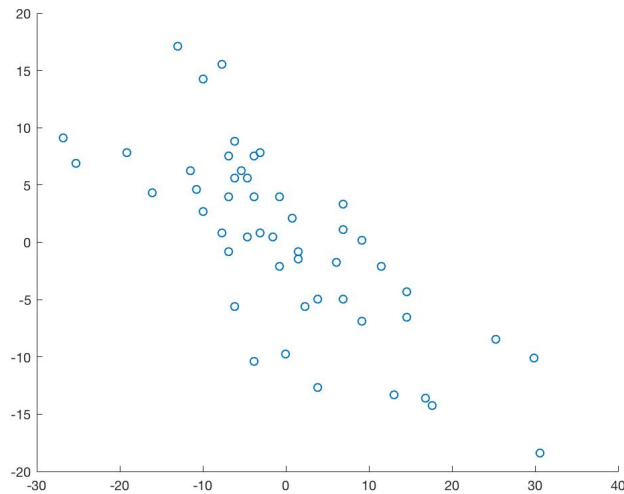


Figure 8: Projecting matrix to 2-D project with SVD

12. Yes.

the PCA algorithm computes the eigenvalues and eigenvectors of covariance matrix of the data, which is equivalently the correlation matrix of the standard normalized data. The covariance matrix of a data matrix X is $X \cdot X^T$ and is symmetric, the matrix is diagonalizable,

and the eigenvectors can be normalized such that they are orthonormal: $X \cdot X^T = W \cdot D \cdot W^T$, where W is the matrix of eigenvectors and D is the diagonal matrix of eigenvalues.

On the other hand, applying SVD to the data matrix X , we have $X = S \cdot \Sigma \cdot U^T$. The covariance matrix from this decomposition, given that $U^T \cdot U = I$

$$X \cdot X^T = S \cdot \Sigma \cdot U^T \cdot (S \cdot \Sigma \cdot U^T)^T = S \cdot \Sigma \cdot U^T \cdot U \cdot \Sigma \cdot S^T = S \cdot \Sigma^2 S^T$$

```

1  % 1. (2 points) Import: Import the data file food.csv in Matlab.
2  rawfood = csvread('food.csv',1,1);
3
4  % 2. Processing: Centering
5  food = rawfood;
6  meanFood = mean(food);
7  meanFood = repmat(meanFood,length(food(:,1)),1);
8  centeredFood = food-meanFood;
9
10 % 2. Processing: Scaling centered data
11 stdFood = std(food);
12 stdFood = repmat(stdFood,length(food(:,1)),1);
13 scaledFood = centeredFood./stdFood;
14
15 %4. customized correlation matrix
16 numDim = length(food(1,:));
17 numObs = length(food(:,1));
18 corrMat = zeros(numDim);
19 n = length(food(:,1));
20 for i=1:numDim
21     for j=1:numDim
22         corrMat(i,j) = (1/(numObs-1)) * scaledFood(:,i)'*scaledFood(:,j);
23     end
24 end
25
26 % 5. eigenvalues and eigen vectors
27 [V,D] = eigs(corrMat);
28
29 % 6. sort eigen values from largest to smalelst
30 [c, ind]=sort(diag(D),'descend');
31 D2=diag(c);
32
33 % 6. accordingly update the order of the eigenvectors in matrix
34 V2 = V(:,ind);
35
36 % Plot the percentage of variance captured by the individual components
37 % in decreasing order
38 allDim = length(D2);

```

```

39 cumG = c/sum(c);
40 figure();
41 plot(cumG,'b-o');
42 xlabel('Principal Component')
43 ylabel('Variance Explained (%)')
44
45 numComponents = 2;
46 projMatPCA = V2(:,1:numComponents);
47 figure();
48 scatter(scaledFood*projMatPCA(:,1),scaledFood*projMatPCA(:,2));
49
50 % 10. SVD
51 [U,S,V] = svd(scaledFood);
52 % sort s matrix from largest to smalelst
53 [c, ind]=sort(diag(S),'descend');
54 S2=diag(c);
55 % project
56 projMatSVD = S2(1:numComponents,:);
57 figure()
58 scatter(scaledFood*projMatSVD(:,1), scaledFood*projMatSVD(:, 2));

```