

Automated Learning and Data Analysis

Team G40

TEAM MEMBER 1 Sujith Tumma,stumma2

TEAM MEMBER 2 Vishnu Challa,vchalla2

TEAM MEMBER 3 Indu Chenchala,ichench

April 26, 2022

1 SVM

1. 1a) Support vector machines (SVM) learn a decision boundary leading to the largest margin between classes. In this question, you will train a SVM on a tiny dataset with 4 data points, shown in Figure 1. This dataset consists of two points with Class 1 ($y = -1$) and two points with Class 2 ($y = 1$). Each data point has two non-class attributes: x_1 and x_2 .
 - i) (5 points) Find the weight vector w and bias w_0 for the decision boundary of a hard-margin SVM. What is the equation corresponding to this decision boundary? Show your work, including the equations you used to derive your answer

2) SVM

1a) Given

Margin Plane Constraint

$$y_i (\omega^T \mathbf{x}_i + \omega_0) - 1 \geq 0$$

Class -1 ($y = -1$)Class -2 ($y = 1$)

Two closest points in the plane

for $y = -1$, $P = (5, 2)$. $y = 1$, $P = (2, 4)$

$$\Rightarrow -1(5\omega_1 + 2\omega_2 + \omega_0) - 1 = 0 \rightarrow \textcircled{7}$$

$$1(2\omega_1 + 4\omega_2 + \omega_0) - 1 = 0 \rightarrow \textcircled{8}$$

on solving both the above equations

$$5\omega_1 + 2\omega_2 + \omega_0 = -1$$

$$2\omega_1 + 4\omega_2 + \omega_0 = 1$$

$$\underline{\begin{array}{r} (-) \\ (-) \end{array}} \quad \underline{\begin{array}{r} (+) \\ (+) \end{array}} \quad \underline{\begin{array}{r} (-) \\ (-) \end{array}}$$

$$3\omega_1 - 2\omega_2 = -2 \rightarrow \textcircled{1}$$

The margin is always perpendicular to the slope of line joining our support vectors.

slope of line joining SVs $\Rightarrow \frac{(y_2 - y_1)}{(\mathbf{x}_2 - \mathbf{x}_1)} \Rightarrow \frac{(4 - 2)}{(2 - 5)}$



$m_1 = \left(-\frac{2}{3}\right)$
 Let slope of margin be m_2

$$\therefore m_1 \times m_2 = -1$$

$$-\frac{2}{3} \times m_2 = -1$$

$$\boxed{m_2 = \frac{3}{2}} \rightarrow \textcircled{2}$$

since the equation of Margin is

$$\omega^T x + \omega_0 = 0$$

$$\omega_1 x_1 + \omega_2 x_2 + \omega_0 = 0$$

$$x_2 = \left(-\frac{\omega_1}{\omega_2}\right) x_1 + \left(-\frac{\omega_0}{\omega_2}\right) \rightarrow \textcircled{3}$$

i.e it is in $y = mx + c$ format. Combining $\textcircled{2}$ & $\textcircled{3}$

$$-\frac{\omega_1}{\omega_2} = \frac{3}{2}$$

$$\boxed{\omega_1 = \left(-\frac{3}{2}\right) \omega_2} \rightarrow \textcircled{4}$$

Substitute $\textcircled{4}$ in $\textcircled{1}$

$$3 \times \frac{-3}{2} \omega_2 - 2 \omega_2 = -2$$

$$= 9 \omega_2 - 4 \omega_2 = -4$$

$$-13 \omega_2 = -4$$

$$\boxed{\omega_2 = \frac{4}{13}} \rightarrow \textcircled{5}$$



$$\therefore \text{from } ④ \text{ & } ⑤$$

$$w_1 = \left(\frac{-6}{13} \right) \rightarrow ⑥$$

in equation ⑦

Substitute w_1 & w_2

$$= 1 \times \left(5 \times \frac{-6}{13} + 2 \times \frac{4}{13} + w_0 \right) - 1 = 0$$

$$-\frac{30}{13} + \frac{8}{13} + w_0 = (-1)$$

$$-\frac{22}{13} + w_0 = -1$$

$$w_0 = -1 + \frac{22}{13}$$

$$w_0 = \left(\frac{9}{13} \right) \rightarrow ⑨$$

\therefore substituting values of w_0, w_1 & w_2

in $w_1 x_1 + w_2 x_2 + w_0 = 0$

$$-\frac{6}{13} x_1 + \frac{4}{13} x_2 + \frac{9}{13} = 0$$

$$(-6x_1 + 4x_2 + 9) = 0$$

Decision boundary equation
margin

\therefore Our final equation for decision boundary is

$$y_i (-6x_1 + 4x_2 + 9) - 1 = 0 \rightarrow ⑩$$

where

$$y_i = \begin{cases} +1 & (\text{class 2}) \\ -1 & (\text{class 1}) \end{cases} \text{ and}$$

- ii) (4 points) Circle the support vectors and draw the decision boundary.

10(ii)

The equation of decision boundary is

$$-6x_1 + 4x_2 + 9 = 0$$

$$-6x_1 + 4x_2 = -9$$

$$\frac{x_1}{\left(\frac{-9}{-6}\right)} + \frac{x_2}{\left(\frac{-9}{4}\right)} = 1$$

$$\frac{x_1}{\left(\frac{3}{2}\right)} + \frac{x_2}{\left(-\frac{9}{4}\right)} = 1$$

\therefore The above equation is in the form of

$$\frac{x}{a} + \frac{y}{b} = 1$$

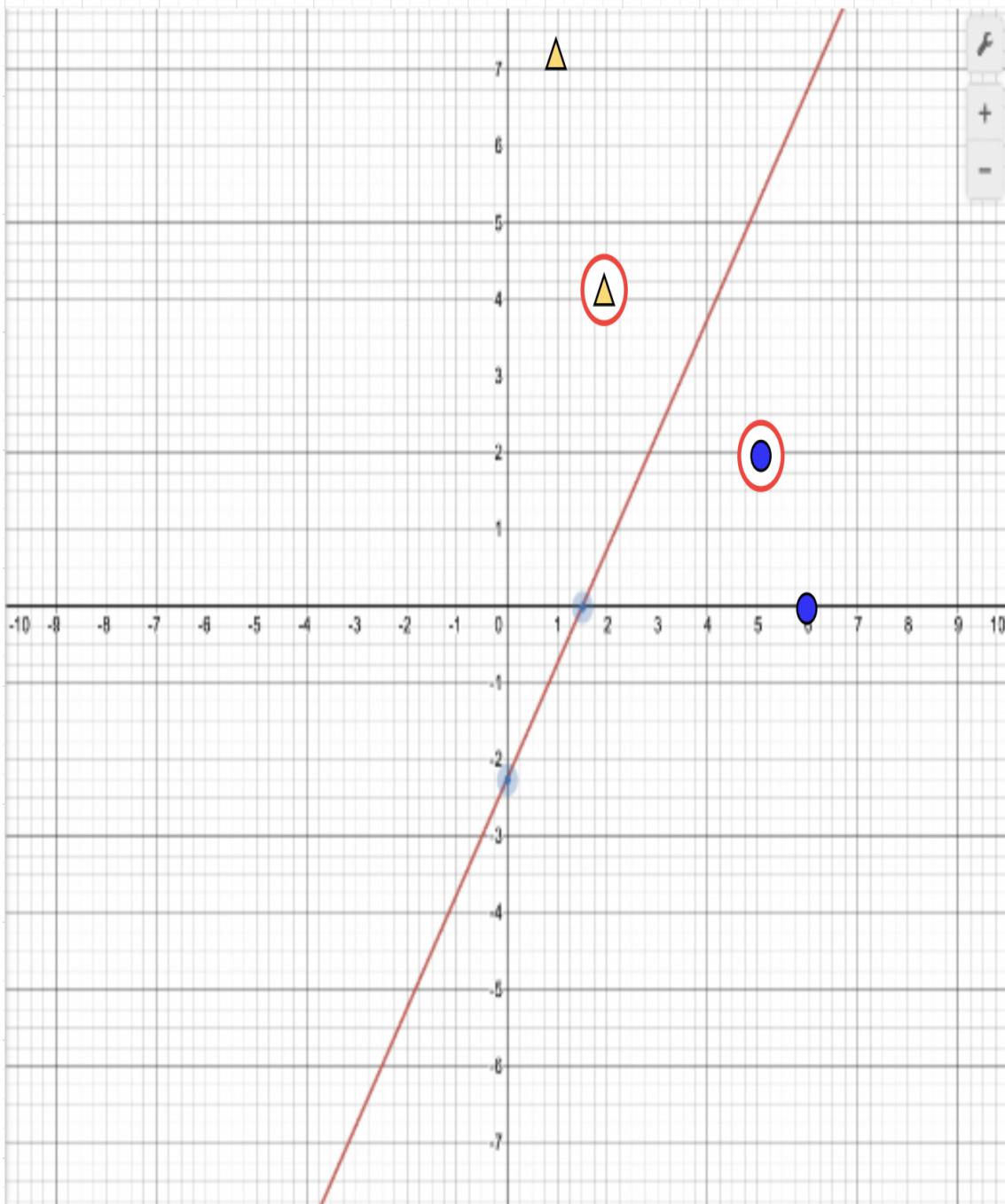
where a is x -intercept & b is y -intercept.

Our line is passing through $(\frac{3}{2}, 0)$ & $(0, -\frac{9}{4})$ and has a slope of $(-\frac{b}{a})$ which is

$$\frac{\frac{+9}{4}}{\frac{3}{2}} = \frac{3}{4} \times \frac{2}{3} = \frac{3}{2}$$

$\text{Slope} = \left(\frac{3}{2}\right)$





2. 1b) This whole question is required for CSC522; Extra Credit for CSC 422: You are given 1- dimensional data points X_i , $i \in [1, 2, 3, 4, 5, 6, 7]$ as shown in Table 1 ,also shown in Figure 2 in this question.

- i)(1 point) Calculate the equation for the decision boundary of a hard-margin SVM, or if this is not possible, explain why in 1-2 sentences.

In one-dimensional space, there is no straight line that rigorously separates the two classes, according to observation. The negative class blue and the positive class orange are not linearly separable. As data is not linearly separable, a hard margin SVM algorithm will not converge.

- ii) (1 point) If you were to train a soft-margin SVM on this data, would you select a C value where $C \rightarrow 0$ or $C \rightarrow \infty$. Explain why in 1 sentence.

(ii) While training a soft-margin SVM our goal is to minimize the below equation.

$$\frac{\|w\|^2}{2} + C \sum_{i=1}^n \xi_i$$

If $C \rightarrow 0$, then we are cancelling out the 2nd term of the equation (the one with $C \sum \xi_i$)

indicating no penalty for having slack. That means we can make margin maximized. Eventually every data point is going to be a support vector. Similar to a Soft-margin.



Scanned with CamScanner

If $C \rightarrow \infty$, then we will have infinite cost for slack penalty which means even for a slightest bit of slack we have gotten a huge error value which is not good. This is similar to a hard margin.

If we were to train soft-margin SVM we would select $C \rightarrow 0$.

- iii) (1 point) Imagine you want to transform the 7 given data points to a higher dimensional space. You decide to use the kernel function $K(X_i, X_j) = ((1 + (9/2) X_i X_j)^2)$

(iii) Given 1-dimensional data.

$$x_i : x_1$$

$$x_j : x_2$$

Given kernel function.

$$K(x_i, x_j) = \left(1 + \frac{9}{2}x_i x_j\right)^2 = \phi(x_i) \cdot \phi(x_j)$$

$$\Rightarrow 1 + \frac{81}{4}x_1^2 x_2^2 + 9x_1 x_2$$

$$\Rightarrow 1 + \left(\frac{3}{\sqrt{2}}x_1\right)^2 \cdot \left(\frac{3}{\sqrt{2}}x_2\right)^2 + (3x_1) \cdot (3x_2)$$

$$\Rightarrow \begin{bmatrix} \left(\frac{3}{\sqrt{2}}x_1\right)^2 \\ 3x_1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} \left(\frac{3}{\sqrt{2}}x_2\right)^2 \\ 3x_2 \\ 1 \end{bmatrix}$$

$$\text{Hence } \phi(x) = \begin{bmatrix} \left(\frac{3}{\sqrt{2}}x_1\right)^2 \\ 3x_1 \\ 1 \end{bmatrix}, \begin{bmatrix} \left(\frac{3}{\sqrt{2}}x_2\right)^2 \\ 3x_2 \\ 1 \end{bmatrix}$$

Hence the dimension of $\phi(x) = 3$,



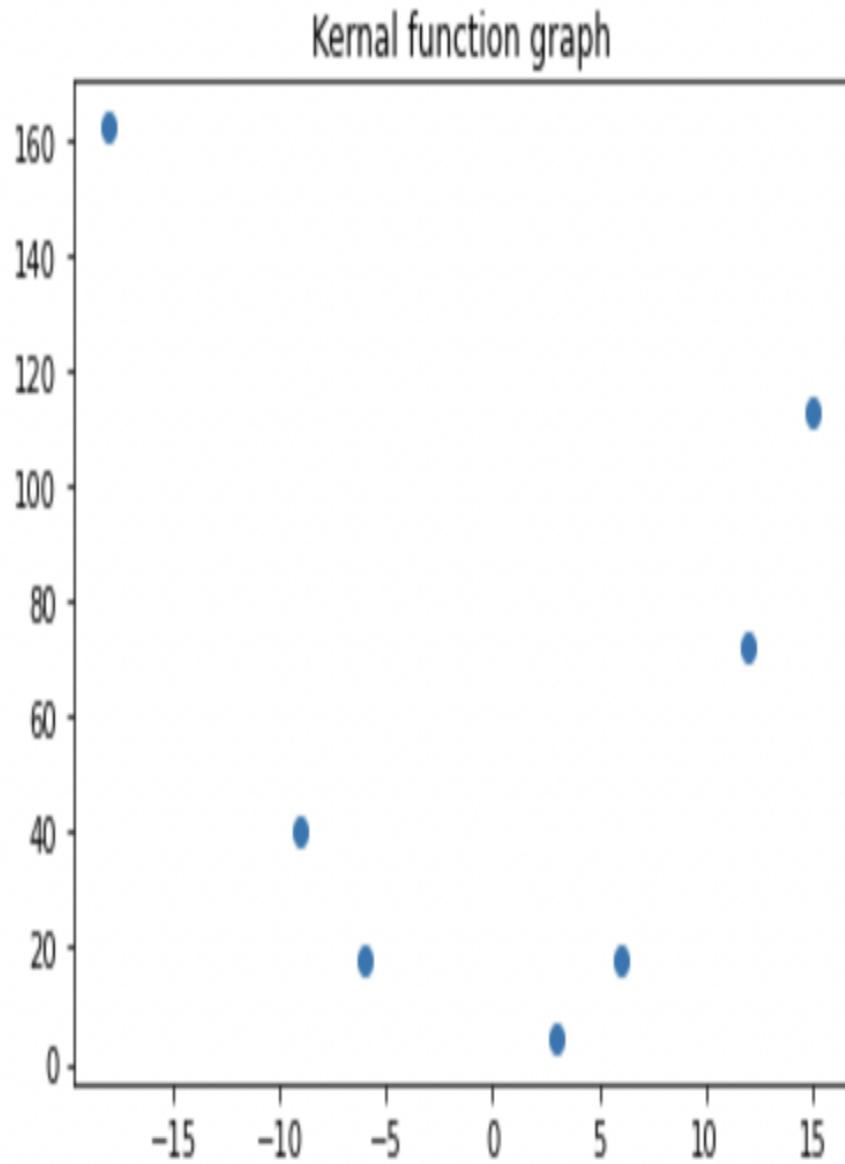
Scanned with CamScanner

- iv) (1 point) Use the function $\Phi(X)$ to calculate $\Phi(x_i)$ for $i \in [1, 2, 3, 4, 5, 6, 7]$. Graph these data points in the higher-dimensional space. (Hint: If the data is more than 2-dimensional, can you simplify your visualization to show it in 2D?).

From the above kernel function 1 is vector constant. So we will be discarding 1 and will only represent $((9/2)^*(X))^2$ and $3X$ in the graph. Below is the coordinates table and the graph.

X	$3X$	$(9/2) * X^2$
-6	-18	162
-3	-9	40.5
-2	-6	18
1	3	4.5
2	6	18
4	12	72
5	15	112.5

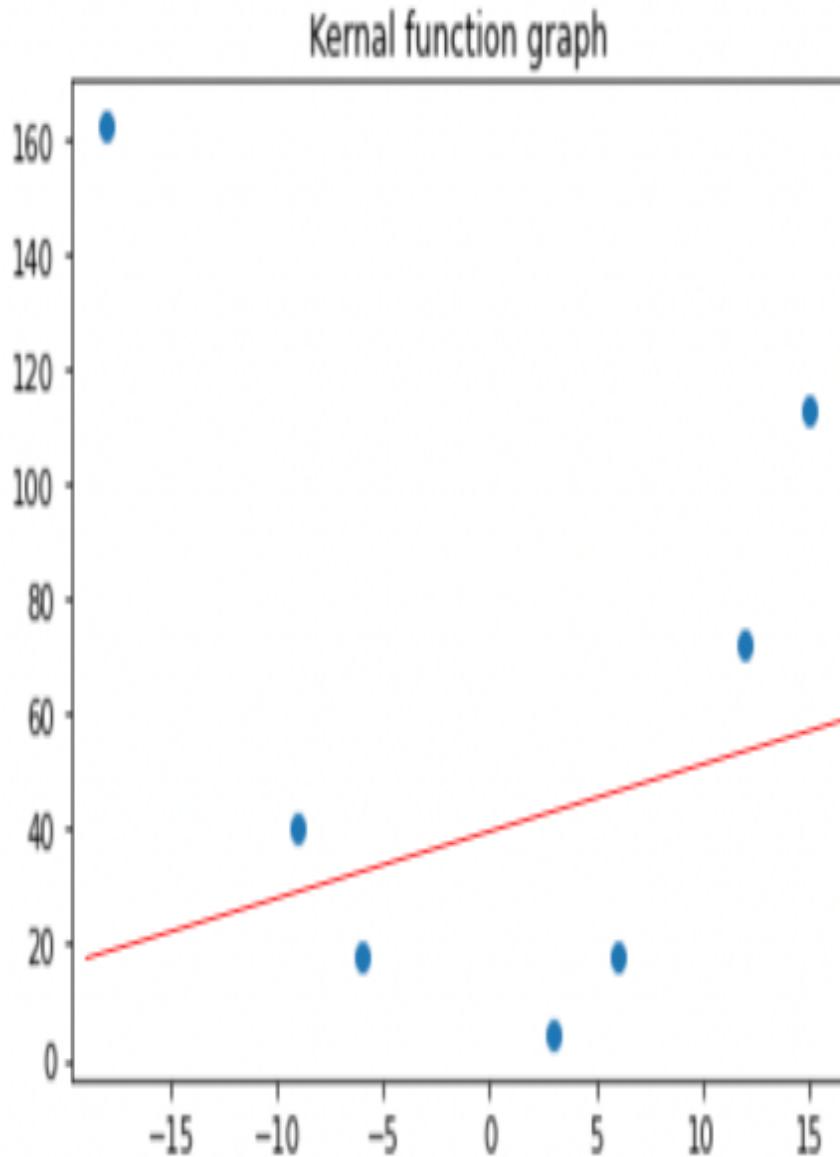
<matplotlib.collections.PathCollection at 0x7f0ce29a8cc0>



v) (1 point) Is it possible to linearly separate the data in the higher-dimensional space? If so, draw the decision boundary in your graph. If not, explain why. Note: You do not have to calculate the weights, just draw the decision boundary.

Yes, they are linearly separable. The decision boundary is shown below.

<matplotlib.collections.PathCollection at 0x7f0ce29a8cc0>



vi) 1 point) You train a hard-margin SVM on the higher-dimensional data using a library and it gives you the following Lagrange multiplier for your data1 : $\alpha_2 = 0.1$, $\alpha_3 = 0.2$, $\alpha_6 = 0.1$. What are the remaining Lagrange multipliers, α_1 , α_4 , α_5 and α_7 ? Justify your answer in 1-2 sentences. (Hint: This should not require any math to calculate.)

According to the graph above X2, X3 and X6 are support vectors as they are on the hard margin of the decision boundary. Since they are the support vectors they have a lagrange multiplier value $\alpha_i > 0$ and for the rest of the points the lagrange multiplier value is $\alpha=0$ as they are not the support vectors. So $\alpha_1=0$, $\alpha_4=0$, $\alpha_5=0$ and $\alpha_7=0$.

vii)(1 point) Recall that the SVM's prediction (using the Kernel transformation) for a data point Z can be defined with the following equation: $f(Z) = \text{sign}(\sum_i \alpha_i y_i (\phi(X_i) \cdot \phi(Z)) + w_0)$. You are given $w_0 = 2$. You are now asked to classify a new test data point, Z, using the SVM defined earlier by the Lagrange

multipliers. You do not know what Z's attributes are, but you do know: $K(X_2, Z) = 36$, $K(X_3, Z) = 12$, $K(X_6, Z) = 9$. Classify Z using the SVM. (Hint: If you find yourself trying to solve for Z's x value, you are doing it wrong.)

(vii) Given equation

$$f(z) = \text{sign}\left(\sum_i \alpha_i y_i (\phi(x_i) \cdot \phi(z)) + w_0\right)$$

$$w_0 = -2$$

$$K(x_2, z) = 36, K(x_3, z) = 12 \text{ and} \\ K(x_6, z) = 9$$

The above equation can be re-written as

$$f(z) = \text{sign}\left(\sum_i \alpha_i y_i K(x_i, z) + w_0\right)$$

On expanding the equation:

$$\text{sign}(\alpha_1 \times y_1 \times K(x_1, z) + \alpha_2 \times y_2 \times K(x_2, z) \\ + \alpha_3 \times y_3 \times K(x_3, z) + \alpha_4 \times y_4 \times K(x_4, z) \\ + \alpha_5 \times y_5 \times K(x_5, z) + \alpha_6 \times y_6 \times K(x_6, z) \\ + \alpha_7 \times y_7 \times K(x_7, z))$$

from above we know that $\alpha_1, \alpha_4, \alpha_5$ and α_7 are '0'. Since ignoring those terms & substituting remaining

$$\Rightarrow \text{sign}((0.1 \times (-1) \times 36 + 0.2 \times 1 \times 12 + 0.1 \times (-1) \times 9) + (-2))$$

$$\Rightarrow \text{Sign}(-4.1)$$

\Rightarrow $-$ (negative class).

z is classified as negative

CS Scanned with CamScanner

2 K-Means Clustering

1. (8 points) After each iteration of k-means, report the coordinates of the new centroids and which cluster each data point belongs to. Stop when the algorithm converges and clearly label on the graph where the algorithm converges. To report your work, give your answer in tabular format with the following attributes: Round (e.g. Round 1, 2, etc), Points (e.g. A, B, C), and Cluster ID (order does not matter). Also report the centroids for each cluster after each round. Please follow the example table format in Table 3.

ANS

(Q) K-Means Clustering

(a) Given Initial seed centroids are at points C, A, E.

Assigning each data points B, D, F to nearest centroid.

Round 1:-Point B \Rightarrow forms cluster with centroid "E" since this is nearest comparatively with A, C.So, $\{B, E\}$ forms 1 cluster.

$$\text{Calculating centroid for } B, E = B(1, 4), E(2, 1) \\ \left(\frac{1+2}{2}, \frac{4+1}{2}\right) = \underline{(1.5, 2.5)}$$

Points D, F \Rightarrow forms cluster with centroid "C"

$$\text{Centroid: } C(7, 4) \quad D(1, 9) \quad F(5, 6)$$

$$\left(\frac{7+1+5}{3}, \frac{4+9+6}{3}\right) \Rightarrow \left(\frac{13}{3}, \frac{19}{3}\right) = (4.33, 6.3)$$

?

Round	Points	cluster-ID	Centroids
1	$\{A\}$	1	$\{9, 5\}$
2	$\{E, B\}$	2	$\{1.5, 2.5\}$
3	$\{C, D, F\}$	3	$\{4.33, 6.33\}$

Round 2 :-Points $\{C, A\}$ form cluster because centroid $(9, 5)$ is nearest to point C.

Now centroids changes for clusters 1 & 3.

$$(A, C) \Rightarrow (9, 5) \text{ and } (7, 4) \Rightarrow \left(\frac{16}{2}, \frac{9}{2}\right) \Rightarrow (8, 4.5)$$

$$(D, F) \Rightarrow (1, 9) \text{ and } (5, 6) \Rightarrow \left(\frac{6}{2}, \frac{15}{2}\right) \Rightarrow (3, 7.5)$$

<u>Round</u>	<u>Points</u>	<u>cluster-ID</u>	<u>centroids</u>
2	{A, C}	1	[8, 4.5]
2	{B, E}	2	[1.5, 2.5]
2	{D, F}	3	[3, 7.5]

Round 3 :- ~~All~~ clusters remain same

<u>Round</u>	<u>Points</u>	<u>cluster-ID</u>	<u>centroids</u>
3	{A, C}	1	[8, 4.5]
3	{B, E}	2	[1.5, 2.5]
3	{D, F}	3	[3, 7.5]

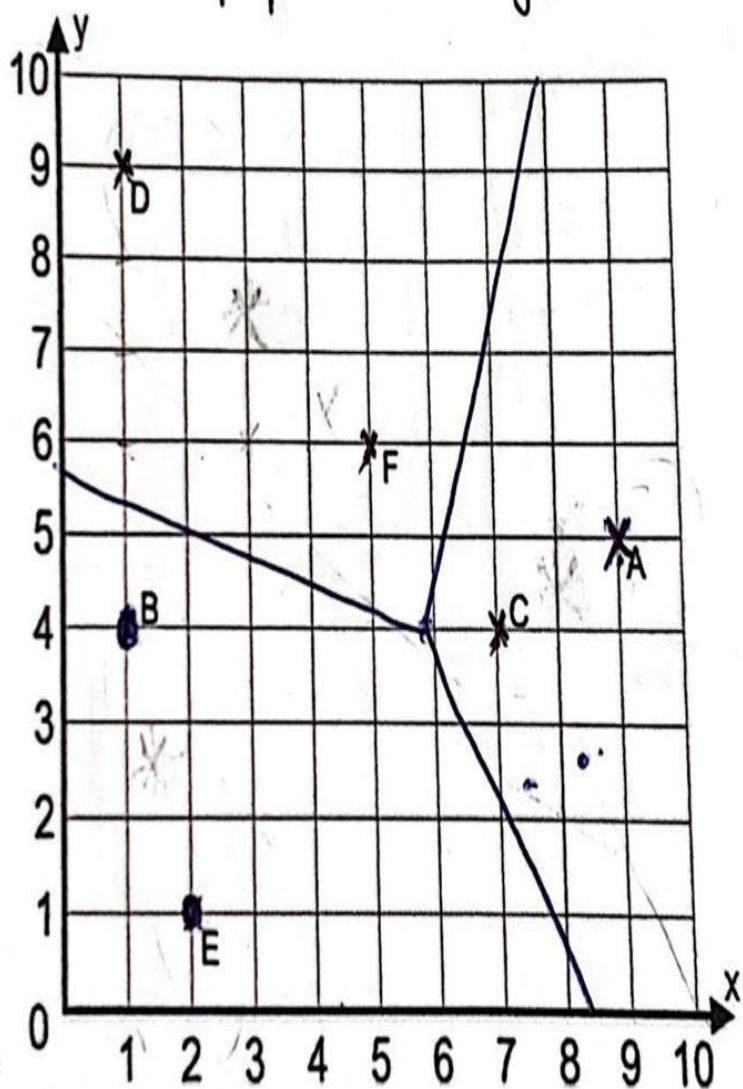
On Round 3, All Data Points are in the same cluster since the centroids are not changing. This states algorithm converged.

Overall Table.

<u>Round</u>	<u>Points</u>	<u>cluster-ID</u>	<u>centroids</u>
1	A	1	[9, 5]
	E, B	2	[1.5, 2.5]
	C, D, F	3	[4.33, 6.33]
2	A, C	1	[8, 4.5]
	B, E	2	[1.5, 2.5]
	D, F	3	[3, 7.5]
3	A, C	1	[8, 4.5]
	B, E	2	[1.5, 2.5]
	D, F	3	[3, 7.5]

data format in Table 3

Graph when algorithm converged.



Scanned with CamScanner

Figure 3: K-means Clustering (a)

2. 2b) (2 points) How many rounds are needed for the K-means clustering algorithm to converge?
ANS

(Qb) 3 rounds are needed for the k-means algorithm to converge. We can conclude this with the last iteration. In the last iteration of none of the points changes its cluster. As we saw in round 3 no point changed its cluster so we can terminate the algorithm.

2c)(4 points) Calculate the Sum of Squared Errors (SSE) for the k-means clustering.
ANS

(Qc)

From Round 3.

Cluster 1 : A, C and centroid (C1)

Distance (A, c1) and Distance (C, c1)

$$(A, c1) = \sqrt{(9-8)^2 + (5-4.5)^2} = \sqrt{1+0.25} = \sqrt{1.25}$$

$$(C, c1) = \sqrt{(7-8)^2 + (4-4.5)^2} = \sqrt{1+0.25} = \sqrt{1.25}$$

$$\text{Sum of squared Error [SSE1]} = (\sqrt{1.25})^2 + (\sqrt{1.25})^2 \Rightarrow \underline{\underline{2.50}}$$

Cluster 2 : B, E and centroid (C2)

Distance B/w (B, c2) and (E, c2)

$$(B, c2) = \sqrt{(1-1.5)^2 + (4-2.5)^2} = \sqrt{2.5}$$

$$(E, c2) = \sqrt{(2-1.5)^2 + (1-2.5)^2} = \sqrt{2.5}$$

$$\text{Sum of squared Error SSE2} = (\sqrt{2.5})^2 + (2.5)^2 \Rightarrow 5$$

Cluster 3 : (D, c3) & (F, c3)

$$(D, c3) = \sqrt{(3-1)^2 + (9-7.5)^2} = \sqrt{6.25}$$

$$(F, c3) = \sqrt{(5-1)^2 + (6-7.5)^2} = \sqrt{14+2.25} = \sqrt{16.25}$$

$$\text{SSE3} = (\sqrt{6.25})^2 + (\sqrt{16.25})^2 = 12.50$$

$$\boxed{\text{SSE} \Rightarrow \text{SSE1} + \text{SSE2} + \text{SSE3} \Rightarrow 2.50 + 5 + 12.50 \Rightarrow \underline{\underline{20}}}$$

$\frac{6.25}{16.25}$

3 Hierarchical Clustering

1. (4 points) Perform single link hierarchical clustering. Show your work at each iteration by giving the inter-cluster distances. Report your results by drawing a corresponding dendrogram. The dendrogram should clearly show the order and the height in which the clusters are merged. If possible, use a program to construct your dendrogram (e.g. PowerPoint, LucidChart2 , or VisualParadigm3). Scanned hand drawings will also be accepted if they are very clear. NOTE: There may be ties (i.e. two clusters have the same distance). In this case, you can choose any order to merge in, and ensure that this is reflected in your dendrogram.

ANS

(3) Hierarchical clustering

Given Distance Matrix

	B	D	E	F	
A	8.06	2.24	8.94	8.06	4.02
B	0	6	5	3.16	4.47
C	6	0	7.81	5.83	2.83
D	5	7.81	0	8.06	5
E	3.16	5.83	8.06	0	5.83

(3a) Single Link hierarchical Clustering

Iteration - 1 :-

Choosing minimum value from matrix (B-E) 2.24

So points A-C can be clustered. Now Euclidean matrix becomes

→ need to update matrix using minimum value.

	B	D	E	F
A	6	7.81	5.83	2.83
B	0	5	3.16	4.47
D	5	0	8.06	5
E	3.16	8.06	0	5.83

Eq'n :- A + C = 2.24

Again choosing minimum from matrix 2.83 and updating matrix with minimum value

	B	D	E
A C F	4.47	5	5.83
B	0	5	3.16
D	5	0	8.06

Eq'n's:

$$A + C = 2.24$$

$$AC + F = 2.83$$

Iteration 3 :- choosing min value: 2.16 & updating matrix with min value.

	BE	D
ACF	(4.47)	5
BE	0	5

Eq'n:

$$A + C = 2.24$$

$$AC + F = 2.83$$

$$B + E = 2.16$$

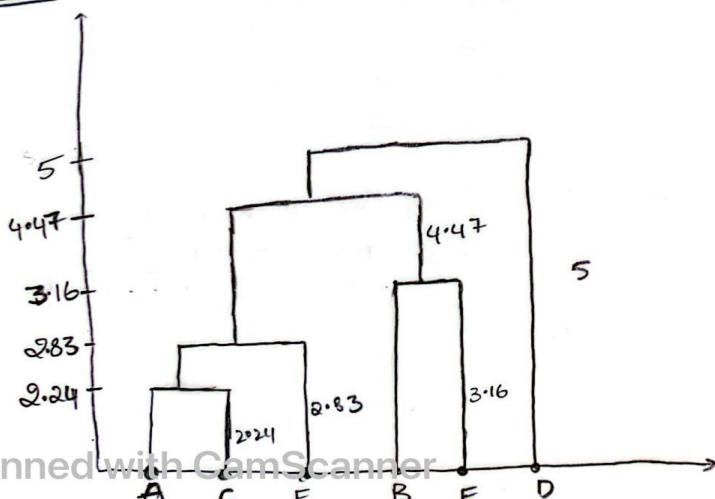
Iteration 4 :-

ACFB	E	D
	(5)	

Eq'n

$A + C = 2.24$ $A C + F = 2.83$ $B + E = 2.16$ $AC + F + BE = 4.47$ AC + F + BE $ACFB + D = 5$

Dendrogram



- 3b) (4 points) Perform complete link hierarchical clustering on the dataset. As above, show your calculations and report the corresponding dendrogram.

(3b) Complete link hierarchical clustering

Given Euclidean matrix

	B	C	D	E	F
A	8.06	(2.24)	8.94	8.06	4.012
B	0	6	5	3.16	4.047
C	6	0	7.81	5.83	2.83
D	5	7.81	0	8.06	5
E	3.16	5.83	8.06	0	5.83

Iteration - 1In Complete link clustering, choose minimum value cluster two points.For updating matrix consider max value minvalue is 2.24 Points A,C can be clustered.

	B	D	E	F
A,C	8.06	8.94	8.06	4.012
B	0	5	(3.16)	4.047
D	5	0	8.06	5
E	3.16	8.06	0	5.83

Eqn : A+C : 2.24

Iteration 2: choose min value (i.e) point which is nearer = 3.16 so B,E can be clustered & update the matrix with max values

	BE	D	F
AC	8.06	8.94	(4.12)
BE	0	8.06	5.83
D	8.06	0	5

Eq'n :-

$A + C = 2.24$
$B + E = 3.16$

Iteration :- 3 :-

	BE	D
ACF	(6.06)	8.94
BE	0	8.06

Eq'n :-

$A + C = 2.24$
 $B + E = 3.16$
 $AC + F = 4.12$

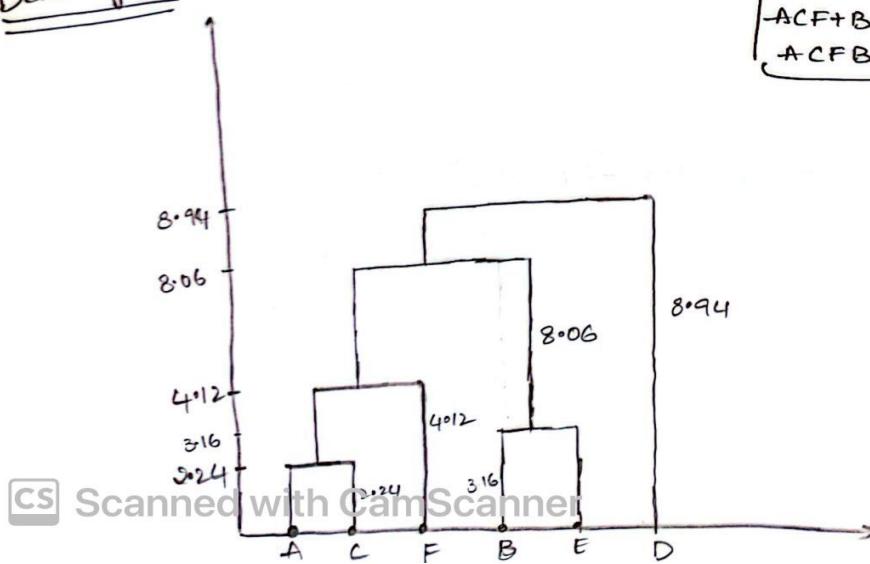
Iteration :- 4 :-

~~ABC~~ ACFBE D [8.94]

Eq'n :-

$A + C = 2.24$
$B + E = 3.16$
$AC + F = 4.12$
$ACF + BE = 8.06$
$ACFBE + D = 8.94$

Dendrogram



- 3c) (2 points) If we assume there are two clusters, will the single or complete link approach give a better clustering? Justify your answer. (Hint: you should not need to calculate SSE to answer the question.)
ANS

In case of only two clusters we think the complete link approach gives a better clustering. Because complete link clustering forms equal number of clusters in each round of iteration. Whereas in single link clustering if one item is cluster in one cluster and all the remaining ones go to the others. There is no uniform split in single link clustering.

4 Association Rule Mining

- 4a) (1 point) What is the maximum number of unique itemsets that can be extracted from this data set (including itemsets that have zero support)? Briefly explain your answer in 1-2 sentences

ANS

(4) Association Rule Mining

(4a) Ans: 63

There are 6 items. So we can form 2^6 unique subsets. So, 2^6 includes an empty transaction.

If we discard that, we get $2^6 - 1 \Rightarrow$ 63 unique subsets

4b) (1 point) What is the maximum number of association rules that can be extracted from this data set (including rules that have zero support)? Briefly explain your answer in 2-3 sentences

ANS

Ans: 602

(4b) Maximum number of Association Rules that can be Extracted from this Dataset $R = 3^d - 2^{d+1} - 1$ [where $d \Rightarrow$ itemset] $d=6$ here.

$$= 3^6 - 2^7 - 1$$

R = 602 \Rightarrow Max number of Association Rules.

4c) (1 point) Compute the support of the itemset: Chicken, Cookie?

ANS

$$(4c) \text{ Support of stemset } \{\text{chicken, cookie}\} = \frac{4}{10}$$

4d) (2 points) Compute the support and confidence of association rule: Cheese \rightarrow Butter?

ANS

$$(4d) \text{ Support of stemset } (\{\text{cheese}\} \rightarrow \{\text{Butter}\}) = \frac{\text{Support } \{\text{cheese} \rightarrow \text{Butter}\}}{11} = \frac{2}{10}$$

$$\text{Confidence of stemset } (\{\text{cheese}\} \rightarrow \{\text{Butter}\}) = \frac{\text{Support } \{\text{cheese}, \text{Butter}\}}{\text{Support } \{\text{cheese}\}} = \frac{2}{4}$$

4e) (3 points) Given min support = 0.2 and min confidence = 0.8, identify all valid association rules of the form A, B \rightarrow C.

ANS

(4e) Association Rules that can be formed & which are of the form $\{A, B\} \rightarrow \{C\}$

$\{\text{cheese, milk, Butter, chicken, cookie}\}$

$\{\text{cheese, milk}\} \rightarrow \text{Butter}$

$\boxed{\{\text{cheese, milk}\} \rightarrow \text{chicken}}$ ✓

$\{\text{cheese, milk}\} \rightarrow \text{cookie}$

$\{\text{cheese, Butter}\} \rightarrow \text{milk}$

$\{\text{cheese, Butter}\} \rightarrow \text{chicken}$

$\{\text{cheese, Butter}\} \rightarrow \text{cookie}$

$\{\text{cheese, chicken}\} \rightarrow \text{milk}$

$\{\text{cheese, chicken}\} \rightarrow \text{Butter}$

$\{\text{cheese, chicken}\} \rightarrow \text{cookie}$

$\{\text{cheese, chicken, cookie}\} \rightarrow \text{milk}$

$\{\text{cheese, cookie}\} \rightarrow \text{Butter}$

$\{\text{cheese, cookie}\} \rightarrow \text{chicken}$

$\{\text{milk, butter}\} \rightarrow \text{chicken}$

$\boxed{\{\text{milk, butter}\} \rightarrow \text{cookie}}$ ✓

$\{\text{milk, butter}\} \rightarrow \text{cheese}$

$\boxed{\{\text{milk, butter}\} \rightarrow \text{Butter}}$

$\{\text{milk, chicken}\} \rightarrow \text{cheese}$

$\{\text{milk, chicken}\} \rightarrow \text{cookie}$

$\{\text{milk, chicken, cookie}\} \rightarrow \text{cheese}$

$\{\text{milk, cookie}\} \rightarrow \text{cheese}$

$\boxed{\{\text{milk, cookie}\} \rightarrow \text{Butter}}$ ✓

$\{\text{milk, cookie}\} \rightarrow \text{chicken}$

$\{\text{butte, chicken}\} \rightarrow \text{cheese}$

$\{\text{butte, chicken}\} \rightarrow \text{cookie}$

$\{\text{butte, chicken, cookie}\} \rightarrow \text{cheese}$

$\{\text{chicken, cookie}\} \rightarrow \text{cheese}$

$\{\text{chicken, cookie}\} \rightarrow \text{milk}$

$\{\text{chicken, cookie}\} \rightarrow \text{Butter}$

$\{\text{butte, cookie}\} \rightarrow \text{cheese}$

$\boxed{\{\text{butte, cookie}\} \rightarrow \text{milk}}$

∴ only the rules

$\{\text{cheese, milk}\} \rightarrow \text{chicken}$

$\{\text{milk, butter}\} \rightarrow \text{cookie}$

$\{\text{milk, cookie}\} \rightarrow \text{Butter}$

$\{\text{Butter, cookie}\} \rightarrow \text{milk}$

These rules satisfies the form $\{A, B\} \rightarrow \{C\}$ which has $\text{minsup} = 0.2$

and $\text{confidence} \geq \text{minconf} = 0.8$

CS Scanned with CamScanner

4f) (2 points) In a different dataset, the support of the rule $a \rightarrow b$ is 0.62, and the support of the rule $a, c \rightarrow b, d$ is 0.31. What can we say for sure about the support of the rule $a \rightarrow b, d$. Explain in 1-2 sentences.

ANS

(4f) Given,

$$\text{Support of rule } \{a\} \rightarrow \{b\} = 0.62$$

$$\text{Support of rule } \{a,c\} \rightarrow \{b,d\} = 0.31$$

find support of rule $\{a\} \rightarrow \{b,d\}$:-

From the Antimonotone Property:-

Property states that support of an itemset never exceeds the support of its subsets.

$$(a \rightarrow \{b\}) \subseteq \{a\} \rightarrow \{b,d\} \stackrel{C}{\leq} \{a,c\} \rightarrow \{b,d\}$$

Therefore:-

$$S(\{a\} \rightarrow \{b\}) > S(\{a\} \rightarrow \{b,d\}) \geq S(\{a,c\} \rightarrow \{b,d\})$$

$$0.62 \geq S(\{a\} \rightarrow \{b,d\}) \geq S(\{a,c\} \rightarrow \{b,d\})$$

From this we can further say that

$$[S(\{a\} \rightarrow \{b,d\}) \text{ lies b/w } 0.62 \text{ & } 0.31]$$



Scanned with CamScanner

5 Apriori algorithm

1. 5a) (5 points) Show (compute) each step of frequent itemset generation process using the apriori algorithm, with a minimum support count of 3.

ANS

(5) Apriori Algorithm

dataset given

TID	Items
t1	A, C, D
t2	A, B, C, D
t3	A, C
t4	A, B, C
t5	B, C
t6	A, D
t7	A, B, D
t8	A, B

One itemsets

Itemset	Sup.
{A}	7
{B}	5
{C}	5
{D}	4

Given minsupcount=3
All has supcount >= 3
So all items need
to be considered

Possible two itemsets

$\{A, B\}$
 $\{A, C\}$
 $\{A, D\}$
 $\{B, C\}$
 $\{B, D\}$
 $\{C, D\}$

two itemsets

Itemset	Sup.
{A, B}	4
{A, C}	4
{A, D}	4
{B, C}	3
{B, D}	2
{C, D}	2

minsup = 3 So
 $\{B, D\}, \{C, D\}$ has minsup=2
 So need not consider these
 itemsets

atlast there is no itemset left which has supcount \geq minsupcount

CS Scanned with CamScanner

- 5b) (5 points) Show the lattice structure for the data given in table above, and mark each node in the lattice as either F: Frequent, IC: Infrequent due insufficient support count, or IP: Infrequent due to pruning (we do not need to calculate the support count). (Scanned hand-drawing is acceptable as long as it is clear. Tip (optional): Applications like CamScanner and Adobe Scan improve the quality of photos and reduce their size.)

ANS

Color green represents = F: Frequent

Color yellow represents = IC: Infrequent due insufficient support count

color pink represents = IP: Infrequent due to pruning

