

Course Code: CS481	Course Name: Data Science
Instructor Name: Dr. Muhammad Atif Tahir and Zeshan Khan	
Student Roll No: Solution	Section No: Solution

#### Instructions:

Start of Exam: 9:00 am; End of Exam: 12:30 pm including submission time  
Read each question completely before answering it. There is **7 questions and 6 pages**.

In case of any ambiguity, you may make assumptions. But your assumption should not contradict any statement in the question paper.

You will attempt this paper **offline**, in your **hand writing**.

You may use **cam-scanner**, **MS lens** or any equivalent application to scan and convert your hand-written answer sheets in a **single PDF file**.

The paper should be submitted using Google Form (link at the end of the paper). You are given 30 minutes for this purpose, which is already included in the exam time mentioned above. Additionally, after submitting, you should email it to your instructor which should be exactly same pdf as uploaded earlier.

**WRITE YOUR ID ON TOP OF EVERY PAGE by your hand. Write also page # on every page. You should also sign on every page.**

Please fill the below table with your details. A sample value for a male student having roll number K16-3689 and name Zeshan Khan Alvi is provided.

Sr#	Key	Description	Sample Value	Value for you
1	@fullname	Your Full Name	Zeshan Khan Alvi	Zeshan Khan Alvi
2	@fname	Your First Name	Zeshan	Zeshan
3	@lname	Your Last Name	Alvi	Alvi
4	@gender	Your Gender [male,female]	Male	Male
5	@nameparts	Number of words/parts in your full name	3	3
6	@serial	The last 4 digits of your roll number	3689 @serial[0]=3 @serial[1]=6 @serial[2]=8 @serial[3]=9	3501 @serial[0]=3 @serial[1]=5 @serial[2]=0 @serial[3]=1

**Time:** 180 minutes

**Max Marks:** 50 Points

**Question 1 [1 (Tokenize)+ 4(VSM)+ 2(ranking) = 7 Points]:** Given paragraphs as documents and your @fullname as a query string. Tokenize the words and take first letter of each word in lower case as a token. Apply word vector space model (VSM), to compute the similarity between query vector and document vector and return the ranked list of documents for the provided query. For the ease of computation, you can

use  $\text{similarity}(q, d) = \sum_{t=1}^{\text{tokens}} (q_t * d_t)$  for similarity between document vector d and query vector q where  $d_t$  is the counting of  $t^{th}$  token in document d.

Document 1) Shahzaib Yousuf Bilal Hyder Saad Umar Imtiaz Ali Issam Ahmed Neha  
 Document 2) Nadeem Hassan Afzal Abdullah Mujeeb Doulat Singh Murtaza Ali  
 Document 3) Dawar Hasnain Hamza Ashfaq Aliakber Madni Hussain Ashar Ali  
 Document 4) Ramchand Muhammad Ushay Murtaza Fakhruddin Ammar Rizwan  
 Document 5) Mujtaba Usama Vasnani Mehdi Raza Subash Kumar Shumail Steve

Note: Tokens for a document/query “Zeshan Khan Alvi” are [z, k, a].

## Solution

### Tokenize

Doc1	s	y	b	h	s	u	i	a	i	a	n
Doc2	n	h	a	a	m	d	s	m	a		
Doc3	d	h	h	a	a	m	h	a	a		
Doc4	r	m	u	m	f	a	r				
Doc5	m	u	v	m	r	s	k	s	s		
Query	z	k	a								

### Vectors

Doc/Tokens	a	b	h	i	n	s	y	m	n	d	r	u	f	r	v	k	z
Doc1	1	1	1	1	1	2	1										
Doc2	3		1			1		2	1	1							
Doc3	4		3					1		1							
Doc4	1							2			1	1	1	1			
Doc5						3		2				1		1	1	1	
Query	1															1	1

### Vector Space Model (VSM)

Similarity(Doc1, Query) =  $0 * 1 + 0 * 1 + 1 * 1 = 1$

Similarity(Doc2, Query) =  $0 * 1 + 0 * 1 + 3 * 1 = 3$

Similarity(Doc3, Query) =  $0 * 1 + 0 * 1 + 4 * 1 = 4$

Similarity(Doc4, Query) =  $0 * 1 + 1 * 1 + 0 * 1 = 1$

Similarity(Doc5, Query) =  $0 * 1 + 1 * 1 + 0 * 1 = 1$

### Ranking

Doc3, Doc2, Doc1, Doc4, Doc5

### Question 2 [1 (Data Completion) + 4 (Training)+ 2 (Testing)= 7 Points]:

Apply the naive bayes classifier on the data provided below. There are five columns (3 for X and one for the output Label and Sr# is only a serial number) of the data. The data-set provided in “Table 1” have 10 training samples and two testing samples. You are required to provide the labels for the test samples (Sr# 11 and Sr# 12) in the test data for classification.

Note: “**if(@fname[0]==vowel)**” is true if your first name starts from a vowel e.g. Owais, Ali, Imran etc.

Table 1: Training data for classification

Sr#	gender	is_even (@serial)	Name starts at Vowel	Grade (SU/Letter)
1	male	Yes	yes	SU
2	male	No	if(@fname[0]==vowel)	SU
3	@gender	Yes	yes	SU
4	male	No	if(@lname[0]==vowel)	Letter
5	male	Yes	no	Letter
6	@gender	No	if(@fname[0]==vowel)	Letter
7	female	Yes	if(@lname[0]==vowel)	SU
8	male	@serial%2==0	yes	SU
9	male	Yes	no	SU
10	female	@serial%2==0	yes	SU

Table 2: Test data for Classification

Sr#	Gender	is_even (@serial)	Name starts at Vowel	Grade (SU/Letter)
11	@gender	@serial%2==0	if(@fname[0]==vowel)	?
12	@gender	@serial%2==0	if(@lname[0]==Consonant)	?

## Solution

### Data Completion

Table 1: Training data for classification

Sr#	gender	is_even (@serial)	Name starts at Vowel	Grade (SU/Letter)
1	male	Yes	yes	SU
2	male	No	no	SU
3	male	Yes	yes	SU
4	male	No	yes	Letter
5	male	Yes	no	Letter
6	male	No	no	Letter
7	female	Yes	yes	SU
8	male	No	yes	SU
9	male	Yes	no	SU
10	female	No	yes	SU

Table 2: Test data for Classification

Sr#	Gender	is_even (@serial)	Name starts at Vowel	Grade (SU/Letter)
11	male	No	no	?
12	male	No	no	?

## Training

Gender(male, SU) = 5/7

Gender(female, SU) = 2/7

Gender(male, Letter) = 3/3

Gender(female, Letter) = 0/3

Is\_even(Yes, SU) = 4/7

Is\_even(No, SU) = 3/7

Is\_even(Yes, Letter) = 1/3

Is\_even(No, Letter) = 2/3

Vowel(Yes, SU) = 5/7

Vowel(No, SU) = 2/7

Vowel(Yes, Letter) = 1/3

Vowel(No, Letter) = 2/3

## Testing

$P(\text{male, no, no: SU}) = \text{Gender}(\text{male, SU}) * \text{Is\_even}(\text{No, SU}) * \text{Vowel}(\text{No, SU})$

$P(\text{male, no, no: SU}) = 5/7 * 3/7 * 2/7 = 30/343 = 0.08$

$P(\text{male, no, no: Letter}) = 3/3 * 2/3 * 2/3 = 12/27 = 0.44$

Letter Grade for both

**Question 3 [6 Points]:** Using hierarchical clustering algorithms (Single link and Distance b/w centroids) and City-block distance ( $d = (|x_2 - x_1|) + (|y_2 - y_1|)$ ) to cluster the following 5 points into 3 clusters. Using  $A1 = (@\text{serial}[1], 10)$ ,  $A2 = (2, @\text{serial}[0])$ ,  $A3 = (8, 4)$ ,  $A4 = (5, @\text{serial}[2])$ ,  $A5 = (7, 5)$ .

Point	X	Y
A1	5	10
A2	2	3
A3	8	4
A4	5	0
A5	7	5

## Centroid

X	y	Distance				
5	10	0	10	9	10	7
2	3	10	10	9	10	7
8	4	9	7	0	7	2
5	0	10	6	7	0	7
7	5	7	7	2	7	0

Point	X	Y	Distance			
3 5	1	5	10	0	10	8
	2	2	3	10	0	7
		7.5	4.5	8	7	0
	4	5	0	10	6	7

Point	X	Y	Distance			
2 3 4 5	1	5	10	0	8	
		4.5	2.5	8	0	

Single Link

	X	Y	Distance					
A1	5	10	0	10	9	10	7	
A2	2	3	10	0	7	6	7	
A3	8	4	9	7	0	7	2	
A4	5	0	10	6	7	0	7	
A5	7	5	7	7	2	7	0	

Point	Distance					
3 5	1	0	10	7	10	
	2		0	7	6	
				0	7	
	4				0	

Point	Distance			
2 4 3 5	1	0	10	7
			0	7
				0

**Question 4 [1.5+1.5+1.5+1.5=6 Points]:** Consider the data set shown in Table below:

$w = \text{ceil}(@\text{serial}[0]/2)$   
 $x = \text{ceil}(@\text{serial}[1]/2)$   
 $y = \text{ceil}(@\text{serial}[2]/2)$   
 $z = \text{ceil}(@\text{serial}[3]/2)$

Customer Number	Items Bought
A	{1, <b>x</b> , 5}
A	{1,2,3,5}
B	{1,2,4,5}
B	{1,3,4,5}
C	{2,3, <b>w</b> }
C	{2,4,5}
D	{3,4}
D	{1, <b>y</b> ,3}
E	{1,4,5}
E	{1,2, <b>z</b> }

- Compute the support for itemsets {5}, {2, 4}, and {2, 4, 5} by treating each transaction ID as automobile shop basket.
- Computer the confidence for the association rule (i) {2,4} -> {5} (ii) {5} -> {2,4}. Is the confidence a symmetric measure?
- Repeat (a) by treating each customer number as automobile parts basket. Similar customer numbers should be treated as one customer.
- What is the maximum number of association rules that can be extracted from this data?

### Solution

Customer Number	Items Bought
A	{1, 4, 5}
A	{1,2,3,5}
B	{1,2,4,5}
B	{1,3,4,5}
C	{2,3,5}
C	{2,4,5}
D	{3,4}
D	{1,2,3}
E	{1,4,5}
E	{1,2,5}

- A) Support ({5}) =  $8 / 10 = 0.8$
- Support ({2,4}) =  $2 / 10 = 0.2$
  - Support ({2,4,5}) =  $2 / 10 = 0.2$
- B) confidence(2,4 -> 5) =  $0.2 / 0.2 = 1$
- confidence(5->2,4) =  $0.2 / 0.8 = 0.25$
  - No, it is not symmetric measure
- C) Support(5) =  $4/5$

a)  $\text{Support}(2,4) = 5/5 = 1$

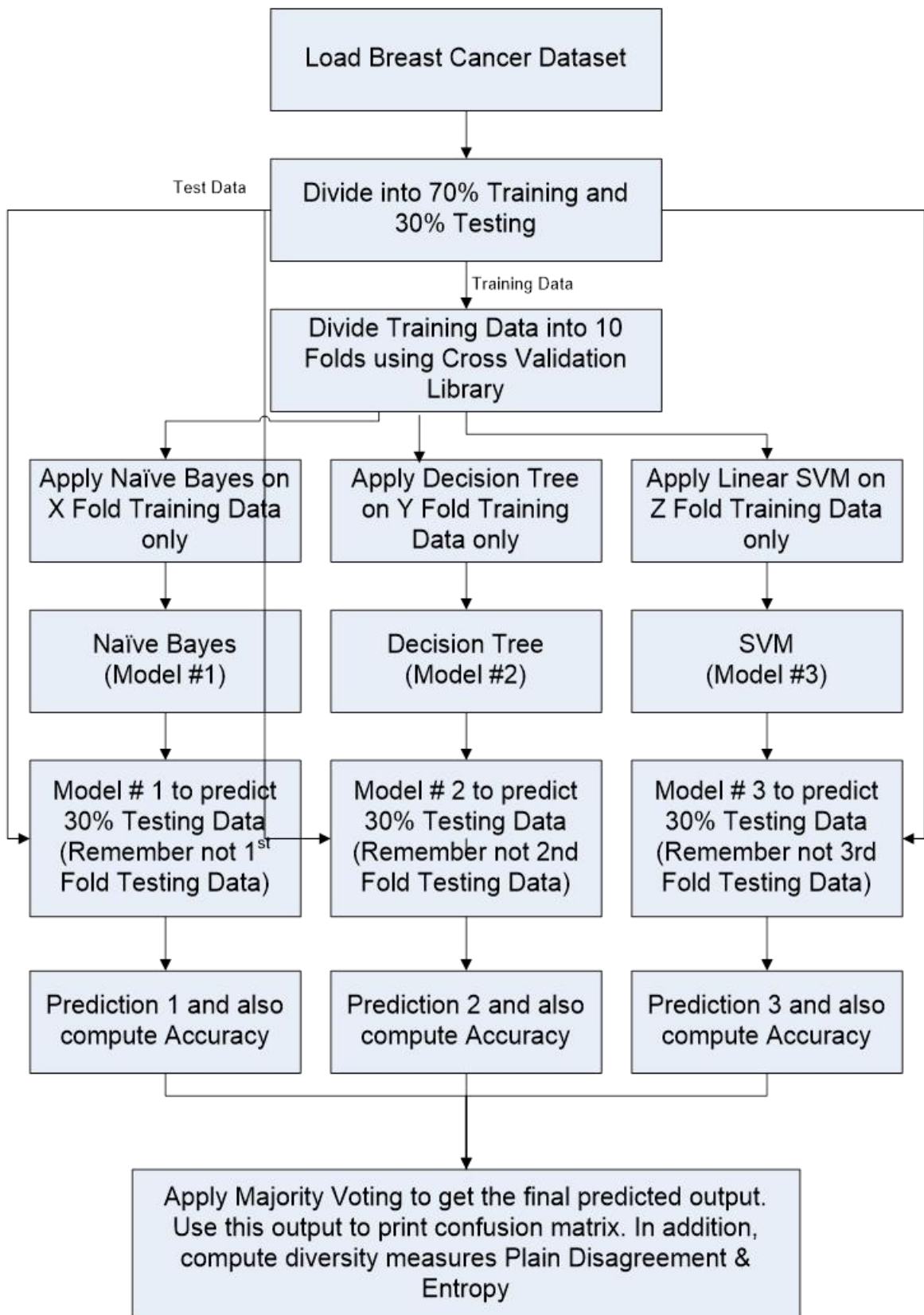
b)  $\text{Support}(2,4,5) = 4/5$

D) Sol:  $3^{10} - 2^{10+1} + 1 = 57002$

**Question 5 [10 Points]:** Implement the model shown in Figure below. Upload the source code only. Here  $X = @serial[1]+1$ ,  $Y = @serial[2]+1$ ,  $Z = @serial[3]+1$

Consider the following Initial Coding:

```
# Load libraries
..... // import necessary libraries
from sklearn.datasets import load_breast_cancer
breast_data = load_breast_cancer()
#..... Complete the program as explained in block diagram
```





**Question 6 [8 Marks]:** The following table shows the value of shares of company in Karachi Stock at the end of last four weeks:

Date	Share Value (Target Variable)
3 <sup>rd</sup> Sept 2017	@serial[0]
27 <sup>th</sup> Oct 2017	@serial[1]
20 <sup>th</sup> Nov 2017	@serial[2]
1 <sup>st</sup> Dec 2017	@serial[3]

The following two events in Table below are responsible for the change of shares of company

**Table: Events**

Date	Event1 in Million Rupees (New Investment)	Event2 in Million Rupees (Loan Return)
3 <sup>rd</sup> Sept 2017	3	4
27 <sup>th</sup> Oct 2017	4	3
20 <sup>th</sup> Nov 2017	2	1
1 <sup>st</sup> Dec 2017	1	2

What is the predicted share value on 1<sup>st</sup> Jan 2018 (show all steps with illustration) if following events are going to happen on 1<sup>st</sup> Jan 2018 [Hint: Use PCA to reduce the dimensions of 2 events to 1, then apply linear regression on 1st dimension as independent variable and share value as target variable]

Event1 in Million Rupees (New Investment)	Event2 in Million Rupees (Loan Return)
@serial[1]	@serial[2]

**Question 7 [6 Points]:**

**(a)[4 Points]** In this Problem, you will work on the **Error Bars** to display error visually in a bar chart. For someone who is learning about the different drink types at Macdonald, a bar chart of milk amounts in each drink may be useful. We have provided the `ounces_of_milk` list, which contains the amount of milk in each 14oz drink in the `drinks` list. According to different barista styles and measurement errors, there might be variation on how much milk actually goes into each drink. We have included a list `error` on each amount of milk. You need to write program in your answer sheet.

```
drinks = ["cappuccino", "latte", "chai", "americano", "mocha", "espresso"]
ounces_of_milk = [w, x, y, z, 9, 10]
error = [1.1, 0.7, 1.2, 1.0, 0, 1.7]
where, w = @serial[0], x = @serial[1], y = @serial[2], z = @serial[3]
```

- Plot this information as a bar chart. [1 point]
- Display this error as error bars on the bar graph and add caps of size 5 to your error bars. [1 point]
- Set the axis to go from 'cappuccino' to 'espresso' on the x-axis and 2 to 14 on the y-axis. [1 point]

iv. Add the title "Drinks to milk ratio", x-axis label "Drinks", and y-axis label "Milk amount in ounces. [1 point]

**(b)[2 Points]** Write in your own words difference between seaborn heatmap, seaborn stripplot, seaborn violin, and seaborn heatmap.

### Solution A

```
from matplotlib import pyplot as plt
```

```
drinks
```

```
=["cappuccino","latte","chai","americano","mocha","espresso"]
```

```
ounces_of_milk = [6, 9, 4, 3, 9, 10]
```

```
error = [1.1, 0.7, 1.2, 1.0, 0, 1.7]
```

```
ax = plt.subplot()
```

```
plt.bar(drinks,ounces_of_milk,  
yerr=error,capsize=10)
```

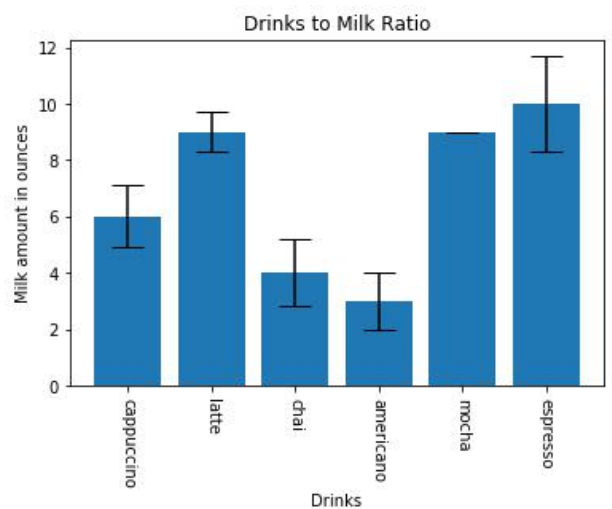
```
ax.set_xticklabels(drinks,  
rotation=270)
```

```
plt.title("Drinks to Milk Ratio")
```

```
ax.set_xlabel("Drinks")
```

```
ax.set_ylabel("Milk amount in ounces")
```

```
plt.show()
```



### Concluding Remarks

You need to prepare a pdf file of all the question as per the question ordering. The orientation should be portrait for each page. It should be clearly visible for each and every text written on the page. You suppose to upload it on the provided form as an assignment submission. You have good 30 minutes for it.

Form URL

[https://docs.google.com/forms/d/e/1FAIpQLSefZ3vTJHuudiEDISu3ok7t1kHXEDIkEkdTUHNxlbVYtby1gw/viewform?usp=sf\\_link](https://docs.google.com/forms/d/e/1FAIpQLSefZ3vTJHuudiEDISu3ok7t1kHXEDIkEkdTUHNxlbVYtby1gw/viewform?usp=sf_link)

***BEST OF LUCK!***