

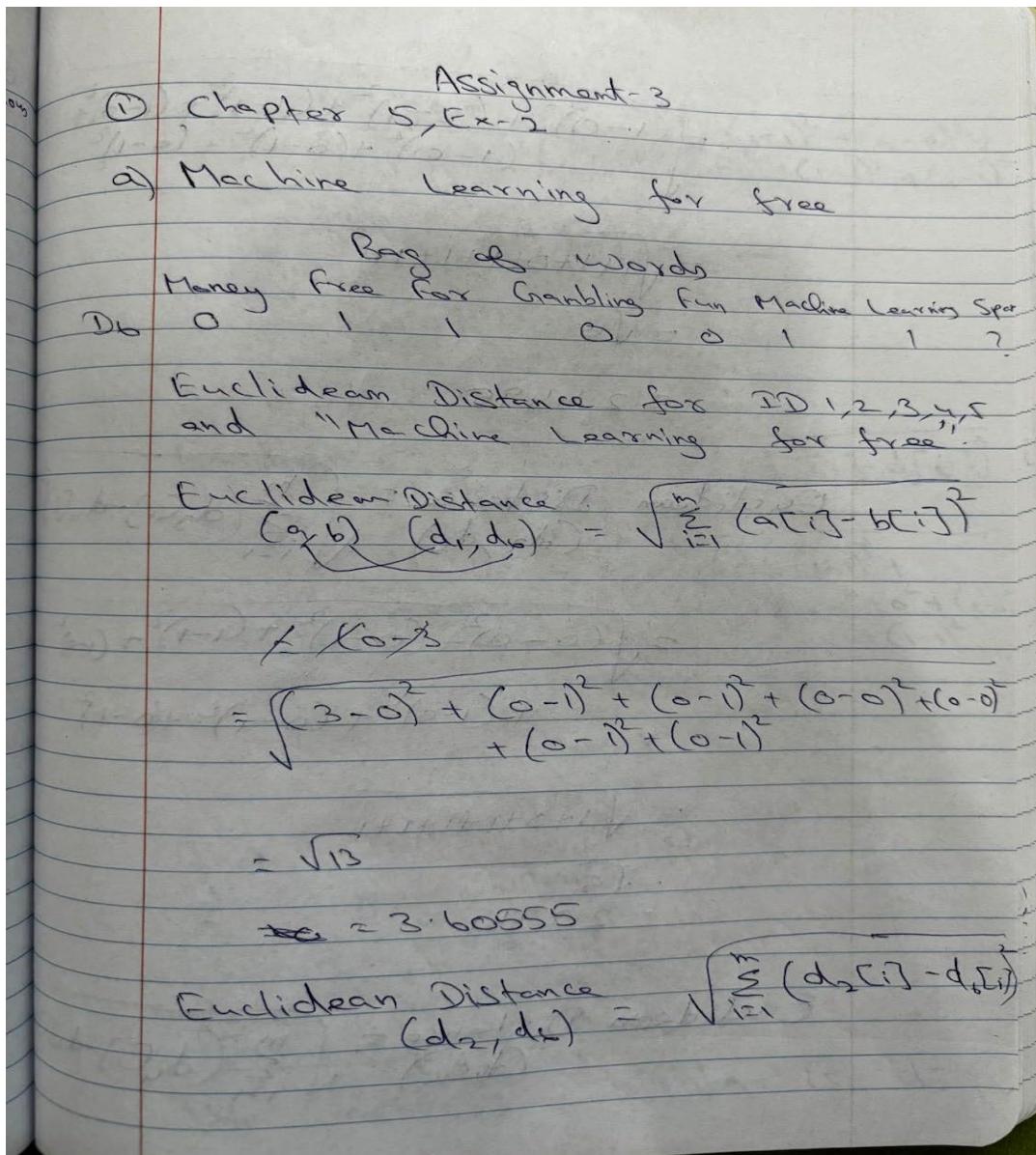
# Assignment-3-Aparna Bharathi Suresh

## 1.(15%) Chapter 5, exercise 2

Question:a)

What target level would a nearest neighbor model using Euclidean distance return for the following email: "machine learning for free"?

Solution:



$$\begin{aligned}
 &= \sqrt{(1-0)^2 + (2-1)^2 + (1-1)^2 + (1-0)^2 + \\
 &\quad (1-0)^2 + (0-1)^2 + (0-1)^2} \\
 &= \sqrt{1+1+0+1+1+1} \\
 &= \sqrt{6}
 \end{aligned}$$

$$= 2.4495$$

$$\text{Euclidean Distance} = \sqrt{\sum_{i=1}^m (d_3(i) - d_6(i))^2}$$

$$\begin{aligned}
 &= \sqrt{(0-0)^2 + (0-1)^2 + (1-1)^2 + (1-0)^2} \\
 &\quad + (0-0)^2 + (1-0)^2 + (1-0)^2 + (1-0)^2 \\
 &= \sqrt{1+0+1+1+1+1}
 \end{aligned}$$

$$= \sqrt{5}$$

$$= 2.2361$$

$$\text{Euclidean Distance} = \sqrt{\sum_{i=1}^m (d_4(i) - d_5(i))^2}$$

$$\begin{aligned}
 d_{5,6} &= \sqrt{(0-0)^2 + (0-1)^2 + (1-1)^2 + (0-1)^2 + (3-0)^2 + (1-1)^2 + (1-1)^2} \\
 &= \sqrt{1+0+0+0+0} \\
 &= \sqrt{5} \\
 &= 2.2361
 \end{aligned}$$

$$\begin{aligned}
 \text{Euclidean Distance} &= \sqrt{\sum_{i=1}^m (d_{5,i} - d_{6,i})^2} \\
 (d_5, d_6) &= \sqrt{(0-0)^2 + (1-1)^2 + (0-1)^2 + (0-0)^2 + (0-0)^2 + (1-1)^2} \\
 &= \sqrt{1+0+0+0+0+1} \\
 &= \sqrt{3.1623}
 \end{aligned}$$

ID	Euclidean Distance
1	3.6055
2	2.4495
3	2.2361
4	3.1623
5	1

The nearest neighbor is  $d_5$ .

The target value of  $d_5$  is false so the target value of "Machine Learning for free" is false.  $\underline{\text{SPAM} = \text{False}}$ .

### Question: b)

What target level would a k-NN model with  $k=3$  and using Euclidean distance return for the same query?

### Solution:

b) From the previous solution we see that the 3 nearest neighbors are  $d_2, d_3$  and  $d_5$ .

SPAM	
$d_2 \Rightarrow$	True
$d_3 \Rightarrow$	True
$d_5 \Rightarrow$	False

Majority is True, So the 3-NN model will return  $\text{SPAM} = \text{true}$ .

### Question: c)

What target level would a weighted k-NN model with  $k=5$  and using a weighting scheme of the reciprocal of the squared Euclidean distance between the neighbor and the query, return for the query?

### Solution:

c) Weighted K-NN,  $K=5$

ID	Euclidean Distance
1	3.6055
2	2.4495
3	2.2361
4	3.1623
5	1

$$\text{Weighted KNN} = \frac{1}{(\text{Euclidean Distance})^2}$$

ID of Weighted KNN

$$\begin{aligned} 1 & \quad \frac{1}{(3.6055)^2} = 0.0769 \\ 2 & \quad \frac{1}{(2.4495)^2} = 0.1667 \\ 3 & \quad \frac{1}{(2.2361)^2} = 0.2 \\ 4 & \quad \frac{1}{(3.1623)^2} = 0.1 \\ 5 & \quad \frac{1}{1^2} = 1 \end{aligned}$$

Total weight of True:  $d_1 + d_2 + d_3$

$$\begin{aligned} & 0.0769 + 0.1667 + 0.2 \\ & = 0.4436 \end{aligned}$$

$$\begin{aligned} \text{Weight of False} &= d_4 + d_5 \\ &= 0.1 + 1 \\ &= 1.1 \end{aligned}$$

$\text{Span} = \text{false}$  has the max. weight, so the model predicts  $\text{SPAN} = \text{false}$ .

## Question: d)

What target level would a k-NN model with k=3 and using Manhattan distance return for the same query?

Solution:

d) Manhattan Distance: K=3

$$\text{Manhattan Distance}(a, b) = \sum_{i=1}^n |\text{abs}(a[i] - b[i])|$$

Manhattan Distance

$$(v, d_1) = |\text{abs}(0-3)| + |\text{abs}(1-0)| + |\text{abs}(1-0)| + |\text{abs}(0-0)| + |\text{abs}(1-0)|$$

$$= 3 + 1 + 1 + 0 + 1 + 1$$

$$= 7$$

Manhattan Distance

$$(v, d_2) = |\text{abs}(0-1)| + |\text{abs}(1-2)| + |\text{abs}(1-1)| + |\text{abs}(0-1)| + |\text{abs}(0-1)| + |\text{abs}(1-0)| + |\text{abs}(1-0)|$$

$$= 1 + 1 + 0 + 1 + 1 + 1 + 1$$

$$= 6$$

Manhattan Distance

$$(v, d_3) = |\text{abs}(0-0)| + |\text{abs}(1-0)| + |\text{abs}(1-1)| + |\text{abs}(0-1)| + |\text{abs}(0-1)| + |\text{abs}(1-0)| + |\text{abs}(1-0)|$$

$$= 0 + 1 + 0 + 1 + 1 + 1 + 1$$

$$= 5$$

Manhattan Distance

$$(v, d_4) = |\text{abs}(0-0)| + |\text{abs}(1-0)| + |\text{abs}(1-1)| + |\text{abs}(0-0)| + |\text{abs}(0-3)| + |\text{abs}(1-1)| + |\text{abs}(1-1)|$$

$$= 0 + 1 + 0 + 0 + 3 + 0 + 0$$

~~= 4~~  
Manhattan Distance

$$(q, d_5) = |0-0| + |1-1| + |1-0| + |0-0| \\ + |0-0| + |1-1| + |1-1| \\ = 0+1+1+0+0+0 \\ = 1$$

$k=3$ ; 3 Nearest Neighbors  
are:  $d_5, d_4, d_3$

$d_5 \Rightarrow$  False  
 $d_4 \Rightarrow$  False  
 $d_3 \Rightarrow$  True

Majority False so 3NN  
model using Manhattan distance  
will return SPAM = False.

### Question: e)

There are a lot of zero entries in the spam bag-of-words dataset. This is indicative of sparse data and is typical for text analytics. Cosine similarity is often a good choice when dealing with sparse non-binary data. What target level would a 3-NN model using cosine similarity return for the query?

### Solution:

R. Cosine Similarity:

$$\text{Sim}_{\text{cosine}}(a, b) = \frac{(a[1] \times b[1]) + \dots + (a[n] \times b[n])}{\sqrt{\sum_{i=1}^n a[i]^2} \times \sqrt{\sum_{i=1}^n b[i]^2}}$$

$$\text{Sim}_{\text{cosine}}(d_1, q) = \frac{(d_1[1] \times q[1]) + \dots}{\sqrt{\sum_{i=1}^m d_1[i]^2} \times \sqrt{\sum_{i=1}^m q[i]^2}}$$

$$\sqrt{\sum_{i=1}^m d_1[i]^2} = \sqrt{3^2 + 0 + 0 + 0 + 0 + 0 + 0} \\ = \sqrt{9}$$

$$\sqrt{\sum_{i=1}^m q[i]^2} = \sqrt{0^2 + 1^2 + 1^2 + 0^2 + 0^2 + 1^2 + 1^2} \\ = \sqrt{4} \\ = 2$$

$$\text{Sim}_{\text{cosine}}(d_1, q) = \frac{(3 \times 0) + (0 \times 1) + (0 \times 1) + (0 \times 0) + (0 \times 0) + (0 \times 1) + (0 \times 1)}{3 \times 2}$$

$$= \frac{0}{6}$$

$$= 0$$

$$\text{Sim cosine } (d_2, g) = \frac{(d_2[1]g[1] + \dots)}{\sqrt{\sum_{i=1}^m d_2[i]^2} \times \sqrt{\sum_{i=1}^m g[i]^2}}$$

$$\sqrt{\sum_{i=1}^m d_2[i]^2} = \sqrt{1^2 + 2^2 + 1^2 + 1^2 + 1^2 + 0^2 + 0^2}$$

$$= \sqrt{8}$$

$$= 2.8284$$

$$\sqrt{\sum_{i=1}^m g[i]^2} = 2$$

$$\text{Sim cosine } (d_2, g) = \frac{(1 \times 0) + (2 \times 1) + (1 \times 1) + 0 + 0 + 0}{2.8284 \times 2}$$

$$= \frac{3}{5.6568}$$

$$= 0.5303$$

$$\text{Sim cosine } (d_3, g) = \frac{(d_3[1]g[1] + \dots)}{\sqrt{\sum_{i=1}^m d_3[i]^2} \times \sqrt{\sum_{i=1}^m g[i]^2}}$$

$$\sum_{i=1}^r d_y[i]^2 = \sqrt{0+0+1^2+1^2+1^2+0+0} \\ = \sqrt{3}$$

$$= 1.732$$

$$\text{Sim}_{\text{cosine}}(d_3, q_r) = \frac{(0)+(0 \times 1)+(1 \times 1)+0+0}{1.732 \times 2}$$

$$= \frac{1}{3.4641}$$

$$(0 \times 1)+(1 \times 1)+(0 \times 1) = (0, 1) \\ = 0.2887$$

$$\text{Sim}_{\text{cosine}}(d_4, q_r) = \frac{(d_y[i] \times q_r[i]) + \dots}{\sqrt{\sum_{i=1}^m d_y[i]^2} \times \sqrt{\sum_{i=1}^n q_r[i]^2}}$$

$$\sqrt{\sum_{i=1}^m d_y[i]^2} = \sqrt{0+0+1^2+0+3^2+1^2+1^2} \\ = \sqrt{12}$$

$$= 3.4641$$

$$\text{Sim cosine } (d_4, q) = \frac{0+0+1+0+0+1+1}{3 \cdot 4641 \times 2}$$

$$= \frac{3}{6.9282}$$

$$= 0.4330$$

$$\text{Sim cosine } (d_5, q) = \frac{(d_5[1] \times q[1]) + \dots}{\sqrt{\sum_{i=1}^m d_5[i]^2} \times \sqrt{\sum_{i=1}^n q[i]^2}}$$

$$\sqrt{\sum_{i=1}^m d_5[i]^2} = \sqrt{0+1+0+0+0+1+1}$$

$$= \sqrt{3}$$

$$= 1.7320$$

$$\text{Sim cosine } (d_5, q) = \frac{0+1+0+0+0+1+1}{1.7320 \times 2}$$

$$= \frac{3}{3.4641}$$

$$= 0.8660$$

In Similerry Index , higher the number , the more Similerry the instances.

## ID Cosine Similarity

1	0
2	0.5303
3	0.2887
4	0.4330
5	0.8660

So from the above  $d_5, d_2$  and  $d_4$  are the 3 most similar instances to the query.

$d_5 \Rightarrow \text{False}$

$d_2 \Rightarrow \text{True}$

$d_4 \Rightarrow \text{False}$

Majority False hence the prediction of SPAM = False.

## 2. (15%) Chapter 5, exercise 3

**Question:** a)

What value would a 3-nearest neighbor prediction model using Euclidean distance return for the CPI of Russia?

**Solution:**

② Chapter - 5 , Ex - 3

a) Euclidean Distance:

$$\text{Euclidean Distance} = \sqrt{\sum_{i=1}^m (a_{i,j} - b_{i,j})^2}$$

(Russia, Afghanistan)

$$= \sqrt{(67.62 - 59.61)^2 + (31.68 - 23.21)^2 + (10 - 74.30)^2 + (3.87 - 4.44)^2 + (12.90 - 0.40)^2}$$

= 66.5354

$$\text{Euclidean Distance} = \sqrt{\sum_{i=1}^m (a_{i,j} - b_{i,j})^2}$$

(Russia, Haiti)

$$= \sqrt{(67.62 - 45)^2 + (31.68 - 47.67)^2 + (10 - 73.10)^2 + (3.87 - 0.09)^2 + (12.90 - 3.40)^2}$$

= 69.6670

Euclidean Distance  
(Russia, Nigeria) =

$$\sqrt{(67.62 - 51.30)^2 + (31.68 - 38.23)^2 + (10 - 82.60)^2 + (3.87 - 1.07)^2 + (12.90 - 4.10)^2}$$
$$= 75.26808$$

Euclidean Distance  
(Russia, Egypt) =

$$\sqrt{(67.62 - 51.30)^2 + (31.68 - 26.58)^2 + (10 - 19.60)^2 + (3.87 - 1.86)^2 + (12.90 - 5.30)^2}$$
$$= 13.7168$$

Euclidean Distance =  
(Russia, Argentina)

$$\sqrt{(67.62 - 75.77)^2 + (31.68 - 32.30)^2 + (10 - 13.30)^2 + (3.87 - 0.76)^2 + (12.90 - 10.10)^2} = 9.7775$$

Euclidean Distance =  
(Russia, China)

$$\sqrt{(67.62 - 74.87)^2 + (31.68 - 29.99)^2 + (10 - 13.70)^2 + (3.87 - 1.95)^2 + (12.90 - 6.40)^2} = \cancel{10.72748}$$

Euclidean Distance =  
(Russia, Brazil)

$$\sqrt{(67.62 - 73.87)^2 + (31.68 - 42.93)^2 + (10 - 14.50)^2 + (3.87 - 1.43)^2 + (12.90 - 7.20)^2} = 14.6861$$

Euclidean Distance  
(Russia, Israel)

$$\sqrt{(67.62 - 81.30)^2 + (31.68 - 28.80)^2 + (10 - 3.60)^2 + (3.87 - 5.71)^2 + (12.90 - 12.50)^2} = 15.6514$$

Euclidean Distance  
(Russia, USA)

$$\sqrt{(67.62 - 78.51)^2 + (31.68 - 29.85)^2 + (10 - 6.30)^2 + (3.87 - 4.72)^2 + (12.90 - 13.70)^2} = 11.7044$$

Euclidean Distance  
(Russia, Ireland)

$$\sqrt{(67.62 - 80.15)^2 + (31.68 - 27.23)^2 + (10 - 3.50)^2 + (3.87 - 0.60)^2 + (12.90 - 11.50)^2} = 15.2219$$

Euclidean Distance  
(Russia, U.K)

$$\sqrt{(67.62 - 80.09)^2 + (31.68 - 28.49)^2 + (10 - 4.40)^2 + (3.87 - 2.59)^2 + (12.90 - 13.00)^2} = 14.0955$$

Euclidean Distance  
(Russia, Germany)

$$\sqrt{(67.62 - 80.24)^2 + (31.68 - 22.07)^2 + (10 - 3.50)^2 + (3.87 - 1.31)^2 + (12.90 - 12.00)^2} = 17.3559$$

Euclidean Distance  
(Russia, Canada)

$$\sqrt{(67.62 - 80.99)^2 + (31.68 - 24.79)^2 + (10 - 4.90)^2 + (3.87 - 1.42)^2 + (12.90 - 14.20)^2} = 16.1223$$

Euclidean Distance =  
(Russia, Australia)

$$\sqrt{(67.62 - 82.09)^2 + (31.68 - 25.40)^2 + (10 - 4.20)^2 + (3.87 - 1.86)^2 + (12.90 - 11.50)^2}$$
$$= 16.9841$$

Euclidean Distance =  
(Russia, Sweden)

$$\sqrt{(67.62 - 81.43)^2 + (31.68 - 22.18)^2 + (10 - 2.40)^2 + (3.87 - 1.27)^2 + (12.90 - 12.80)^2}$$
$$= 18.5261$$

Euclidean Distance =  
(Russia, New Zealand)

$$\sqrt{(67.62 - 80.67)^2 + (31.68 - 27.81)^2 + (10 - 4.90)^2 + (3.87 - 1.13)^2 + (12.90 - 12.30)^2}$$
$$= 14.804$$

## Euclidean Distance order:

ID	Euclidean	CPI
Argentina	9.7775	2.9961
China	10.7248	3.6356
USA	11.7044	7.1357
Egypt	13.7168	2.8622
UK	14.0955	7.7751
Brazil	14.6801	3.7741
New Zealand	14.804	9.4627
Ireland	15.2219	7.5360
Israel	15.6514	5.8069
Canada	16.1233	8.6725
Australia	16.9841	8.8442
Germany	17.3559	8.844
Sweden	18.5261	9.2985
Afghanistan	66.5354	1.5171
Haiti	69.6670	1.7999
Nigeria	75.26808	2.4493

The 3 NN are Argentina  
China + USA.

CPI value will be the  
average of the neighbors.  

$$\frac{2.9961 + 3.6356 + 7.1357}{3}$$

$\frac{2.9961 + 3.6356 + 7.1357}{3} = 4.5891$

CPI of Russia = 4.5891

### Question: b)

What value would a weighted k-NN prediction model return for the CPI of Russia? Use k\_16 (i.e., the full dataset) and a weighting scheme of the reciprocal of the squared Euclidean distance between the neighbor and the query.

### Solution:

b) Weighted KNN :

Weighted KNN (Argentina) =  $\frac{1}{(9.7715)^2} \times 2.9961$   
= 0.01045

Weight x CPI :

$$0.01045 \times 2.9961 = 0.0313$$

Weighted KNN (China) =  $\frac{1}{(10.7248)^2} \times 3.6356$   
= 0.0086

Similarly

ID	Euclidean	CPI	Weight	Weight x CPI
Argentina	9.7775	2.9961	0.01045	0.0313
China	10.7248	3.6356	0.0086	0.0316
USA	11.7044	7.1357	0.0072	0.0520
Egypt	13.7168	2.8622	0.0053	0.0152
UK	14.0955	7.7751	0.0050	0.0391
Brazil	14.6801	3.7741	0.0046	0.0172
New Zealand	14.804	9.4627	0.0045	0.043
Ireland	15.2219	7.5360	0.0043	0.0325
Israel	15.6514	5.8069	0.0040	0.0237
Canada	16.1223	8.6725	0.0038	0.0333
Australia	16.9841	8.8442	0.0035	0.0306
Germany	17.3559	8.0461	0.0033	0.0267
Sweden	18.5261	9.2985	0.0029	0.0270
Afghanistan	66.5354	1.5171	0.0002	0.0003
Haiti	69.6670	1.7999	0.0002	0.0004
Nigeria	75.26808	2.4493	0.0002	0.0004
			0.06805	0.4043

$$\frac{\sum \text{Weight} \times \text{CPI}}{\sum \text{Weight}} = \frac{0.4043}{0.06805}$$

$$= 5.9412$$

## Question: c)

The descriptive features in this dataset are of different types. For example, some are percentages, others are measured in years, and others are measured in counts per 1,000. We should always consider normalizing our data, but it is particularly important to do this when the descriptive features are measured in different units. What value would a 3-nearest neighbor prediction model using Euclidean distance return for the CPI of Russia when the descriptive features have been normalized using range normalization?

## Solution:

Microsoft Excel - Normalization

File Home Insert Page Layout Formulas Data Review View Automate Help Acrobat

AutoSave On Calculations Saved Search

Clipboard Font Alignment Number Styles Cells Editing Add-ins Analyze Data Create PDF and Share via Outlook Adobe Acrobat

G2: = (B2 - MIN(\$B\$2:\$B\$17)) / (MAX(\$B\$2:\$B\$17) - MIN(\$B\$2:\$B\$17))

**COUNTRY ID LIFE Exp TOP-10 INCOME INFANT MORT MILSPEND SCHOOL YEARS Normalized Life Expectancy Normalized TOP-10 INCOME Normalized INFANT MORT Normalized MIL SPEND Normalized SCHOOL YEARS Euclidean Distance between Russia and other countries**

COUNTRY ID	LIFE Exp	TOP-10 INCOME	INFANT MORT	MILSPEND	SCHOOL YEARS	Normalized Life Expectancy	Normalized TOP-10 INCOME	Normalized INFANT MORT	Normalized MIL SPEND	Normalized SCHOOL YEARS	Euclidean Distance between Russia and other countries
2 Afghanistan	59.61	23.21	74.3	4.44	0.4	0.393906713	0.04453125	0.896508728	0.651197605	0	1.275402958
3 Haiti	45	47.67	73.1	0.09	3.4	0	1	0.881546135	0	0.217391304	1.474660231
4 Nigeria	51.3	38.23	82.6	1.07	4.1	0.169857104	0.63125	1	0.146706587	0.268115942	1.288744017
5 Egypt	70.48	26.55	19.6	1.86	5.3	0.686977622	0.176171875	0.21446384	0.26497006	0.355072464	0.673646529
6 Argentina	75.77	32.3	13.3	0.76	10.1	0.829603667	0.399609375	0.135910224	0.100299401	0.702898551	0.55541483
7 China	74.87	29.98	13.7	1.95	6.4	0.805338366	0.308984375	0.140897756	0.278443114	0.434782609	0.590943968
8 Brazil	73.12	42.93	14.5	1.43	7.2	0.758155837	0.81484375	0.150872818	0.200598802	0.492753623	0.722691746
9 Israel	81.3	28.8	3.6	6.77	12.5	0.978700458	0.262890625	0.014962594	1	0.876811594	0.586832479
10 USA	78.51	29.85	6.3	4.72	13.7	0.903478026	0.30390625	0.048628429	0.693113772	0.963768116	0.33615085
11 Ireland	80.15	27.23	3.5	0.6	11.5	0.947894796	0.2015625	0.013715711	0.076347305	0.804347826	0.633115028
12 UK	80.09	28.49	4.4	2.59	13	0.946077111	0.25078125	0.024937656	0.374251497	0.913043478	0.412564205
13 Germany	80.24	22.07	3.5	1.31	12	0.950121327	0	0.013715711	0.182634731	0.84057971	0.643723893
14 Canada	80.99	24.79	4.9	1.42	14.2	0.97034241	0.10625	0.03117207	0.199101796	1	0.59145093
15 Australia	82.09	25.4	4.2	1.86	11.5	1	0.130078125	0.02244389	0.26497006	0.804347826	0.564307574
16 Sweden	81.43	22.18	2.4	1.27	12.8	0.982205446	0.004296875	0	0.176646707	0.898550725	0.660962784
17 New Zealand	80.67	27.81	4.9	1.13	12.3	0.961714748	0.22421875	0.03117207	0.155688623	0.862318841	0.566419158
18 Russia	67.62	31.68	10	3.87	12.9	0.609867889	0.375390625	0.094763092	0.5656868263	0.905797101	0
<b>Total</b>											
20											
21											
22											
23											

Normalization Good to go

Display Settings

Microsoft Excel - Normalization

File Home Insert Page Layout Formulas Data Review View Automate Help Acrobat

AutoSave On Calculations Saved Search

Clipboard Font Alignment Number Styles Cells Editing Add-ins Analyze Data Create PDF and Share via Outlook Adobe Acrobat

L10: =SQRT((G10-G\$18)^2 + (H10-H\$18)^2 + (I10-I\$18)^2 + (J10-J\$18)^2 + (K10-K\$18)^2)

**COUNTRY ID INFANT MORT MILSPEND SCHOOL YEARS Normalized Life Expectancy Normalized TOP-10 INCOME Normalized INFANT MORT Normalized MIL SPEND Normalized SCHOOL YEARS Euclidean Distance between Russia and other countries CPI**

COUNTRY ID	INFANT MORT	MILSPEND	SCHOOL YEARS	Normalized Life Expectancy	Normalized TOP-10 INCOME	Normalized INFANT MORT	Normalized MIL SPEND	Normalized SCHOOL YEARS	Euclidean Distance between Russia and other countries	CPI
2 Afghanistan	74.3	4.44	0.4	0.393906713	0.04453125	0.896508728	0.651197605	0	1.275402958	1.5171
3 Haiti	73.1	0.09	3.4	0	1	0.881546135	0	0.217391304	1.474660231	1.7999
4 Nigeria	82.6	1.07	4.1	0.169857104	0.63125	1	0.146706587	0.268115942	1.288744017	2.4493
5 Egypt	19.6	1.86	5.3	0.686977622	0.176171875	0.21446384	0.26497006	0.355072464	0.673646529	2.8622
6 Argentina	13.3	0.76	10.1	0.829603667	0.399609375	0.135910224	0.100299401	0.702898551	0.55541483	2.9961
7 China	13.7	1.95	6.4	0.805338366	0.308984375	0.140897756	0.278443114	0.434782609	0.590943968	3.6356
8 Brazil	14.5	1.43	7.2	0.758155837	0.81484375	0.150872818	0.200598802	0.492753623	0.722691746	3.7741
9 Israel	3.6	6.77	12.5	0.978700458	0.262890625	0.014962594	1	0.876811594	0.568532479	5.8069
10 USA	6.3	4.72	13.7	0.903478026	0.30390625	0.048628429	0.693113772	0.963768116	0.33615085	7.1357
11 Ireland	3.5	0.6	11.5	0.947894796	0.2015625	0.013715711	0.076347305	0.804347826	0.633115028	7.536
12 UK	4.4	2.59	13	0.946077111	0.25078125	0.024937656	0.374251497	0.913043478	0.412564205	7.7751
13 Germany	3.5	1.31	12	0.950121327	0	0.013715711	0.182634731	0.84057971	0.643723893	8.0461
14 Canada	4.9	1.42	14.2	0.97034241	0.10625	0.03117207	0.199101796	1	0.59145093	8.6725
15 Australia	4.2	1.86	11.5	1	0.130078125	0.02244389	0.26497006	0.804347826	0.564307574	8.8442
16 Sweden	2.4	1.27	12.8	0.982205446	0.004296875	0	0.176646707	0.898550725	0.660962784	9.2985
17 New Zealand	4.9	1.13	12.3	0.961714748	0.22421875	0.03117207	0.155688623	0.862318841	0.566419158	9.4627
18 Russia	10	3.87	12.9	0.609867889	0.375390625	0.094763092	0.5656868263	0.905797101	0	?
<b>Total</b>										
20										
21										
22										
23										

Normalization Good to go

Display Settings

c) From the table above we  
see that 3NN are:  
USA, UK and Argentina

CPI will be the average of  
the 3 NN:

$$\frac{7.1357 + 7.7751 + 2.9961}{3} = 5.9689$$

$$= 5.9689$$

#### Question: d)

What value would a weighted k-NN prediction model—with k = 16 (i.e., the full dataset) and using a weighting scheme of the reciprocal of the squared Euclidean distance between the neighbor and the query—return for the CPI of Russia when it is applied to the range-normalized data?

## Solution:

Excel screenshot showing the formula  $=1/(L2)^2$  entered in cell M2. The formula is highlighted with a green border.

	A	H	I	J	K	L	M	N	O	P	Q	R	S
1	COUNTRY ID	Normalized TOP-10 INCOME	Normalized INFANT MORT	Normalized MIL SPEND	Normalized SCHOOL YEARS	Eculidean Distance between Russia and other countries	Weight	CPI	Weight * CPI				
2	Afghanistan	0.04453125	0.896508728	0.651197605	0	1.275402958	0.614759375	1.5171	0.93265145				
3	Haiti	1	0.881546135	0	0.217391304	1.474860231	0.459725157	1.7999	0.82745931				
4	Nigeria	0.63125	1	0.146706587	0.268115942	1.288744017	0.602097295	2.4493	1.4747169				
5	Egypt	0.176171875	0.21446384	0.264697006	0.355072464	0.673646529	2.203615646	2.8622	6.3071887				
6	Argentina	0.399609375	0.135910224	0.100299401	0.702898651	0.55541483	3.241642042	2.9961	9.71228372				
7	China	0.308984375	0.140897756	0.278443114	0.434782609	0.590943968	2.863567281	3.6356	10.4107852				
8	Brazil	0.81484375	0.150872818	0.200598802	0.492753623	0.722691746	1.914669462	3.7741	7.22615402				
9	Israel	0.262890625	0.014962594	1	0.876811594	0.586832479	2.903833529	5.8069	16.8622709				
10	USA	0.30390625	0.048628429	0.693113772	0.963768116	0.33615085	8.849761619	7.1357	63.149244				
11	Ireland	0.2015625	0.013715711	0.076347305	0.804347826	0.633115028	2.494794372	7.536	18.8007704				
12	UK	0.25078125	0.024937656	0.374251497	0.913043478	0.412564205	5.875122282	7.7751	45.6796633				
13	Germany	0	0.013715711	0.182634731	0.84057971	0.643723893	2.413241261	8.0461	19.4171805				
14	Canada	0.10625	0.03117207	0.199101796	1	0.59145093	2.856660379	8.6725	24.7917321				
15	Australia	0.130078125	0.02244389	0.264697006	0.804347826	0.564307574	3.140279046	8.8442	27.7732559				
16	Sweden	0.004296875	0	0.176646707	0.898550725	0.660962784	2.2890102	9.2985	21.284276				
17	New Zealand	0.22421875	0.03117207	0.155688623	0.862318841	0.566419158	3.11690906	9.4627	29.4943754				
18	Russia	0.375390625	0.094763092	0.565686263	0.905797101	0	?						
19	Total						45.84167882		304.144008				

Excel screenshot showing the formula  $=M2*N2$  entered in cell O2. The formula is highlighted with a green border.

	A	H	I	J	K	L	M	N	O	P	Q	R	S
1	COUNTRY ID	Normalized TOP-10 INCOME	Normalized INFANT MORT	Normalized MIL SPEND	Normalized SCHOOL YEARS	Eculidean Distance between Russia and other countries	Weight	CPI	Weight * CPI				
2	Afghanistan	0.04453125	0.896508728	0.651197605	0	1.275402958	0.614759375	1.5171	0.93265145				
3	Haiti	1	0.881546135	0	0.217391304	1.474860231	0.459725157	1.7999	0.82745931				
4	Nigeria	0.63125	1	0.146706587	0.268115942	1.288744017	0.602097295	2.4493	1.4747169				
5	Egypt	0.176171875	0.21446384	0.264697006	0.355072464	0.673646529	2.203615646	2.8622	6.3071887				
6	Argentina	0.399609375	0.135910224	0.100299401	0.702898651	0.55541483	3.241642042	2.9961	9.71228372				
7	China	0.308984375	0.140897756	0.278443114	0.434782609	0.590943968	2.863567281	3.6356	10.4107852				
8	Brazil	0.81484375	0.150872818	0.200598802	0.492753623	0.722691746	1.914669462	3.7741	7.22615402				
9	Israel	0.262890625	0.014962594	1	0.876811594	0.586832479	2.903833529	5.8069	16.8622709				
10	USA	0.30390625	0.048628429	0.693113772	0.963768116	0.33615085	8.849761619	7.1357	63.149244				
11	Ireland	0.2015625	0.013715711	0.076347305	0.804347826	0.633115028	2.494794372	7.536	18.8007704				
12	UK	0.25078125	0.024937656	0.374251497	0.913043478	0.412564205	5.875122282	7.7751	45.6796633				
13	Germany	0	0.013715711	0.182634731	0.84057971	0.643723893	2.413241261	8.0461	19.4171805				
14	Canada	0.10625	0.03117207	0.199101796	1	0.59145093	2.856660379	8.6725	24.7917321				
15	Australia	0.130078125	0.02244389	0.264697006	0.804347826	0.564307574	3.140279046	8.8442	27.7732559				
16	Sweden	0.004296875	0	0.176646707	0.898550725	0.660962784	2.2890102	9.2985	21.284276				
17	New Zealand	0.22421875	0.03117207	0.155688623	0.862318841	0.566419158	3.11690906	9.4627	29.4943754				
18	Russia	0.375390625	0.094763092	0.565686263	0.905797101	0	?						
19	Total						45.84167882		304.144008				

Microsoft Excel - [New Book] [Untitled Workbook]

File Home Insert Page Layout Formulas Data Review View Automate Help Acrobat

AutoSave On

Calculations + Saved

Search

M19 : =SUM(M2:M17)

Clipboard

Font

Alignment

Number

Styles

Cells

Editing

Comments Share

General Conditional Formatting Insert Delete Sort & Filter Select Cells Cell Styles Format Cells Editing Sensitivity Add-ins Analyze Data Create PDF and Share link Create PDF and Share via Outlook Adobe Acrobat

A H I J K L M N O P Q R S

**Ecclidean Distance between Russia and other countries**

COUNTRY ID	Normalized TOP-10 INCOME	Normalized INFANT MORT	Normalized MIL SPEND	Normalized SCHOOL YEARS	Ecclidean Distance between Russia and other countries	Weight	CPI	Weight * CPI
1 Afghanistan	0.04453125	0.896508728	0.651197605	0	1.275402958	0.614759375	1.5171	0.93265145
2 Haiti	1	0.881546135	0	0.217391304	1.474860231	0.459725157	1.7999	0.82745931
3 Nigeria	0.63125	1	0.146706587	0.268115942	1.288744017	0.602097295	2.4493	1.4747169
4 Egypt	0.176171875	0.21446384	0.26497006	0.355072464	0.673646529	2.203615646	2.8622	6.3071887
5 Argentina	0.399609375	0.135910224	0.100299401	0.702896551	0.55541483	3.241642042	2.9961	9.71228372
7 China	0.308984375	0.140897756	0.278443114	0.434782609	0.590943968	2.863567281	3.6356	10.4107852
8 Brazil	0.81484375	0.150872818	0.200598802	0.492753623	0.722691746	1.914669462	3.7741	7.2615402
9 Israel	0.262890625	0.014962594	1	0.876811594	0.586832479	2.903833529	5.8069	16.8622709
10 USA	0.30390625	0.048628429	0.693113772	0.963768116	0.33615085	8.849761619	7.1357	63.149244
11 Ireland	0.2015625	0.013715711	0.076347305	0.804347826	0.633115028	2.494794372	7.536	18.8007704
12 UK	0.25078125	0.024937656	0.374251497	0.913043478	0.412564205	5.87512282	7.7751	45.6796633
13 Germany	0	0.013715711	0.182634731	0.84057971	0.643723893	2.413241261	8.0461	19.4171805
14 Canada	0.10625	0.03117207	0.199101796	1	0.59145093	2.856660379	8.6725	24.7917321
15 Australia	0.130078125	0.02244389	0.26497006	0.804347826	0.564307574	3.140279046	8.8442	27.7732559
16 Sweden	0.004296875	0	0.176646707	0.898550725	0.660962784	2.28900102	9.2985	21.284276
17 New Zealand	0.22421875	0.03117207	0.155688623	0.862318841	0.566419158	3.11690906	9.4627	29.4943754
18 Russia	0.375390625	0.094763092	0.565868263	0.905797101	0	?		
19 Total					45.84167882		304.144008	

20

21

22

23

< > Normalization +

Display Settings

Ready Accessibility: Good to go

100%

Microsoft Excel - [New Book] [Untitled Workbook]

File Home Insert Page Layout Formulas Data Review View Automate Help Acrobat

AutoSave On

Calculations + Saved

Search

O19 : =SUM(O2:O17)

Clipboard

Font

Alignment

Number

Styles

Cells

Editing

Comments Share

General Conditional Formatting Insert Delete Sort & Filter Select Cells Cell Styles Format Cells Editing Sensitivity Add-ins Analyze Data Create PDF and Share link Create PDF and Share via Outlook Adobe Acrobat

A H I J K L M N O P Q R S

**Ecclidean Distance between Russia and other countries**

COUNTRY ID	Normalized TOP-10 INCOME	Normalized INFANT MORT	Normalized MIL SPEND	Normalized SCHOOL YEARS	Ecclidean Distance between Russia and other countries	Weight	CPI	Weight * CPI
1 Afghanistan	0.04453125	0.896508728	0.651197605	0	1.275402958	0.614759375	1.5171	0.93265145
2 Haiti	1	0.881546135	0	0.217391304	1.474860231	0.459725157	1.7999	0.82745931
3 Nigeria	0.63125	1	0.146706587	0.268115942	1.288744017	0.602097295	2.4493	1.4747169
4 Egypt	0.176171875	0.21446384	0.26497006	0.355072464	0.673646529	2.203615646	2.8622	6.3071887
5 Argentina	0.399609375	0.135910224	0.100299401	0.702896551	0.55541483	3.241642042	2.9961	9.71228372
7 China	0.308984375	0.140897756	0.278443114	0.434782609	0.590943968	2.863567281	3.6356	10.4107852
8 Brazil	0.81484375	0.150872818	0.200598802	0.492753623	0.722691746	1.914669462	3.7741	7.2615402
9 Israel	0.262890625	0.014962594	1	0.876811594	0.586832479	2.903833529	5.8069	16.8622709
10 USA	0.30390625	0.048628429	0.693113772	0.963768116	0.33615085	8.849761619	7.1357	63.149244
11 Ireland	0.2015625	0.013715711	0.076347305	0.804347826	0.633115028	2.494794372	7.536	18.8007704
12 UK	0.25078125	0.024937656	0.374251497	0.913043478	0.412564205	5.87512282	7.7751	45.6796633
13 Germany	0	0.013715711	0.182634731	0.84057971	0.643723893	2.413241261	8.0461	19.4171805
14 Canada	0.10625	0.03117207	0.199101796	1	0.59145093	2.856660379	8.6725	24.7917321
15 Australia	0.130078125	0.02244389	0.26497006	0.804347826	0.564307574	3.140279046	8.8442	27.7732559
16 Sweden	0.004296875	0	0.176646707	0.898550725	0.660962784	2.28900102	9.2985	21.284276
17 New Zealand	0.22421875	0.03117207	0.155688623	0.862318841	0.566419158	3.11690906	9.4627	29.4943754
18 Russia	0.375390625	0.094763092	0.565868263	0.905797101	0	?		
19 Total					45.84167882		304.144008	

20

21

22

23

< > Normalization +

Display Settings

Ready Accessibility: Good to go

100%

d) From the above table the value returned by the model is

$$\frac{\sum \text{weight} \times \text{CPI}}{\sum \text{weight}} = \frac{304.1440}{45.8416}$$

$$= 6.6346$$

### Question: e)

The actual 2011 CPI for Russia was 2.4488. Which of the predictions made was the most accurate? Why do you think this was?

### Solution:

The most accurate prediction is the **3-Nearest Neighbor without normalization** (4.59), though it still overestimates the actual CPI for Russia. The normalization process seems to have increased the predicted CPI in both models.

## 3. (15%) Chapter 5, exercise 4

### Question: a)

The company has decided to use a similarity-based model to implement the recommender system. Which of the following three similarity indexes do you think the system should be based on?

### Solution:

For a large data set considering co-absences are not meaningful, as there will be more sparse data. Jaccard similarity index is ideal as it ignores co-absences. **Jaccard Similarity** is likely the best choice because:

- The data is binary (either bought or not bought).

- Jaccard focuses on the overlap in purchases between two customers, which is important for making recommendations based on shared interests.
- It only considers the items both customers have bought (true-true matches), which aligns well with the goal of recommending similar items.

### Question: b)

What items will the system recommend to the following customer? Assume that the recommender system uses the similarity index you chose in the first part of this question and is trained on the sample dataset listed above. Also assume that the system generates recommendations for query customers by finding the customer most similar to them in the dataset and then recommending the items that this similar customer has bought but that the query customer has not bought.

Solution:

③ Chapters , Ex-4

b)

ID	Item	Item	Item	Item	Item
2	107	498	7256	28063	7534
1	true	true	true	false	false
2	true	false	false	true	true

Query true false true and false false

For previous question I chose  
Jaccard( $x, y$ ) =  $\frac{C_P(x, y)}{C_P(x, y) + P_A(x, y) + A(y)}$

Jaccard( $q, d_1$ ) =  $\frac{2}{2+1+0} = 0.6667$

$$\text{Jaccard}(q, d_2) = \frac{1}{1+2+1} = 0.25$$

$$\text{Jaccard}(q, d_1) > \text{Jaccard}(q, d_2)$$

So  $d_1$  is more similar to  $q$ .

There is only one item (498) the query customer has not bought and  $d_1$  has bought.

So the system will recommend item 498 to the query customer.

for this example in all similarity metrics  $d_1$  is more similar to  $q$ .

$$\text{Russell-Rao}(q, d_1) = \frac{2}{5} = 0.4$$

$$\text{Russell-Rao}(q, d_2) = \frac{1}{5} = 0.2$$

$$(q, d_1) > (q, d_2)$$

$$\text{Sokal-Michener}(q, d_1) = \frac{4}{5} = 0.8$$

$$\text{Sokal-Michener}(q, d_2) = \frac{2}{5} = 0.4$$

So the system will recommend 498 regardless of similarity metrics.

#### 4. (15%) Chapter 5, exercise 5

##### Question: a)

A good measure of distance between two instances with categorical features is the overlap metric (also known as the hamming distance), which simply counts the number of descriptive features that have different values. Using this measure of distance, compute the distances between the mystery animal and each of the animals in the animal dataset.

##### Solution:

The overlap metric between the query instance and each instance in the dataset by counting the number of feature values that are different.

ID	CLASS	Overlap Metric
1	Mammal	6
2	Amphibian	1
3	Mammal	6
4	Bird	2

##### Question b):

If you used a 1-NN model, what class would be assigned to the mystery animal?

##### Solution:

D2 is the first nearest neighbor, so the mystery animal would be Amphibian.

##### Question c):

If you used a 4-NN model, what class would be assigned to the mystery animal? Would this be a good value for k for this dataset?

##### Solution:

4-NN covers all the instances in the dataset as there are only 4 records. As a result, any query would be assigned the majority class in the dataset, that is mammal in this case. So, for this dataset, a 4-NN model would underfit the dataset. 4 is not a good value of k for this dataset.

## 5. (15%) Chapter 5, exercise 6

Question a):

Create a k-d tree for this dataset. Assume the following order over the features: RENT then SIZE.

Solution:

5) Chapter 5, Ex-6				
ID	SIZE	RENT	PRICE	
1	2700	9235	2,000,000	
2	1315	1800	820,000	
3	1050	1250	800,000	
4	2200	7000	1,750,000	
5	1800	3800	1,450,500	
6	1900	4000	1,500,500	
7	960	800	720,000	

First level (SIZE):

Sort by size

ID	SIZE	RENT	PRICE
7	960	800	720,000
3	1050	1250	800,000
2	1315	1800	820,000
5	1800	3800	1,450,000
6	1900	4000	1,500,500
4	2200	7000	1,750,000
1	2700	9235	2,000,000

Median => ID 5    1800    3800    1,450,000

Left Subtree    Size < 1800

Right Subtree    Size > 1800

Left Subtree:

ID	Size	Rent	Price
7	960	800	720000
3	1050	1250	800000
2	1315	1800	820000

Sort by RENT.

Medium:

3 1050 1250 800000

Left of Left Subtree:

Size < 1050

7 960 800 720000

This is the leaf node.

Right of Left Subtree:

Size > 1050

2 1315 1800 820000

Leaf node.

Right Subtree: ( $> 1800$ )

6	1900	4000	1500500
4	2200	7000	1750000
1	2700	9235	2000000

Medium:

4 2200 7000 1750000

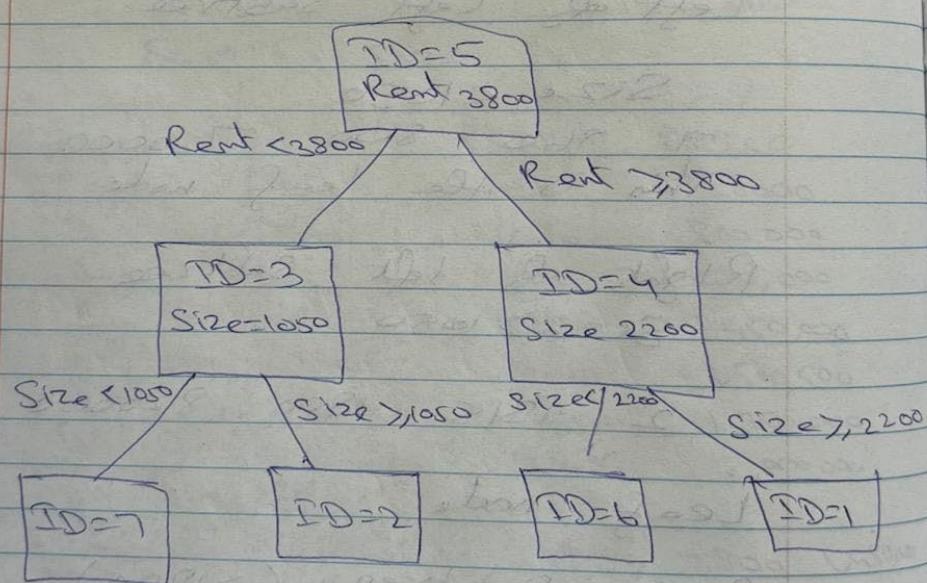
Left of Right Subtree ( $< 2200$ )

6 1900 4000 1500500

Leaf node.

Right of Right Subtree ( $> 2200$ )

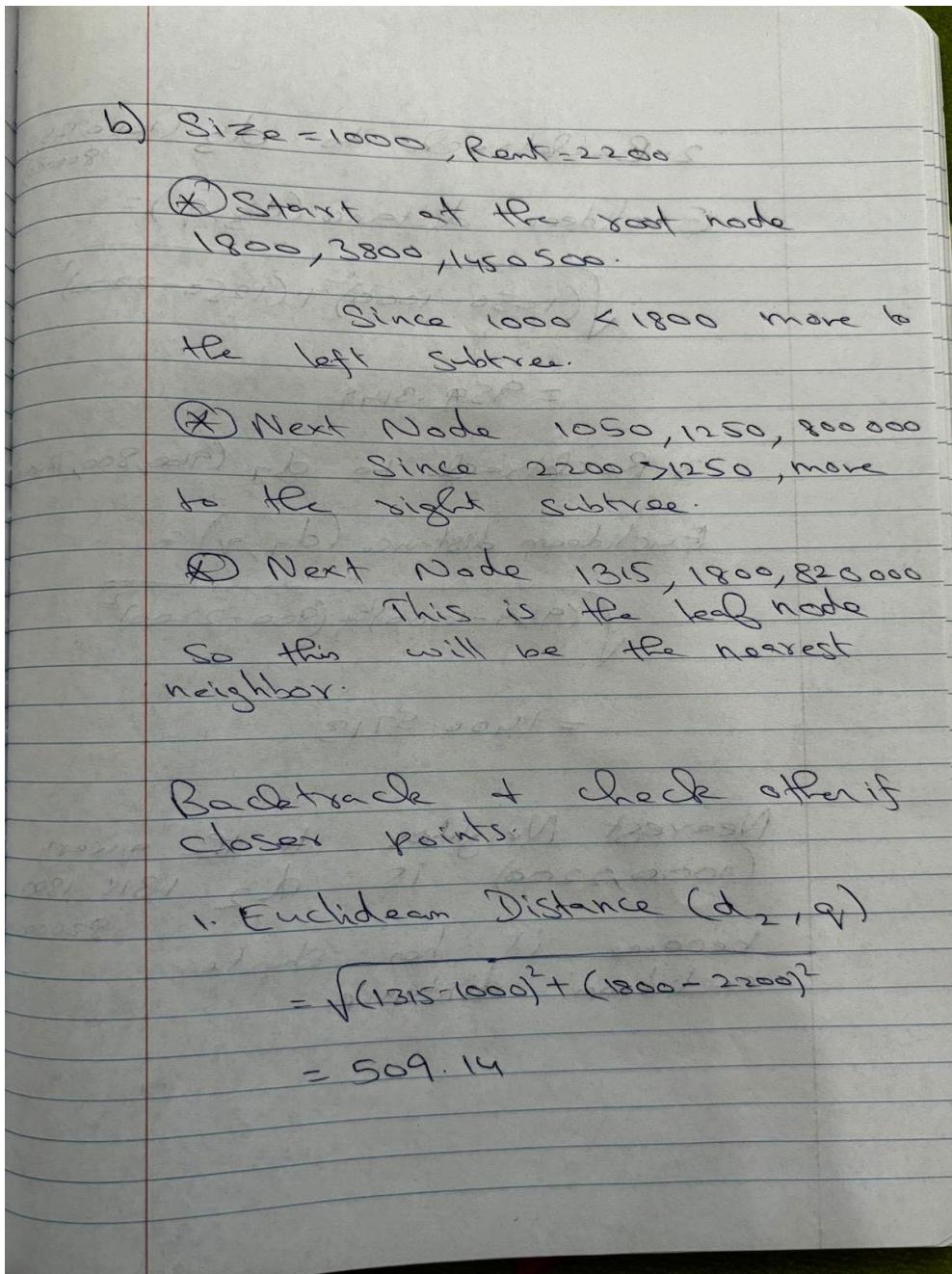
2700 9235 2000000



### Question b):

Using the k-d tree that you created in the first part of this question, find the nearest neighbor to the following query: SIZE = 1;000, RENT = 2;200.

### Solution:



2. Backtrack to  $d_3$ :  $(1050, 1250, 800000)$

Euclidean distance ( $d_3, v$ ) =

$$\sqrt{(1050 - 1000)^2 + (1250 - 2200)^2}$$
$$= 959.3149$$

3. Backtrack to  $d_7$ :  $(960, 800, 72000)$

Euclidean distance ( $d_7, v$ ) =

$$\sqrt{(960 - 1000)^2 + (800 - 2200)^2}$$
$$= 1406.5713$$

Nearest Neighbor to the query  
 $(1000, 2200)$  is  $d_2$  1315 1800  
8200000

because it has shorter  
euclidean distance.