

# CIDM 6355 Data Mining Methods HW3 Template

(100 points; Due 11:59 PM Central Time, October 20, 2024)

Requirements: This homework is open book, open slides, and open notes, but no collaboration or discussion is permitted before the due time. Any questions about the homework should be directed to the instructor. You must adhere to the instructions, completing all questions and deliverables. This is an individual assignment, so sharing your processes, scripts, screenshots, or answers with others constitutes cheating and will be reported. Additionally, ensure your answers meet the required format to avoid point deductions. Screenshots without date and time will receive a penalty of 50% of points. Identical screenshots will be considered academic dishonesty and will be reported to the university authority. Please acknowledge your understanding and agreement to these requirements by typing your name below.

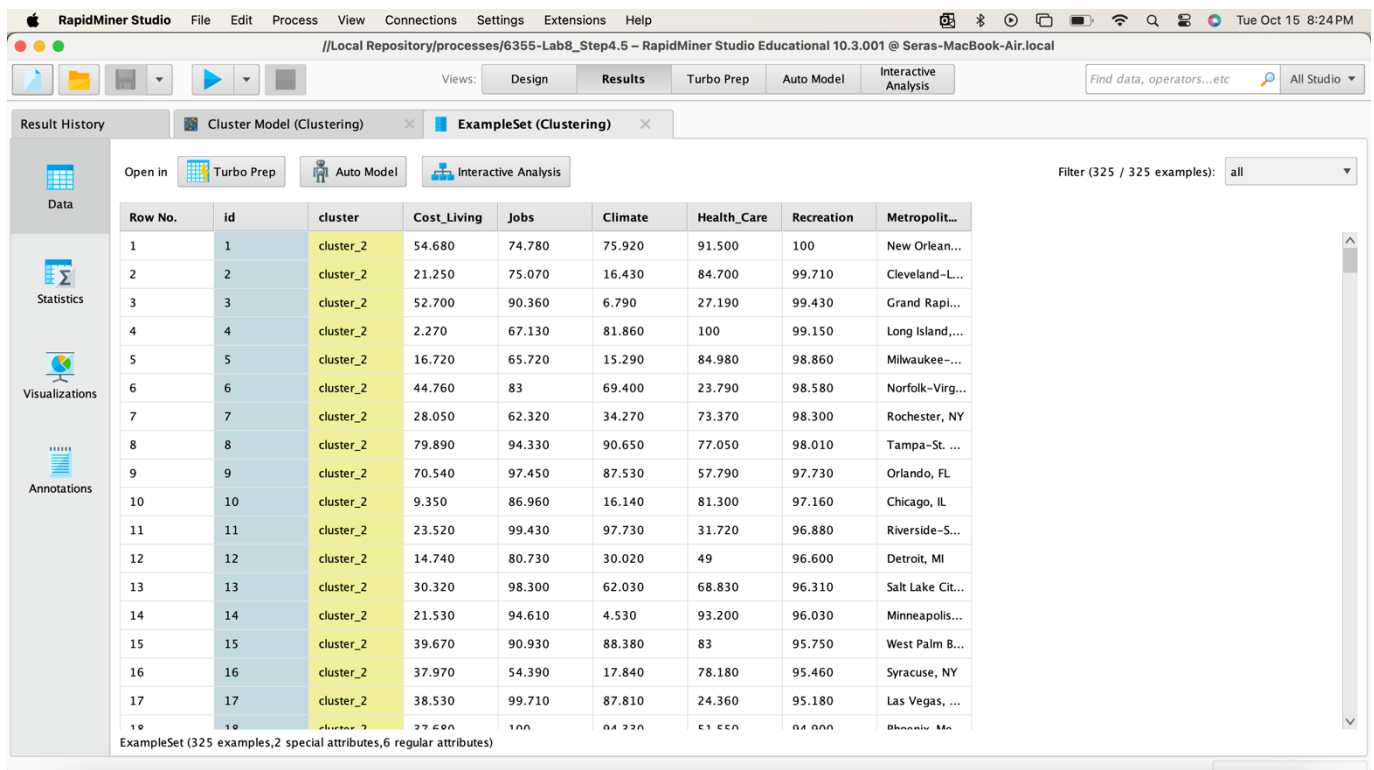
Type your name: Sera Hill

Instruction: Please compile all the deliverables with the required format as below.

1. **Deliverable 1 (Step 1):** Please write down the average for all the five attributes (round them the third decimal place). All these numbers below are the overall centroid for all 325 cities. [5 points]

Attributes	Cost_living	Jobs	Climate	Health_Care	Recreation
Average	51.910	51.023	52.035	47.865	50.227

2. **Deliverable 2 (Step 4.5):** Take a screenshot of your ExampleSet (Screenshot 1) [5 points]



Row No.	Id	cluster	Cost_Living	Jobs	Climate	Health_Care	Recreation	Metropolit...
1	1	cluster_2	54.680	74.780	75.920	91.500	100	New Orlean...
2	2	cluster_2	21.250	75.070	16.430	84.700	99.710	Cleveland-L...
3	3	cluster_2	52.700	90.360	6.790	27.190	99.430	Grand Rapi...
4	4	cluster_2	2.270	67.130	81.860	100	99.150	Long Island,...
5	5	cluster_2	16.720	65.720	15.290	84.980	98.860	Milwaukee-...
6	6	cluster_2	44.760	83	69.400	23.790	98.580	Norfolk-Virg...
7	7	cluster_2	28.050	62.320	34.270	73.370	98.300	Rochester, NY
8	8	cluster_2	79.890	94.330	90.650	77.050	98.010	Tampa-St. ...
9	9	cluster_2	70.540	97.450	87.530	57.790	97.730	Oriando, FL
10	10	cluster_2	9.350	86.960	16.140	81.300	97.160	Chicago, IL
11	11	cluster_2	23.520	99.430	97.730	31.720	96.880	Riverside-S...
12	12	cluster_2	14.740	80.730	30.020	49	96.600	Detroit, MI
13	13	cluster_2	30.320	98.300	62.030	68.830	96.310	Salt Lake Cit...
14	14	cluster_2	21.530	94.610	4.530	93.200	96.030	Minneapolis...
15	15	cluster_2	39.670	90.930	88.380	83	95.750	West Palm B...
16	16	cluster_2	37.970	54.390	17.840	78.180	95.460	Syracuse, NY
17	17	cluster_2	38.530	99.710	87.810	24.360	95.180	Las Vegas, ...

3. Deliverable 3 (Step 4.8): based on the results in 4.5-4.8, please discuss the characteristics in each cluster and find an appropriate name for each cluster. For example, Cluster 0 includes 128 cities such as New Orleans, LA and Long Island, NY have highest scores in job opportunities, climate, healthcare, and recreation. However, this group of cities have quite high living cost. We can name this group of cities Metropolitan Luxury..... [21 points: 7 points for each cluster, including this cluster's sample size (1 pt.), sample cities (1 pt.), comparison on each dimension (4 pts), and name for this cluster (1 pt.)]

Cluster 0 has 97 cities which include Houma, LA, Punta Gorda, FL, Panama City, FL. They have the highest score in climate, but it is the most expensive for cost of living. This group will be named Metropolitan Luxury.

Cluster 1 has 101 cities which include Des Moines, IA, Fort Wayne, IN, Springfield, MA. They have the lowest score for jobs, and climate. This group will be called Metropolitan Austerity.

Cluster 2 has 127 cities which include Orlando, FL, Chicago, IL, Detroit, MI. They have the highest score for jobs, health care, and recreation while having a low cost of living. This group will be called Metropolitan Exemplar.

4. Deliverable 4 (Step 6.2): Take a screenshot of your Result History page (Screenshot 2) [5 points]

The screenshot shows the RapidMiner Studio interface with the 'Result History' tab selected. The interface displays three clusters of results:

- ExampleSet (Multiply)**: Result not stored in repository. Data Table: Number of examples = 325, 8 attributes. Attributes: Role, Name, Type, Range. Values: Cost\_Living (real, [7..7]; mean =?), Jobs (real, [7..7]; mean =?), Climate (real, [7..7]; mean =?).
- Performance Vector (Performance)**: Result not stored in repository. PerformanceVector: Avg. within centroid distance: -2460.728, Avg. within centroid distance\_cluster\_0: -235, Avg. within centroid distance\_cluster\_1: -242, Avg. within centroid distance\_cluster\_2: -256, Davies Bouldin: -1.570.
- Centroid Cluster Model (Clustering)**: Result not stored in repository. Cluster 0: 97 items, Cluster 1: 101 items, Cluster 2: 127 items, Total number of items: 325.

The bottom of the screenshot shows a list of process results for '6355-Lab8\_Step5.5' (2 results, Process results) with completion times and execution times.

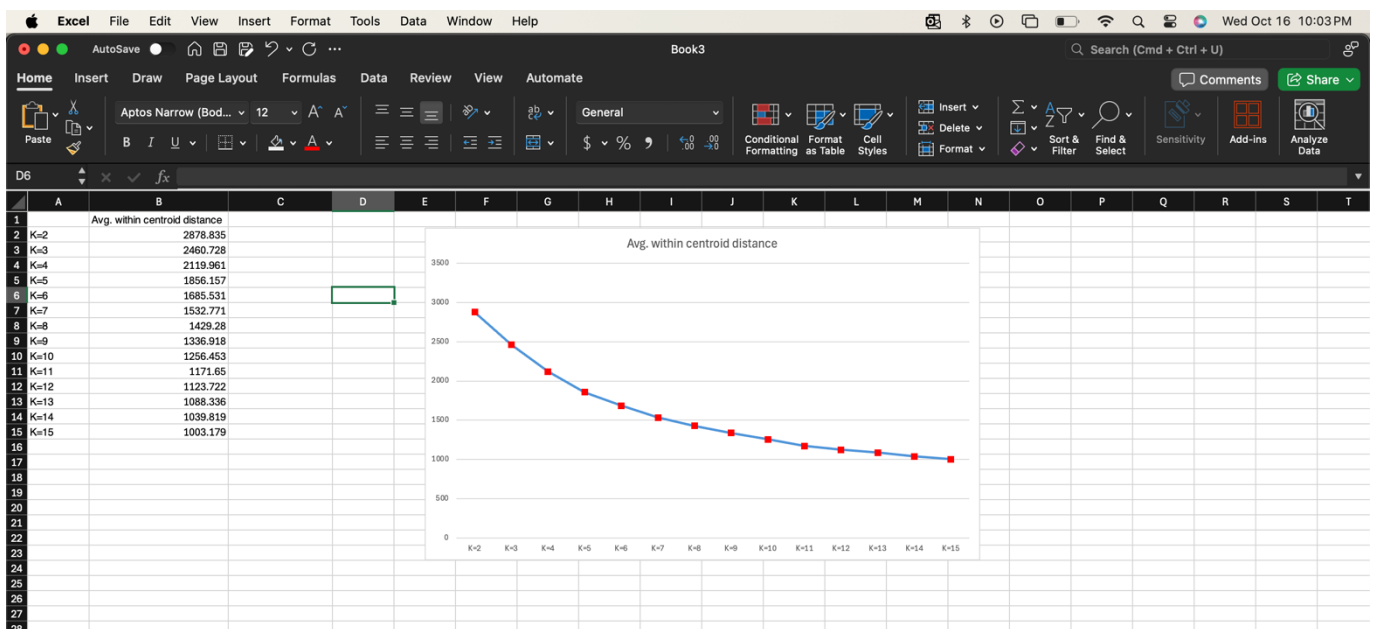
5. Deliverable 5 (Step 6.2): Please answer all the question in this deliverable [9 points]

- Based on the table above, when k increases, what happen to Avg. within centroid distance

(increasing or decreasing)? [2 points] It decreases.

- What about Davies Bouldin Index when  $k$  increases? [2 points] It alternates between decreasing and increasing.
- Imagine an extreme case, when  $k=325$ , what would Avg. within centroid distance be? [2 points] Zero.
- What potential problem will we encounter if we only use Avg. within centroid distance as the main criterion for evaluating clustering models? [3 points] We would be ignoring the size and shape of the cluster, and kind of ignores the possible presence of outliers.

6. Deliverable 6 (Step 7.1): Draw an elbow chart using either average within centroid distance or DBI for  $k=2-15$ . Take a screenshot of your elbow chart with date and time (Screenshot 3). Observe your elbow chart and discuss which  $k$  is the best and why. [10 points: 5 points for screenshot and 5 points for your discussion]



The best  $k$  would be one between 4 and 5 because this is when the elbow starts to flatten out.

7. Deliverable R1: take a screenshot of the result after running the script in Line 19 with date and time (Screenshot 4) and time and briefly interpret the result, explaining what each portion of results means. Your interpretation should cover the following five portions:
- K-means clustering
  - Cluster means
  - Clustering vector
  - Within Cluster Sum of squares
  - Available components

[illegible]

The clustering vector shows us which group each datapoint belongs to. With a value of 1 belonging to group 1, value of 2 belonging to group 2 and value of 3 belonging to group 3. The within cluster sum of squares by cluster, we can see that group 1's datapoints are closer to the centroid than the other groups, with group 3 having the furthest distance between from the centroid.

The available components show us which values that were used during our k-means analysis. Cluster tells us where each point belongs, centers show us the average position of each point in each cluster, totss shows us total variance, and withinss, and tot.withinss shows how well the points fit in each cluster. Betweenss shows how much separation each cluster has, and size tells us how many points are in each cluster. Iter is how many iterations the algorithm went through and ifault will let us know of any problems in the algorithm's execution.

8. Deliverable R2: take a screenshot of the result after running the script in Line 24 with date and time (Screenshot 5) and time and briefly interpret the result, explaining what the result is about and what each column means. [10 points: 5 points for screenshot and 5 points for your interpretation]

The screenshot shows the RStudio interface. The script editor on the left contains R code for performing k-means clustering on the 'Lab8Data' dataset. The console on the right displays the output of the clustering process, including the assignment of cities to clusters and a summary table of cluster sizes and means.

```

# view the top ten rows in the dataset
4 head(Lab8Data)
# view the descriptive statistics for all attributes in Lab8Data
5 summary(Lab8Data)
# view the structure of the dataset
6 str(Lab8Data)
# generate a correlation matrix for numerical attributes to see if we have high
7 cor(Lab8Data[,c(2:6)])
# set the seed to make sure you can get the same result as lab instruction
8 set.seed(100)
# use kmeans function for clustering, which includes three parameters: data, n
9 # data: use columns 2-6 for clustering
10 # we use 3 as initial k
11 # specify nstart as 100, r will try 100 different random starting assignments
12 CityCluster <- kmeans(Lab8Data[, 2:6], 3, nstart = 100)
13 # check the clustering result
14 CityCluster
15 # generate a table to see which city belongs to which cluster
16 table(CityCluster$cluster, Lab8Data$Metropolitan_Area)
17 # display a dataframe to show the number of observations in each cluster and the
18 data.frame(CityCluster$size, CityCluster$center)
19
20
21
22
23
24
25
26
27

```

The console output shows the following results:

```

Topeka, KS Trenton, NJ Tucson, AZ Tulsa, OK Tuscaloosa, AL Tyler, TX Utica-Rome, NY
1 0 1 0 0 0 0 1
2 1 0 0 0 0 1 1
3 0 0 1 1 0 0 0

Vallejo-Fairfield-Napa, CA Ventura, CA Victoria, TX Vineland-Millville-Bridgeton, NJ
1 0 0 0 0 0
2 0 0 1 1
3 1 1 0 0

Visalia-Tulare-Porterville, CA Waco, TX Washington, DC-MD-VA-WV Waterbury, CT
1 0 0 0 1
2 1 1 0 0
3 0 0 1 0

Waterloo-Cedar Falls, IA Wausau, WI West Palm Beach-Boca Raton, FL Wheeling, WV-OH
1 1 1 0 1
2 0 0 0 0
3 0 0 1 0

Wichita Falls, TX Wichita, KS Williamsport, PA Wilmington-Newark, DE-MD Wilmington, NC
1 0 1 0 0 0
2 1 0 1 0 0
3 0 0 0 1 1

Worcester, MA-CT Yakima, WA Yolo, CA York, PA Youngstown-Warren, OH Yuba City, CA Yuma, AZ
1 1 0 0 1 1
2 0 1 0 0 1
3 0 0 1 0 0

```

The final output is a data frame with 3 columns: CityCluster.size, Cost\_Living, Jobs, Climate, Health\_Care, and Recreation. The rows represent the three clusters.

CityCluster.size	Cost_Living	Jobs	Climate	Health_Care	Recreation
92	41.04761	33.73391	23.47891	45.55793	50.02924
105	75.78324	34.22686	56.92467	33.51457	21.94552
128	40.13391	77.22789	68.54938	61.29539	73.56758

The result shows us how many observations are in each cluster as well as the mean of each attribute in the cluster. First column shows us cluster size, the following columns are the mean of the datapoints for each category: cost of living, jobs, climate, health care, and recreation.

9. Deliverable R3: take a screenshot of the result after running the script in Line 28 with date and time (Screenshot 6) and time, and briefly interpret the result, explaining what the result is about and what each column means. [10 points: 5 points for screenshot and 5 points for your interpretation]

```

7 # view the structure of the dataset
8 str(Lab8Data)
9 # generate a correlation matrix for numerical attributes to see if we have high
10 cor(Lab8Data[,c(2:6)])
11 # set the seed to make sure you can get the same result as lab instruction
12 set.seed(100)
13 # use kmeans function for clustering, which includes three parameters: data, n
14 # data: use columns 2-6 for clustering
15 # we use 3 as initial k
16 # specify nstart as 100, n will try 100 different random starting assignments
17 CityCluster <- kmeans(Lab8Data[, 2:6], 3, nstart = 100)
18 # check the clustering result
19 CityCluster
20
21 # generate a table to see which city belongs to which cluster
22 table(CityCluster$cluster, Lab8Data$Metropolitan_Area)
23 # display a dataframe to show the number of observations in each cluster and the
24 data.frame(CityCluster$size, CityCluster$center)
25 # create a new dataframe to contain clusterID and the five attributes for each
26 CityRecords <- data.frame(CityCluster$cluster, Lab8Data[,c(1:6)])
27 # check the first few rows of cityrecords
28 head(CityRecords)
29
30

```

```

> data.frame(CityCluster$size, CityCluster$center)
CityCluster.size Cost_Living Jobs Climate Health_Care Recreation
1 92 41.04761 33.73391 23.47891 45.55793 50.02924
2 105 75.78324 34.22686 56.92467 33.51457 21.94552
3 128 40.13391 77.22789 68.54938 61.29539 73.56758
> CityRecords <- data.frame(CityCluster, Lab8Data[,c(1:6)])
Error in as.data.frame.default(x[[i]], optional = TRUE, stringsAsFactors = stringsAsFactors) :
cannot coerce class "kmeans" to a data.frame
> CityRecords <- data.frame(CityCluster$cluster, Lab8Data[,c(1:6)])
> head(CityRecords)
CityCluster.cluster Metropolitan_Area Cost_Living Jobs Climate Health_Care Recreation
1 3 New Orleans, LA 54.68 74.78 75.92 91.50
2 3 Cleveland-Lorain-Elyria, OH 21.25 75.07 16.43 84.70
3 3 Grand Rapids-Muskegon-Holland, MI 52.70 90.36 6.79 27.19
4 3 Long Island, NY 2.27 67.13 81.86 100.00
5 3 Milwaukee-Waukeasha, WI 16.72 65.72 15.29 84.98
6 3 Norfolk-Virginia Beach-Newport News, VA-NC 44.76 83.00 69.40 23.79

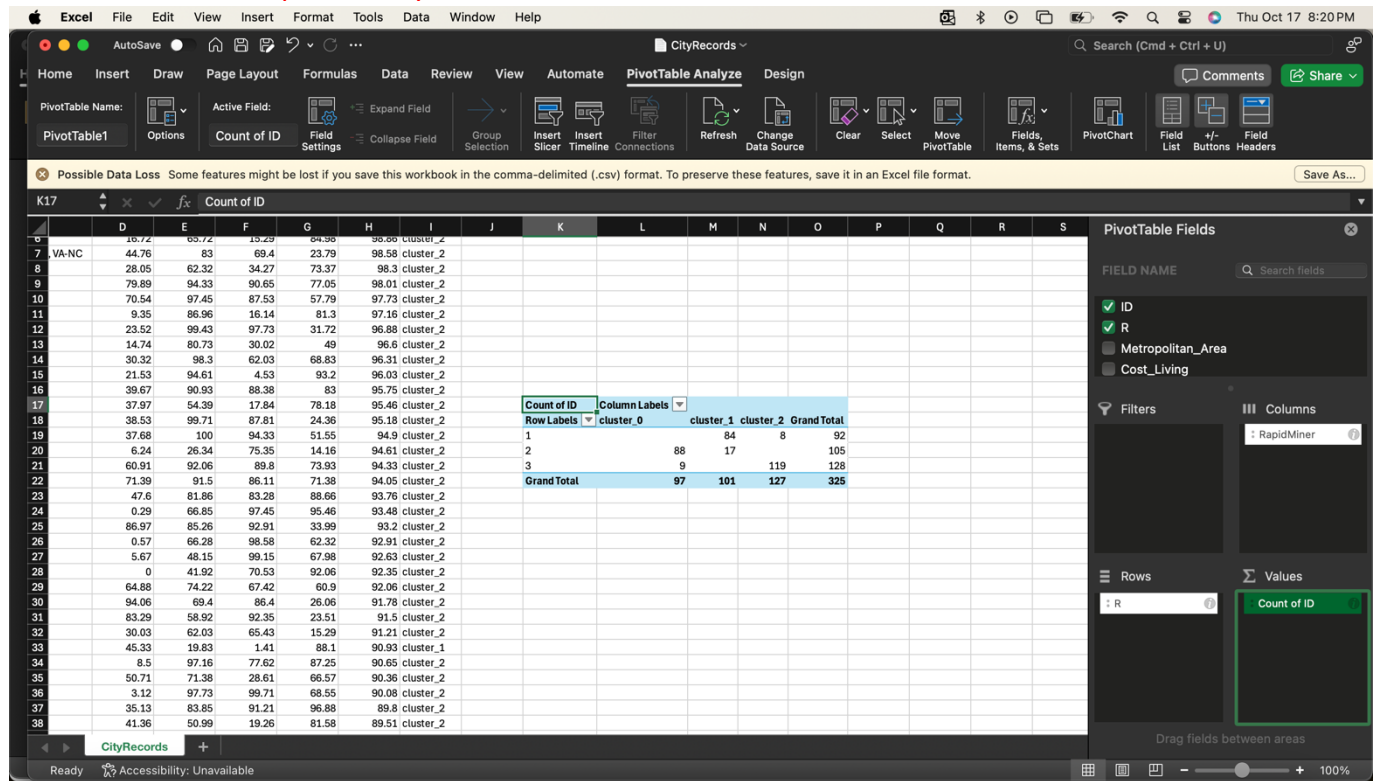
```

This result shows the first few rows of the CityCluster dataframe. This dataframe contains information about the different cities and what cluster they belong to (CityCluster.cluster and Metropolitan\_Area columns), as well as the mean of each attribute in each cluster (Cost\_Living, Jobs, Climate, Health\_Care, Recreation columns).

10. **Deliverable** R4: Compare the clustering result for each observation in R (which is saved in CityRecords.csv) and that in RapidMiner (k=3 only). **Compare the two clustering results and answer the question: Are the two clustering results in R and RM the same or not? Why?** You may follow the instruction in the next slide and take a screenshot of your PivotTable with date and time to support your answer (**Screenshot 7**). Attention: you cannot just simply compare the cluster name because R and RM may label each cluster differently. For example, New Orleans, LA is labeled as cluster\_0 in RM, but Cluster 3 in R, but cluster\_0 in RM might be the same with Cluster 3 in R. **[10 points: 5 points]**



for screenshot and 5 points for your answer]



The two clustering results are not the same because R and RapidMiner may be using different algorithms for the k-means analysis. They also may have different ways of handling preprocessing, how they handle equal distances, or how they scale the data.