

# **Project 2: Big Data Management**

“Advanced Hadoop”

Ehu Shubham Shaw

Supreeth

Roee shraga

October 10, 2024

## 1- Revisiting LinkBook Network with Pig [16 Points]

The way I chose to implement the task is by using Pig latin. For Task A Pig Requires fewer lines of code and is easier to maintain compared to the custom Mapper and Reducer classes in Java. Wheres in MapReduce Offers more control over the logic and performance but requires more boilerplate code. Task B pig is easier for aggregations like top N. The TOP function simplifies this task, making it more intuitive. in MapReduce You need a custom reducer to track and sort page accesses, which adds complexity. in Task C Pig Uses Simple filtering operation with a straightforward query. Where as MapReduce Requires writing a custom Mapper to filter the data, which is less intuitive than Pig's syntax. In Task D Pig Simplifies joins and counting. Pig's high-level syntax saves time and effort.where MapReduce Requires custom Mapper and Reducer for counting and handling relationships, which increases complexity. Task E pig supports counting and distinct operations in a straightforward way. On the other hand MapReduce Requires more work to implement distinct counting, needing additional logic.Task F Pig uses calculation of averages is easier and more abstracted. And in MapReduce Requires a custom reducer to calculate the average and compare each user's relationship count. Task G Pig Date comparison and filtering are easier.compared to MapReduce Requires custom logic to handle dates and filtering. Task H Joins and filtering based on relationships and accesses are handled easily in pig. And for in MapReduce Requires complex logic to join and filter the data.

Pig Implementation: To calculate the frequency of each education level for Task a, utilise GROUP and COUNT on the HighestEdu field. Utilize GROUP, COUNT, and TOP to determine the ten most frequently accessed pages for Task B. Use FILTER in Task c to pick persons whose HighestEdu matches yours. For Task d, count relationships and calculate the pleasure factor by using JOIN on LinkBookPage and Associates. To determine the total and unique accesses in Task e, utilise GROUP and COUNT on the AccessLog. Determine the average number of relationships for Task f and eliminate people with more relationships than the norm. To find people who haven't accessed LinkBook in ninety days, use FILTER on the AccessLog in Task g.Lastly, for Task h, find individuals who have relationships but have never visited their friends' LinkBook sites by using JOIN on Associates and AccessLog.

**TASK A**

FOLDERS	
└── pig	
└── output_A	
└── _SUCCESS.crc	
└── .part-r-00000.crc	
└── _SUCCESS	
└── part-r-00000	
└── output_B	
└── output_C	
└── output_D	
└── output_G	
└── output_task_H	
└── output_taskf	
└── Task_A.pig	
└── Task_B.pig	
└── Task_C.pig	
└── Task_D.pig	
└── Task_E.pig	
└── Task_F.pig	
└── Task_G.pig	
└── Task_H.pig	
part-r-00000	
1	PHD 28571
2	Masters 28572
3	Primary 28572
4	Bachelors 28571
5	Secondary 28572
6	Elementary 28572
7	PlaySchool 28571

**Task B**

1	51870	51870	Michaela Murray Engineer
2	74635	74635	Kathleen Hartman Electrician
3	80846	80846	Samantha Mack Designer
4	89979	89979	Joanna Nolan Mechanic
5	89997	89997	Alison Powers Analyst
6	102383	102383	Priscilla Mason Nurse
7	108547	108547	Samantha Thomas Analyst
8	126009	126009	Karen Wagner Mechanic
9	154271	154271	David Lynch Scientist
10	182560	182560	Jennifer Fox Engineer

## Task B

**Task C**

1	3	James Bennett	Nurse
2	10	Maria Lee	Engineer
3	17	Audrey Harris	Analyst
4	24	Jon Collins	Manager
5	31	Crystal Johnson	Scientist
6	38	Eileen Kim	Chef
7	45	Alexander Yang	Electrician
8	52	Cindy Garcia	Teacher
9	59	Jean Ramos	Mechanic
10	66	Kimberly Pena	Designer
11	73	David Graham	Nurse
12	80	Maria Long	Engineer
13	87	Kelly Herrera	Analyst
14	94	Michael Hensley	Manager
15	101	Joyce Thomas	Scientist
16	108	Samantha Rice	Chef
17	115	Matthew Kline	Electrician
18	122	Amy Frazier	Teacher
19	129	Randy Herrera	Mechanic
20	136	Ms. Alison Jackson	Designer

**Task D**

1	Charles Anderson	1	
2	Patrick Miller	99	
3	Dale Sweeney	103	
4	James Bennett	99	
5	Bradley Peterson	MD 98	
6	Anthony Williams	107	
7	Jonathan Williams	92	
8	Tonya Davis	103	
9	Ashley York	84	
10	Donna Jones	87	
11	Maria Lee	89	
12	Kyle Cooper	103	
13	Tiffany Wong	107	
14	Patrick Contreras	94	
15	Jeffrey Shaffer	101	
16	Cynthia Cooper	112	
17	Michael Contreras	93	
18	Audrey Harris	114	
19	Kyle Morgan	98	
20	Kimberly Ford	103	

## Task D

# Task E

FOLDERS		part-r-00000			
▼	pig	1	1	49	49
►	output_A	2	2	58	58
►	output_B	3	3	54	54
►	output_C	4	4	50	50
►	output_D	5	5	60	60
►	output_G	6	6	51	51
►	output_task_H	7	7	58	58
►	output_taskf	8	8	49	49
▼	OutputE	9	9	38	38
	._SUCCESS.crc	10	10	56	56
	.part-r-00000.crc	11	11	53	53
	._SUCCESS	12	12	49	49
	part-r-00000	13	13	58	58
	Task_A.pig	14	14	50	50
	Task_B.pig	15	15	51	51
	Task_C.pig	16	16	54	54
	Task_D.pig	17	17	47	47
	Task_E.pig	18	18	57	57
	Task_F.pig	19	19	47	47
	Task_G.pig	20	20	44	44
	Task_H.pig	21	21	45	45

# Task F

FOLDERS		part-r-00000			
▼	pig	1	Patrick Miller,99		
►	output_A	2	Dale Sweeney,103		
►	output_B	3	James Bennett,99		
►	output_C	4	Bradley Peterson MD,98		
►	output_D	5	Anthony Williams,107		
►	output_G	6	Jonathan Williams,92		
►	output_task_H	7	Tonya Davis,103		
▼	output_taskf	8	Ashley York,84		
	._SUCCESS.crc	9	Donna Jones,87		
	.part-r-00000.crc	10	Maria Lee,89		
	._SUCCESS	11	Kyle Cooper,103		
	part-r-00000	12	Tiffany Wong,107		
	Task_A.pig	13	Patrick Contreras,94		
	Task_B.pig	14	Jeffrey Shaffer,101		
	Task_C.pig	15	Cynthia Cooper,112		
	Task_D.pig	16	Michael Contreras,93		
	Task_E.pig	17	Audrey Harris,114		
	Task_F.pig	18	Kyle Morgan,98		
	Task_G.pig	19	Kimberly Ford,103		
	Task_H.pig	20	Anthony Harrison,105		

# Task G

FOLDERS		part-r-00000			
▼	pig	1	Charles Anderson		
►	output_A	2	Patrick Miller		
►	output_B	3	Dale Sweeney		
►	output_C	4	James Bennett		
►	output_D	5	Bradley Peterson MD		
►	output_G	6	Anthony Williams		
▼	output_task_H	7	Jonathan Williams		
	._SUCCESS.crc	8	Tonya Davis		
	.part-r-00000.crc	9	Ashley York		
	._SUCCESS	10	Donna Jones		
	part-r-00000	11	Maria Lee		
	output_taskf	12	Kyle Cooper		
	Task_A.pig	13	Tiffany Wong		
	Task_B.pig	14	Patrick Contreras		
	Task_C.pig	15	Jeffrey Shaffer		
	Task_D.pig	16	Cynthia Cooper		
	Task_E.pig	17	Michael Contreras		
	Task_F.pig	18	Audrey Harris		
	Task_G.pig	19	Kyle Morgan		
	Task_H.pig	20	Kimberly Ford		

FOLDERS
pig
output_A
output_B
output_C
output_D
output_G
_SUCCESS.crc
.part-r-00000.crc
.SUCCESS
part-r-00000
output_task_H
output_taskf
Task_A.pig
Task_B.pig
Task_C.pig
Task_D.pig
Task_E.pig
Task_F.pig
Task_G.pig
Task_H.pig

  

part-r-00000			
1	0	Charles Anderson	
2	1	Patrick Miller	
3	10	Maria Lee	
4	100	Gary Zimmerman	
5	1000	Shelby Martinez	
6	10000	Shaun Schmidt	
7	100000	Kristine Silva	
8	100001	Kristen Benson	
9	100002	Angela Lopez	
10	100003	Andrew Morgan	
11	100004	Michael Novak	
12	100005	Joseph Moore	
13	100006	Rebecca Rice	
14	100007	Donna Gentry	
15	100008	Sarah Anderson	
16	100009	Tonya Chen	
17	10001	Derek Roberts	
18	100010	Sophia Ross	
19	100011	Mary Joyce	
20	100012	Cory Stein	

# Task H

## 2- Clustering [44 Points]

For Task 1 (Single-Iteration K-Means, R=1), the solution runs only one iteration of the K-means algorithm, assigning each data point to the nearest centroid and updating the centroids once. I tested with different values of K (the number of clusters) to see how the clustering changed depending on the initial centroid pick. Because this task only requires one repetition, it is computationally fast. However, the quality of the clusters is strongly influenced by the centroids' initial placements.

In Task 2 (Multi-Iteration K-Means with K=10), the centroids are modified after each iteration using the points provided to them. Running studies with different K values revealed how the centroid positions change. However, if convergence comes quickly, the method may perform unneeded computations, rendering subsequent iterations redundant. Performance improves with decreasing K, but plateaus after a few repetitions.

In Task 3 (Multi-Iteration K-Means with Early Termination), the method employs early stopping when centroids converge below a certain threshold ( $\text{epsilon} = 0.001$ ) before reaching the maximum number of iterations. I experimented with different K and R numbers to see how often the algorithm converged before reaching the maximum number of iterations. In one experiment, it converged at 36 iterations, significantly reducing computation time. This strategy improves performance when the dataset allows for rapid convergence.

For Task 4 (K-Means using Hadoop MapReduce Combiner), I employed a combiner to improve data shuffling and aggregation throughout the MapReduce process. The combiner conducts local aggregation at the mapper level before transferring data to the reducer, which reduces computation time. For example, when K=100 was used, the combiner reduced runtime by around one minute compared to running without it.

Task 5 (Clustered Points Evaluation) displayed the final centroid coordinates and the number of points allocated to each cluster. I tested the technique with K values of 5, 10, and 100 and noticed that the points allocated to clusters changed with each iteration. The distance between centroid locations rarely remained constant between iterations. As K grows, clusters get more granular and exact, but calculation time increases as well. More iterations (higher R) improve accuracy while increasing computing time, whereas early termination balances accuracy and efficiency by ending when the clusters stabilize.

# Kaggle dataset

# Python dataset

```

== [ASK (a): Single-Iteration K-Means (R=1) ==

== Centroids ==
Centroid 0: -17740.987917420223, -32.1684657347778, 10.57629271312487
Centroid 1: -2514.97129953554, -151.5264521804564, -1,20855839548794
Centroid 2: -32860.765504206676, 183.78441897614778, 14.959081837988
Centroid 3: -378944979171, 41.8136213937129, 7.84525874245664
Centroid 4: 34674.399623087164, 69.9959088174459, -31.619689426619725
Task (b): Basic Multi-Iteration K-Means (R=10) ==
== Centroids ==
Centroid 0: -21733.330053328043, -457.4049553091517, 8.6757130334633772
Centroid 1: -34760.4881788816, 1803.35901299771, 17.12832984998042
Centroid 2: -8625.72843781891, 773.468178839218, -1,89619167510315
Centroid 3: -6722.956248959, 1899.977173189216, 6.0517634851625
Centroid 4: -3451.323245235176, -187.397781123793, 4.075324956860145
Centroid 5: -11137.323245235176, -187.397781123793, 4.075324956860145
Centroid 6: -34510.718071358, 291.176244520446, 31.63038574371951
Centroid 7: -8346.29437810853, 18.45436594523257, -13.859947289412
Centroid 8: -21517.8074270475, 37.81124656469398, -6.62160585498191
Centroid 9: -8081.29437810853, 73.7621338347052, 1.5060189853865192
Centroid 10: -21137.929245250176, -107.3977811217893, 4.075321958686245
Centroid 11: -34510.718071358, 291.176244520446, 32.68398574371951
Centroid 12: -8081.29437810853, 73.7621338347052, 1.5060189853865192
Centroid 13: -8081.29437810853, 73.7621338347052, 1.5060189853865192
Centroid 14: -20648.2739878316, 82.1986556717474, 19.2686723554664
Centroid 15: -34626.73946794997, 229.47298177198, 23.67997437262663
Centroid 16: -7611.999946586286, -172.5073968788499, -25.98839724561157
Centroid 17: -30188.474270809257, -28.28806586801584, -19.160693929207488
Centroid 18: -9094.8526565958, 21.859141393814, 12.971339977884417
Centroid 19: -20648.6329789616, 82.1986556717474, 19.2686723554664
Centroid 20: -10163.1862479934, 35.4648541678795, 1.5060189853865192
Centroid 21: -561.999946586286, 172.5073968788499, -2.083973452461157
Centroid 22: -30188.474270809257, -28.28806586801584, -19.160693929207488
Centroid 23: -9094.8526565958, 21.859141393814, 12.971339977884417
Centroid 24: -20220.39666899853, 67.758788842898, 19.8543817978822
Centroid 25: -34012.4214729934, 25.615654878215, 21.826605497515
Centroid 26: -6697.544985108782, -203.9679855495343, -31.98402753735079
Centroid 27: -31895.4693835957, -8.76335898152163, 1.539402816243875
Centroid 28: -18386.144985108782, -203.9679855495343, -31.98402753735079
Centroid 29: -30285.4693835957, 57.758788842898, 19.8543817978822
Centroid 30: -34012.4214729934, 25.615654878215, 21.826605497515
Centroid 31: -6697.544985108782, -203.9679855495343, -31.98402753735079
Centroid 32: -31895.4693835957, -8.76335898152163, -21.538402816243875
Centroid 33: -10163.1862479934, 35.4648541678795, 1.506018985386223
Centroid 34: -33846.73618293, 335.449811377719, 26.4088723398285
Centroid 35: -1545.799946586286, -10.087272501598850
Centroid 36: -10163.1862479934, 35.4648541678795, 1.506018985386223
Centroid 37: -30285.116148576544, -33.19399539765349, -36.368184287047
Centroid 38: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 39: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 40: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 41: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 42: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 43: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 44: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 45: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 46: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 47: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 48: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 49: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 50: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 51: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 52: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 53: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 54: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 55: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 56: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 57: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 58: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 59: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 60: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 61: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 62: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 63: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 64: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 65: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 66: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 67: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 68: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 69: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 70: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 71: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 72: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 73: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 74: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 75: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 76: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 77: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 78: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 79: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 80: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 81: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 82: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 83: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 84: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 85: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 86: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 87: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 88: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 89: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 90: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 91: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 92: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 93: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 94: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 95: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 96: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 97: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 98: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 99: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 100: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 101: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 102: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 103: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 104: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 105: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 106: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 107: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 108: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 109: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 110: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 111: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 112: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 113: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 114: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 115: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 116: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 117: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 118: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 119: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 120: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 121: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 122: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 123: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 124: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 125: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 126: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 127: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 128: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 129: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 130: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 131: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 132: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 133: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 134: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 135: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 136: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 137: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 138: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 139: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 140: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 141: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 142: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 143: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 144: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 145: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 146: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 147: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 148: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 149: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 150: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 151: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 152: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 153: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 154: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 155: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 156: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 157: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 158: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 159: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 160: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 161: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 162: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 163: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 164: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 165: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 166: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 167: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 168: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 169: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 170: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 171: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 172: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 173: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 174: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 175: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 176: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 177: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 178: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 179: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 180: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 181: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 182: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 183: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 184: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 185: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 186: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 187: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 188: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 189: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 190: -5754.799602435291, -234.75787118867082, -10.0721255198857
Centroid 191: -11283.466633834987, 64.684787172533, 16.91013300836534
Centroid 192: -32652.16914857644, -33.19399539765349, -30.368184287047
Centroid 193: -10163.1862479934, 35.4791352623431, 1.613194985630223
Centroid 194: -33846.73618293, 335.4498113777194, 26.4088723398285
Centroid 
```

## Task C (Convergence)

```
8603 Centroid 77: 26507.098306342945, -3754.6024568275384, -114.81572853612407
8604 Centroid 78: 2664.3911349585655, 4633.546605849815, -72.78405420841924
8605 Centroid 79: 27160.010547531307, -144.74624394884535, 10.748196661672313
8606 Centroid 80: 29959.98514230785, 3750.339168797623, -140.14671298387756
8607 Centroid 81: 30935.60244327642, -3958.940828767624, -65.62045726931798
8608 Centroid 82: 31913.506919508884, 302.5114671711909, -49.655570810210826
8609 Centroid 83: 38230.21816946891, -1683.4176948318511, 159.15279709763237
8610 Centroid 84: 3858.8245529511332, -2154.246852337512, -72.45768656955322
8611 Centroid 85: 39044.93347018802, 3899.1424275668874, -182.2934928498707
8612 Centroid 86: 3954.886713311318, -4658.733940294479, -44.50593759896657
8613 Centroid 87: 49140.88620105533, -630.948216568444, -54.60059699006024
8614 Centroid 88: 5360.055118609245, 1184.3395827764414, 297.7363504140519
8615 Centroid 89: 5783.591089654008, 3554.8133677466026, 268.58739993092223
8616 Centroid 90: 5957.501790322647, -619.3604285673913, 143.60941819884488
8617 Centroid 91: 5970.260299364924, -3766.8256469550556, -11.492241266061269
8618 Centroid 92: 6024.662859297017, 4838.551606226644, 457.515770632767
8619 Centroid 93: 614396.844620168, 2801.998755726393, 910.6786307685968
8620 Centroid 94: 720.1306302987209, -4702.304573767187, 160.5557635332215
8621 Centroid 95: 8133.39314167297, 3734.445966481086, -272.21723403121
8622 Centroid 96: 8185.385249447876, -1969.4558528353314, -102.69047185746001
8623 Centroid 97: 8547.212690979919, 722.4271422460675, 226.1320156561345
8624 Centroid 98: 8703.85173399999, -4040.262461915782, 174.69338620870212
8625 Centroid 99: 9684.949590794306, -5157.0303212563595, -472.6634819501443
8626 === Task (c): Multi-Iteration K-Means with Early Termination ===
8627 Converged after 42 iterations.
8628 === Centroids ===
8629 Centroid 0: -10191.538973596045, -4671.745749539993, 73.42963581640163
8630 Centroid 1: -11722.617003402802, -2457.2831874257054, 99.90201334700126
8631 Centroid 2: -11751.723502073253, 650.6529379299419, -4.834190346333643
8632 Centroid 3: -12394.414231308006, 4245.869821736836, 84.7327297794913
8633 Centroid 4: -12569.334961473749, -4395.683769664479, -174.74691656365803
8634 Centroid 5: -13870.709175047663, -1416.901626085707, -26.20258264915002
8635 Centroid 6: -14237.126886254311, 1056.6213440075578, -150.68847469788722
8636 Centroid 7: -15380.955428899828, 2334.37038049334, 211.94432563114637
8637 Centroid 8: -15537.553927370547, 4624.97434873775, -8.621699840605238
8638 Centroid 9: -15588.36367570346, -4095.183558409099, -11.308557770107804
8639 Centroid 10: -1633.9322433939901, -4375.6091647099625, -174.72169401593783
8640 Centroid 11: -16663.864356244198, -816.5651351911104, 6.627421788698217
8641 Centroid 12: -18067.037562194528, -3250.1667042152276, 1.5647122717507544
8642 Centroid 13: -18253.51450399208, -4865.845324412149, -364.4971366506231
8643 Centroid 14: -18473.811262452833, 3722.46659351445, 388.8404160798132
8644 Centroid 15: -18574.97773541348, -472.9127084780152, 201.1962527405747
8645 Centroid 16: -1950.0985453618673, -460.1025353547972, 226.58697172910485
8646 Centroid 17: -19723.720211476175, 1342.4913041179398, -209.7449826595355
```

### 3- Get Creative! [40 Points]

#### a. BYOD (Bring Your Own Data):

- **Customer Personality Analysis** is a technique that allows businesses to acquire a better knowledge of their ideal customers. Businesses can tailor their products and services to the unique demands, behaviors, and concerns of each client segment by analyzing them. This targeted approach enables businesses to optimize marketing efforts by focussing on clients who are most likely to buy a product rather than selling to the complete customer base.
- For Pre-Processing the Dataset - We began using a Kaggle dataset containing patient data and 29 features. To make our Java code easier to process and build, we decreased the dataset's dimensionality to three characteristics, denoted by x, y, and z. This dimensionality reduction was done in Task 2 to simplify the dataset for computational efficiency.
- Data Cleaning and Standardization, Before lowering dimensionality, we cleaned the data. This stage required dealing with missing or incorrect values in the matrix's rows

and columns (r and c). After cleaning the dataset, we used standardization to guarantee that the features had a uniform scale, which is required for effective dimensionality reduction. Standardization rescales the data so that each feature has a mean of 0 and a standard deviation of 1, allowing the Principal Component Analysis (PCA) procedure to function more reliably.

- Lastly, Dimensionality Reduction Using PCA After cleaning and standardizing. The dataset, we reduced it to three dimensions using Principal Component Analysis (PCA). PCA is a method for converting data into a set of orthogonal components with the greatest variance. By condensing the dataset's information to three primary components, we were able to keep the most important parts. We used NumPy to reduce dimensionality. PCA created a dense matrix of updated data that we could easily input into our Java code for further processing. This stage streamlined the data, making it easier to handle and assess in future operations for our code.

b. Evaluate the output of the clustering algorithm “*Silhouette*”.

**K = 3** : The clustering technique performs well with a Silhouette score of 0.6186 (K = 3). This implies that the clusters are well-separated, with data points closer to their particular clusters than to others, implying that the K=3 solution is extremely effective for this dataset.

**Updated Clustering Evaluation Report:**

**1. K-Means Clustering Results:**

==== Centroids ===

Centroid 0: -21802.68390426056, 28.37865096345726, 11.36707411598497  
Centroid 1: 27797.449153536334, 44.017145706239766, -3.5132460876444624  
Centroid 2: 3734.459960024626, -63.86339219188164, -10.024505995564208

**2. Silhouette Score Evaluation:**

Overall Silhouette Score: 0.6186280655201676

Interpretation: [See interpretation above]

**K = 5** : A Silhouette score of 0.476 indicates acceptable grouping. Though the score is lower than for K=3, it still suggests adequate cluster separation and cohesiveness, with a little decrease in performance as the number of clusters grows.

The average silhouette score of 0.5473 for K=3 and K=5 indicates balanced performance, but further optimization may improve results. Both K=3 and K=5 offer plausible solutions, with K=3 marginally exceeding. Finally, the Silhouette evaluation demonstrates how the method performs

```

Updated Clustering Evaluation Report:
1. K-Means Clustering Results:
==== Centroids ===
Centroid 0: -18833.225458581812,-982.6554890294938,16.612895922471644
Centroid 1: -28512.2279590318,2602.1598554772518,12.389644658679542
Centroid 2: -40261.364989388196,2484.992129075358,60.42818034064102
Centroid 3: -43247.327029960456,-2914.1440401112536,-121.86159084876817
Centroid 4: 12773.271936976354,9.914676986659785,-8.328557297933298

2. Silhouette Score Evaluation:
Overall Silhouette Score: 0.4761188000180853

```

with various K values. The higher score for K=3 indicates better-defined clusters, but K=5 still performs pretty well.

**K=7 :** Total Silhouette Score: 0.4858 ,Interpretation: The clustering for K=7 lacks a strong or distinguishable structure. It suggests that the clusters may not be well-separated, with data points potentially being assigned to arbitrary clusters, implying that the algorithm or dataset requires further optimisation or changes.

```

==== Task (e): Silhouette Score Evaluation ===
Overall Silhouette Score: 0.4857700718891826
Silhouette Score Interpretation:
Weak structure found. The clustering may be arbitrary.

Updated Clustering Evaluation Report:
1. K-Means Clustering Results:
==== Centroids ===
Centroid 0: -12856.263837201734,328.46249210690564,-4.909393313138605
Centroid 1: -20506.928850287666,6.042531813107754,8.732081232098766
Centroid 2: -33016.19551870957,147.32488912040483,17.69195631424592
Centroid 3: -6355.963149896548,-695.3588878048778,-1.9881955558847955
Centroid 4: 14744.31559905834,-26.52936129732165,-5.356933516850448
Centroid 5: 2141.679289988163,120.51979055729242,28.399263803192554
Centroid 6: 33225.25814256535,45.73608158741366,-36.66985585672809

2. Silhouette Score Evaluation:
Overall Silhouette Score: 0.4857700718891826
Interpretation: [See interpretation above]

```

**K=50 :** Overall silhouette score: 0.4253. Interpretation: Clustering with K=50 has a weaker structure than prior evaluations, implying that the huge number of clusters may have resulted in overfitting or breaking data into numerous arbitrary or less meaningful groups. This result suggests that the clustering algorithm may not be well-suited to this value of K, and additional tweaking may be required.

## Updated Clustering Evaluation Report:

### 1. K-Means Clustering Results:

==== Centroids ===

Centroid 0: -10706.278296834293, -118.16595110416331, 18.02258017119816  
Centroid 1: -11561.25652663303, -4095.382410092767, -432.8849354258641  
Centroid 2: -12832.766978364254, 4594.925726892419, 3.4371348518488922  
Centroid 3: -13386.432647361027, -2486.6194702649814, 36.09976866113162  
Centroid 4: -15028.414955594382, 2265.6986467592055, 77.26048613719486  
Centroid 5: -16015.866204127122, -4800.664247424829, -264.12183133507324  
Centroid 6: -1785.7830956584623, 2809.1885031107363, -110.18146137546853  
Centroid 7: -17535.66136976127, -3021.9428744777333, 24.24768072712916  
Centroid 8: -18248.37893270136, 770.794754599433, 49.650654146422646  
Centroid 9: -20174.941101670738, 4588.838172925966, 179.1565049114489  
Centroid 10: -20920.74981996778, -1262.89730109506, -52.269624474570186  
Centroid 11: -21717.098288737423, 2154.6713609435606, 40.05568164992845  
Centroid 12: -22913.523202861616, -4076.8094176921104, 109.41209639282776  
Centroid 13: -23528.308846226933, -1829.087259672445, -517.5183493689922  
Centroid 14: -24952.14707727689, 1387.8833223302177, 63.448870943725446  
Centroid 15: -27973.51692653096, -3516.043815176994, -8.777182637064989  
Centroid 16: -28320.002628790426, 3898.519108456147, -26.18314336844569  
Centroid 17: -31545.36504114053, -1839.0116634227843, 157.57744531998327  
Centroid 18: -35559.70503906677, 1216.4181889225242, -16.22376538215205  
Centroid 19: -44350.99786288196, 1277.6314768693528, -34.589342060235765  
Centroid 20: -5373.616127427295, -2452.3167881994177, -15.319796744570233  
Centroid 21: -6155.649774323215, 4304.291824712866, 68.99412861149978  
Centroid 22: -8802.8954102653, 2524.180138978442, -80.80251506929311  
Centroid 23: -9250.348041000076, -4379.589040448941, 166.63058376629934  
Centroid 24: 9930.749631314118, -4236.239484743317, -84.81219095715522  
Centroid 25: 11180.213090964167, -1705.4919992022667, -73.51358861385756  
Centroid 26: 12388.402725555128, 683.551752224276, -60.83108730447821  
Centroid 27: 12826.477273909248, 4330.85903149043, 129.35418292119687  
Centroid 28: 14368.78640847944, 3881.43649302593, -203.70954653566397  
Centroid 29: 457.80135411870816, -1360.7601301067255, 35.15063543255052  
Centroid 30: 13875.438641809937, -4864.872445999399, 115.82251856359495  
Centroid 31: 14798.777915851335, -2441.5077565973193, 30.998256244193335  
Centroid 32: 16101.333061196869, 449.709678401825, -88.63801707440889  
Centroid 33: 17592.241367642397, 3866.360203547392, 41.29521841542299  
Centroid 34: 18302.13440610193, -3241.6312919602046, 267.1823419665626  
Centroid 35: 20212.96912341134, 8.271609988607828, -40.30801793056302  
Centroid 36: 20979.391673827315, 4323.199506709797, -216.67568620167367  
Centroid 37: 21961.885652014895, -3925.112750248122, 28.954863636090952  
Centroid 38: 25260.11928165965, 2534.9324014303866, -75.93075630738143  
Centroid 39: 25855.51785114567, -2050.733337077978, 50.02975032231433  
Centroid 40: 2831.192151336269, -4200.452564123226, -73.5070254005828  
Centroid 41: 29653.60143961574, 1409.6293154732416, 9.789436456203335  
Centroid 42: 31824.97577507028, -3738.2311049198865, -64.24417725589957  
Centroid 43: 3020.038513632848, 3207.1182497728287, 76.38907838068533  
Centroid 44: 37188.69676606274, 1816.1323887001945, 5.927393196543938  
Centroid 45: 114180.24895362444, 144.669014862147, -204.40019381472968  
Centroid 46: 6849.81354039599, 3662.936202753248, 75.29099362078107  
Centroid 47: 7414.846390557808, -640.6270909953435, 146.6128893808934  
Centroid 48: 10065.246435551326, 2462.8556708525693, -29.862473616317324  
Centroid 49: 10096.204322739286, 4441.151119724976, -55.46818460144751

### 2. Silhouette Score Evaluation:

Overall Silhouette Score: 0.42533327961968476