

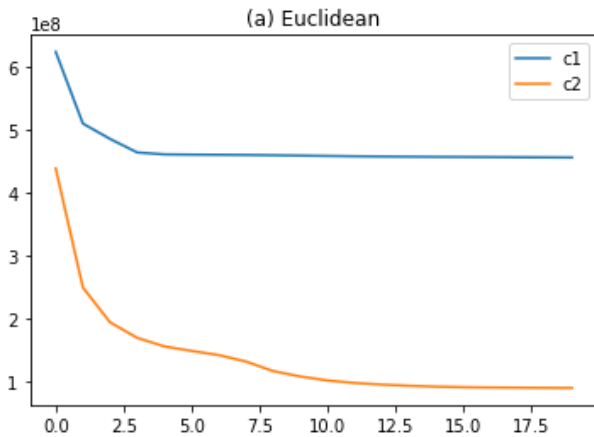
## Introduction to Massive Data Analysis

### HW3 - K-means Report

106070038 杜葳葳

#### (a) Initialization strategies with **Euclidean** distance

- A plot of cost vs iteration



	C1	C2
Round 1	623660345	438747790
Round 2	509862908.3	249803933.6
Round 3	485480682	194494814.4
Round 4	463997011.7	169804841.5
Round 5	460969266.6	156295748.8
Round 6	460537848	149094208.1
Round 7	460313099.7	142508531.6
Round 8	460003523.9	132303869.4
Round 9	459570539.3	117170969.8
Round 10	459021103.3	108547377.2
Round 11	458490656.2	102237203.3
Round 12	457944232.6	98278015.75
Round 13	457558005.2	95630226.12
Round 14	457290136.4	93793314.05
Round 15	457050555.1	92377131.97
Round 16	456892235.6	91541606.25
Round 17	456703630.7	91045573.83
Round 18	456404203	90752240.1
Round 19	456177800.5	90470170.18
Round 20	455986871	90216416.18

- Percentage improvement values and explanation

- Euclidean

- C1: 26.885%

- C2: **79.438%**

- Explanation

- c2 的表現比 c1 佳，初始選擇距離較遠的點作為 cluster 的中心相較隨機選取來的理想，推測如果是隨機選取 (c1)，可能會選到過於靠近的點，導致在過程中掉入區域最佳解、而非全域的最佳解。

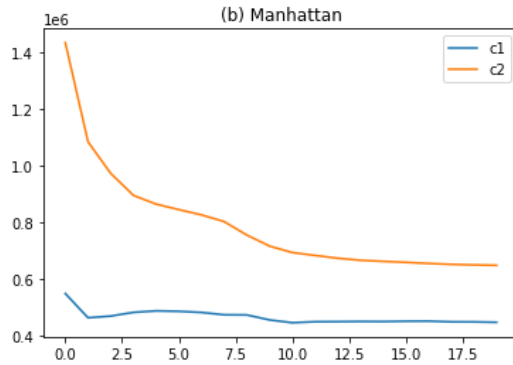
- The Euclidean and Manhattan Distances for all pairs of centroids, with 2 initialization strategies.





(b) Initialization strategies with **Manhattan** distance

- A plot of cost vs iteration



	C1	C2
Round 1	550117.142	1433739.31
Round 2	464869.276	1084488.777
Round 3	470897.382	973431.715
Round 4	483914.409	895934.593
Round 5	489216.071	865128.335
Round 6	487629.669	845846.647
Round 7	483711.923	827219.583
Round 8	475330.773	803590.346
Round 9	474871.239	756039.517
Round 10	457232.92	717332.903
Round 11	447494.386	694587.925
Round 12	450915.013	684444.502
Round 13	451250.367	674574.748
Round 14	451974.596	667409.47
Round 15	451570.364	663556.628
Round 16	452739.011	660162.777
Round 17	453082.73	656041.322
Round 18	450583.671	653036.754
Round 19	450368.749	651112.426
Round 20	449011.364	649689.013

- Percentage improvement values and explanation

- Manhattan

- C1: 18.379%

- C2: 54.686%

- Explanation

- C1 和 C2 相比，與(a)相同，跑多次 iteration 後 C2 皆比 C1 來的好，初始時隨

機選取容易誤掉到區域最佳解。

- 與(a)綜合比較，我認為用 **Euclidean** 的結果較好，C1 和 C2 的 improvement percentage 皆比較高，推測是 **Manhattan** 在高維資料（尤其此次作業為 58 維）的處理上不太理想，用 **Euclidean** 處理較合適。

- The Euclidean and Manhattan Distances for all pairs of centroids, with 2 initialization strategies.

● Manhattan – c1

Manhattan	1	2	3	4	5	6	7	8	9	10
1	0	2341.017	11929.3	651.187	496.332	947.743	770.737	1056.8	1260.511	737.714
2		0	9597.441	2778.946	2830.145	3280.359	3104.286	3388.983	2380.461	1605.27
3			0	12323.288	12421.263	12871.483	12695.554	12979.133	10775.939	11196.787
4				0	335.951	558.469	382.463	667.533	1653.826	1379.165
5					0	452.861	276.326	561.849	1755.106	1226.66
6						0	177.593	110.218	2205.307	1677.667
7							0	287.43	2028.902	1500.993
8								0	2314.667	1786.811
9									0	1006.368
10										0

● Euclidean – c1

Euclidean	1	2	3	4	5	6	7	8	9	10
1	0	2219.177	9948.044	528.7	413.365	827.719	681.035	917.127	832.147	729.056
2		0	7767.946	2734.05	2628.491	3044.478	2898.713	3133.46	1812.455	1491.357
3			0	10433.061	10361.367	10773.531	10626.489	10862.966	9340.275	9236.84
4				0	221.373	375.156	249.379	457.26	1156.583	1251.158
5					0	415.99	270.749	505.071	1171.964	1137.135
6						0	147.047	89.491	1529.464	1553.124
7							0	236.515	1391.55	1407.404
8								0	1613.556	1642.129

