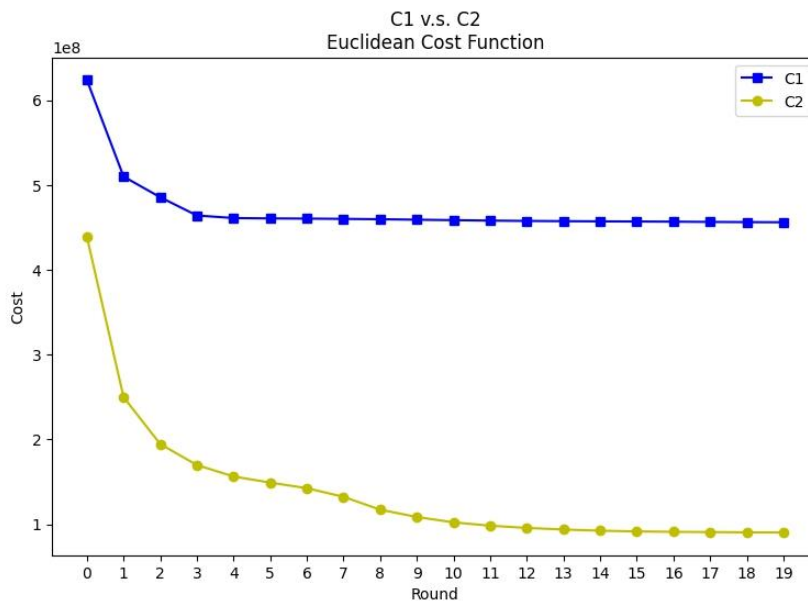


k-means on MapReduce

105072123 黃海茵

(a) Euclidean

	C1	C2
Round 1	6.236603e+08	4.387478e+08
Round 2	5.098629e+08	2.498039e+08
Round 3	4.854807e+08	1.944948e+08
Round 4	4.639970e+08	1.698048e+08
Round 5	4.609693e+08	1.562957e+08
Round 6	4.605378e+08	1.490942e+08
Round 7	4.603131e+08	1.425085e+08
Round 8	4.600035e+08	1.323039e+08
Round 9	4.595705e+08	1.171710e+08
Round 10	4.590211e+08	1.085474e+08
Round 11	4.584907e+08	1.022372e+08
Round 12	4.579442e+08	9.827802e+07
Round 13	4.575580e+08	9.563023e+07
Round 14	4.572901e+08	9.379331e+07
Round 15	4.570506e+08	9.237713e+07
Round 16	4.568922e+08	9.154161e+07
Round 17	4.567036e+08	9.104557e+07
Round 18	4.564042e+08	9.075224e+07
Round 19	4.561778e+08	9.047017e+07
Round 20	4.559869e+08	9.021642e+07



- Percentage improvement : C1 是 26.885% , C2 是 79.438% 。
C2 各點相距較遠，所以 data 在最初的距離會較 C1 來的大，因此 C2 的 percentage improvement 會較大。使用 Euclidean 時，是計算平方和開根號，也就是兩點間的距離，所以 centroid 相距較遠，cluster 分布會較平均，每個點到 centroid 的距離也會越小，所以 C2 的 cost 會小於 C1 。

- C1 Euclidean

[illegible]

- C1 Manhattan

[illegible]

- C2 Euclidean

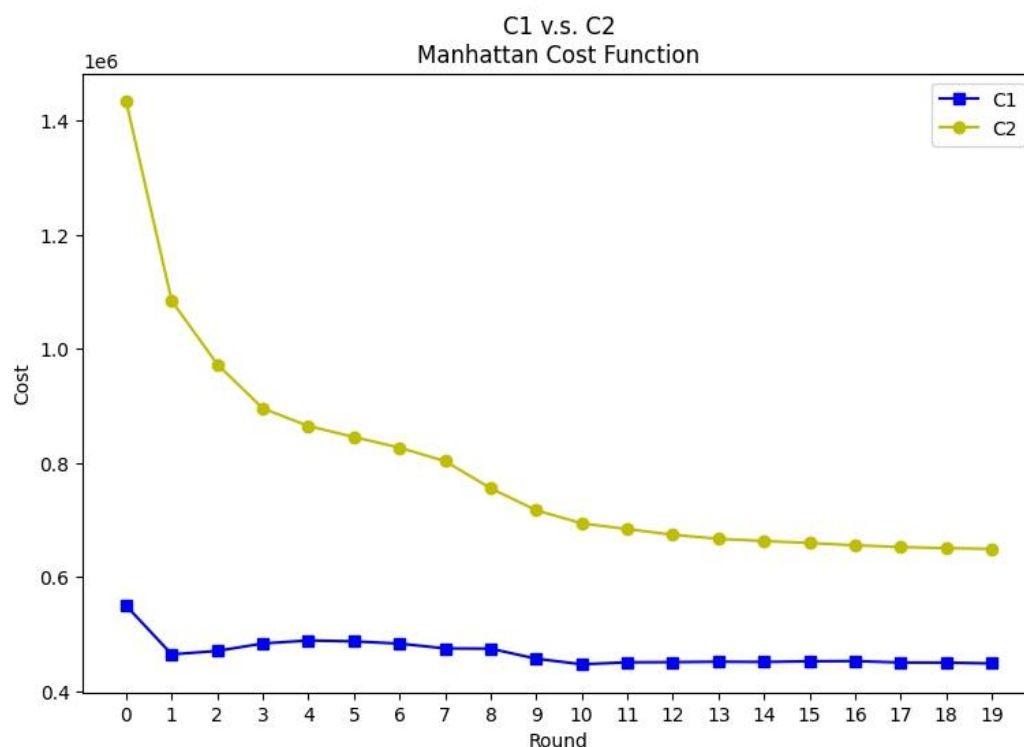
[illegible]

- C2 Manhattan

[illegible]

(b) Manhattan

	C1	C2
Round 1	550117.142000	1.433739e+06
Round 2	464869.275879	1.084489e+06
Round 3	470897.382277	9.734317e+05
Round 4	483914.409173	8.959346e+05
Round 5	489216.071003	8.651283e+05
Round 6	487629.668550	8.458466e+05
Round 7	483711.923214	8.272196e+05
Round 8	475330.773493	8.035903e+05
Round 9	474871.238846	7.560395e+05
Round 10	457232.920115	7.173329e+05
Round 11	447494.386197	6.945879e+05
Round 12	450915.012577	6.844445e+05
Round 13	451250.367073	6.745747e+05
Round 14	451974.595540	6.674095e+05
Round 15	451570.364070	6.635566e+05
Round 16	452739.011366	6.601628e+05
Round 17	453082.730287	6.560413e+05
Round 18	450583.670860	6.530368e+05
Round 19	450368.749317	6.511124e+05
Round 20	449011.363726	6.496890e+05



- Percentage improvement : C1 是 18.379% , C2 是 54.686% 。

同(a)小題，C2 各點相距較遠，所以 data 在最初的距離會較 C1 來的大，因此 C2 的 percentage improvement 會較大。使用 Manhattan 時，是計算每個維度距離的總和，而非兩點間的距離，所以不一定會比較小。

- C1 Euclidean

[illegible]

- C1 Manhattan

[illegible]

- C2 Euclidean

[illegible]

- C2 Manhattan

[illegible]