

MDA_HW3_KMeans

107062321 王劭元

在這次的作業中我主要共用到了 6 個 mapper，2 個 reducer

mapperE_dist：在這個 mapper 中，我是將每個點去與 centroid 算出他們之間的 Euclidean distance 後，以我希望的形式包好後回傳

mapperE_min：在這個 mapper 中，傳進來的資料會是 data 中的點，後面接著十個 centroid，我用迴圈掃過 10 個 centroid 後，會得到是哪個 centroid 距離最短，並且距離為多少，最後要回傳時，我除了回傳是哪個 centroid 之外，我還會回傳這個 centroid 的 index 以利於後面可以 sortByKey，讓最後 csv 檔的輸出可以順利一點。並且我在回傳時，我在最後的 KV pair 就會直接回傳這個點所產生的 cost 為多少，這樣最後我對 data 所有點的這個值用 api sum() 起來，就可以得到這次的總 cost 為多少。除此之外，由於發現 data 中有重複出現的點會影響到 cost，所以我還會去算若同樣的點後面還接了幾個，最後的總 cost 會去乘上點的數量。

mapperM_dist：在這個 mapper 中，我是將每個點去與 centroid 算出他們之間的 Manhattan distance 後，以我希望的形式包好後回傳

mapperM_min：在這個 mapper 中，與 mapperE_min 實作的方式幾乎一樣，我

只是為了要回傳不同的 `cost`。所以這邊只有在回傳的 `cost` 算法不一樣，其他都跟 `mapperE_min` 一樣。

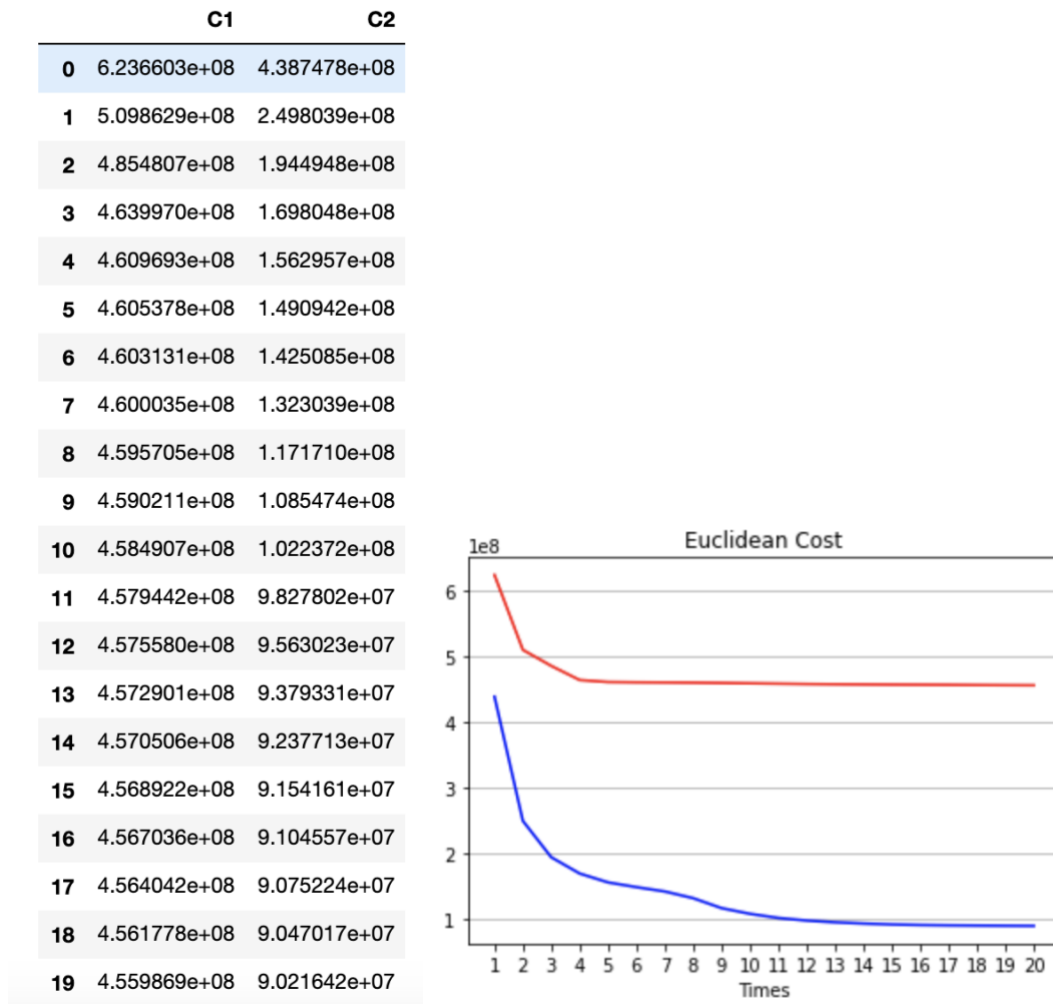
mapperNew_cen：在這個 `mapper` 中，傳進來的資料會是其中一個 `centroid` 接上，現在屬於這個 `centroid`，或者說新的這個 `cluster` 的點。所以要計算新的 `cluster` 的中心時，就把後面用迴圈全部掃過，並且把各 `dimension` 的值加起來平均後，即可得到新的 `centroid`。

mapperinit：在這個 `mapper` 中，我為了讓每次 `iteration` 完，紀錄 `centroid` 的 `rdd` 回復成我一開始的型態所用的。會回傳 `(0, centroid)`。

Reducer：兩個 `reducer` 的功用都蠻像的，都是為了讓我可以將資料串接在一起。一次是在每次 `iteration` 一開始，我的 `centroid` 和 `data` 這兩個 `rdd` 的 `kv pair` 都長成 `(0, centroid)` `(0, data)`，的形式，所以我先把他們 `join` 起來後，算完每一個點和 `centroid` 的距離後，為了要去計算 `min distance`，要把同一個點對不同的 `centroid` 的這些 `pair` 串起來，故我這裡用了一次 `reduceByKey(lambda a, b:a+b)`。第二次是在把每個 `centroid` 的與某個 `data point` 的 `cost` 算完後，為了要算出總 `cost` 要把所有屬於這個 `cluster` 得 `data point` 串在 `centroid` 的後面，於是我在這邊用了第二次的 `reduceByKey(lambda a, b:a+b)`。

(a)

(1) A plot of cost vs. iteration for 2 initialization strategies



(b) Percentage improvement values and your explanation

Change with c1 26.885383292517258

Change with c2 79.437750291599

由於 Euclidean 計算 cost 的方式為距離的平方，所以對於 cluster 距離遠近、離群值較為敏感，故在這種情況，使用 c2 事先刻意挑好讓最遠的點可以自己形成一種 cluster 的方式，最後得到結果以及下降的幅度也較好。

(1) Euclidean distance for all pairs of centroids, with c1 started

(2) Manhattan distance for all pairs of centroids, with c1 started

(3) Euclidean distance for all pairs of centroids, with c2 started

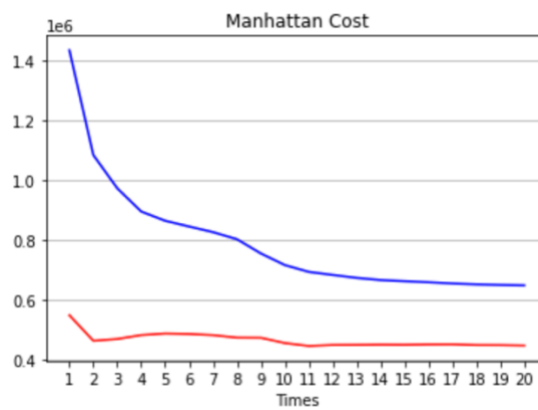
(4) Manhattan distance for all pairs of centroids, with c2 started

[illegible]

(b)

(1) A plot of cost vs. iteration for 2 initialization strategies

	C1	C2
0	550117.142000	1.433739e+06
1	464869.275879	1.084489e+06
2	470897.382277	9.734317e+05
3	483914.409173	8.959346e+05
4	489216.071003	8.651283e+05
5	487629.668550	8.458466e+05
6	483711.923214	8.272196e+05
7	475330.773493	8.035903e+05
8	474871.238846	7.560395e+05
9	457232.920115	7.173329e+05
10	447494.386197	6.945879e+05
11	450915.012577	6.844445e+05
12	451250.367073	6.745747e+05
13	451974.595540	6.674095e+05
14	451570.364070	6.635566e+05
15	452739.011366	6.601628e+05
16	453082.730287	6.560413e+05
17	450583.670860	6.530368e+05
18	450368.749317	6.511124e+05
19	449011.363726	6.496890e+05



(b) Percentage improvement values and your explanation

Change with c1: 18.378954327236887

Change with c2: 54.685694348134085

由於 Manhattan 在計算 cost 的方式相對來說比較不敏感，Manhattan 更希望你
能將集中很多點的大群即使距離很近，也可以分成不同 cluster，故在這種情
況，一開始還特地將距離拉遠來看，最後達到的效果沒有比較好。

(1) Euclidean distance for all pairs of centroids, with c1 started

(2) Manhattan distance for all pairs of centroids, with c1 started

(3) Euclidean distance for all pairs of centroids, with c2 started

(4) Manhattan distance for all pairs of centroids, with c2 started

[illegible]