

Homework 3 : K-Means

108062608 黃柏翰

Map Reduce Algorithm and Code

演算法根據老師第七章 p.34

- **1)** For each point, place it in the cluster whose current centroid it is nearest
- **2)** After all points are assigned, update the locations of centroids of the ***k*** clusters
- **3)** Reassign all points to their closest centroid
 - Sometimes moves points between clusters
- **Repeat 2 and 3 until convergence**

因為題目要求使用不同的距離量測方式，所以在距離計算上面會寫成不同的 function

a. Map

此 Mapper 的功能為做計算各 node 到各 centroid 的距離，並

assign 各 node 到各 centroid，為題目要求使用不同的距離量測方

式所以，會有兩種 mapper

程式碼

```
def Mapper_Euclidean(data):  
    # calculate_Euclidean_distance  
    input_data = data.map(calculate_Euclidean_distance)  
  
    # Assign to new centroid  
    key_value = Assign2new_centroid(input_data)  
  
    return key_value  
  
def Mapper_Manhattan(data):  
    # calculate_Euclidean_distance  
    input_data = data.map(calculate_Manhattan_distance)  
  
    # Assign to new centroid  
    key_value = Assign2new_centroid(input_data)  
  
    return key_value
```

所呼叫的副程式程式碼

```
def Assign2new_centroid(input_data):  
    # find the minimum distance to each centroid  
    # Assign to new cluster  
    # (cluster, node)  
    # value, key, then reduce by key & average, you will get new centroid  
    key_value = input_data.map(lambda x : (x[0][1].index(min(x[0][1])), x[0][0]))  
  
    return key_value
```

b. Reduce

此 Reducer 的功能為計算新的 centroid，因為找 centroid 跟距離量

測方式無關，所以只需一個 Reducer

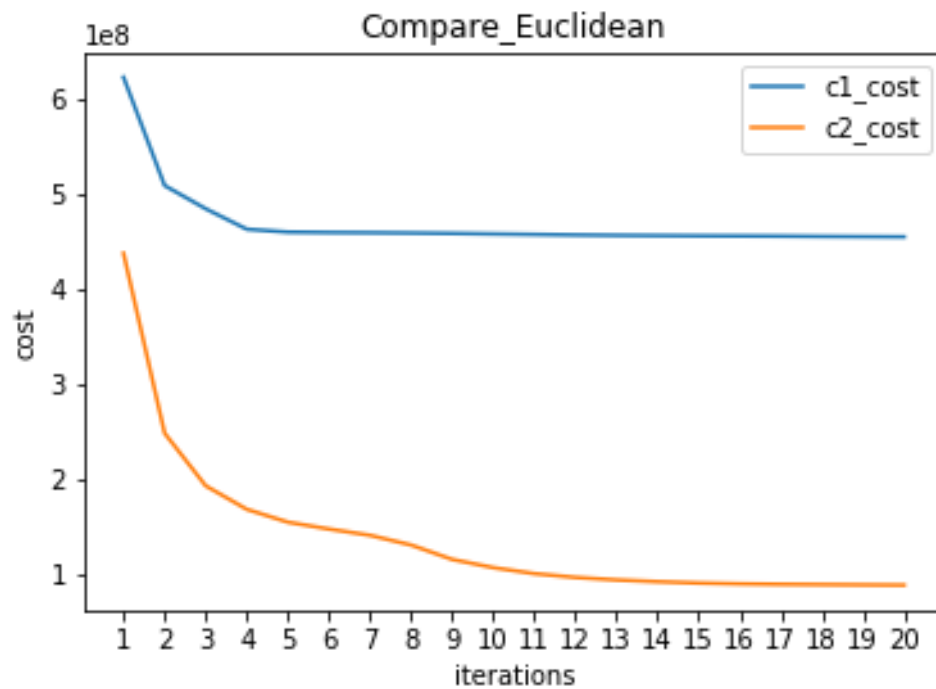
程式碼

```
def Reducer(key_value):  
    global countsByKey, centroid  
    # find_new_centroid  
    # find the number of value for each key  
    countsByKey = key_value.countByKey()  
    # accumulate same key value from same cluster  
    key_value = key_value.reduceByKey(lambda x,y : add(x,y))  
  
    key_value = key_value.map(find_new_centroid)  
  
    #update centroid  
    centroid = key_value.collect()  
  
    return centroid
```

Homework 回答

以 Euclidean distance 做 K-means

1. Performance 比較



2. After 10 iterations, about centroid distance

2.1 探討進步程度

Using Euclidean with c1 as initial centroid improves 0.26 %

Using Euclidean with c2 as initial centroid improves 0.75 %

2.2 c1, c2 performance 比較

c2 在整體 performance 上不管是改善 cost，或者 cost 的大小上都比 c1 來的更好，因為 kmeans 的初始化很重要，所以隨機的初始化的結果很可能比較差

2.3 Centroid 距離比較表格

使用 c1 下用 Euclidean distance

[illegible]

使用 c1 下用 Manhattan distance

[illegible]

使用 c2 下用 Euclidean distance

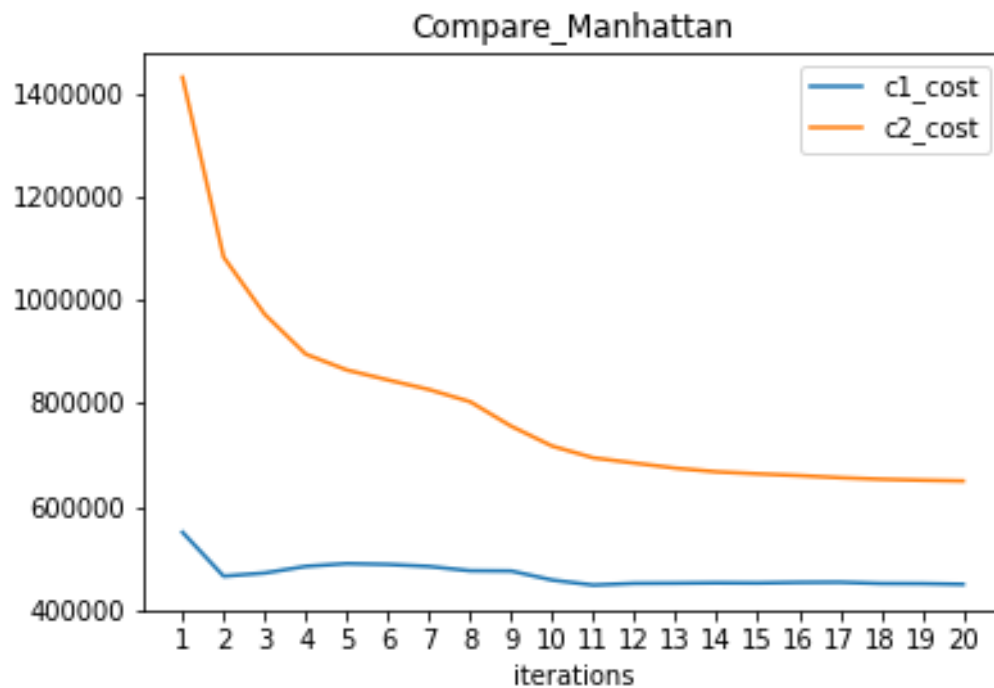
[illegible]

使用 c2 下用 Manhattan distance

[illegible]

以 Manhattan distance 做 K-means

1. Performance 比較



2. After 10 iterations, about centroid distance

2.1 探討進步程度

Using Manhattan with c1 as initial centroid improves 0.17 %

Using Manhattan with c2 as initial centroid improves 0.50 %

2.2 c1, c2 performance 比較

c1 performance 上(cost)表現比 c2 佳，進步上比 c2 差因為 kmeans

的初始化很重要，所以隨機的初始化的結果很可能使進步程度比

較差，但是因為 Manhattan 距離量測方法的考量為各維度的點差

距和，所以 cost 結果可能跟 Euclidean 不同

2.3 centroid 距離比較表格

使用 c1 下用 Euclidean distance

[illegible]

使用 c1 下用 Manhattan distance

[illegible]

使用 c2 下用 Euclidean distance

[illegible]

使用 c2 下用 Manhattan distance

[illegible]