

# Apriori + Hash Tree 程式碼簡介

```
1 file = 'data\\BreadBasket_list.csv'
2 file_freq = 'data\\Bread_freq_HashTree.csv'
3 file_rule = 'data\\Bread_rule_HashTree.csv'
4
5 with open(file,newline='') as file:
6     contents = csv.reader(file)
7     transactions = [row for row in contents if row]
8
9 ...
```

↑ 載入BreadBasket的csv檔 ( 已將原BreadBasket\_DMS檔轉換成每筆transaction的型態 )

```
1 def get_L1(_transactions, min_support):
2     candidate_items = {} #stores support for each item
3     for transaction in _transactions:
4         for item in transaction:
5             candidate_items[item] = candidate_items.get(item,0) + 1
6     L1_support = []
7     L1_list = []
8     for item in candidate_items:
9         if candidate_items[item] >= min_support:
10             L1_support.append([item],candidate_items[item]))
11             L1_list.append([item])
12     return L1_support, L1_list
```

↑ 給入 **transactions dataset** 及 **minimum support**，函式中建立**candidate\_items**字典，並儲存交易中每筆**item**對應的**support**，並篩選 **support** 值大於 **minium support** 的 **item**。

```
1 def get_k_subsets(datasets, length):
2     subsets = []
3     for itemset in datasets:
4         subsets.extend([sorted(itemset) for itemset in list(itertools.combi
5     subsets = [val for val in subsets]
6     return subsets
```

↑ 給定**transaction dataset**及**itemset**長度，以便後續將**transactions dataset**中取出長度為**k**的所有組合並插入建立好的**hash tree**中計算各**pattern**的**support**值

```

1  def ck_generator(Lk,k):
2      ck_set = []
3      #join Ck=(Lk-1 self join Lk-1)
4      lenlk = len(Lk)
5      for i in range(lenlk):
6          for j in range(i+1,lenlk):
7              L1 = list(Lk[i])[:k - 2] #該Lk倒數第一前的元素
8              L2 = list(Lk[j])[:k - 2]
9              if L1 == L2:
10                 ck_set.append(sorted(list(set(Lk[i]).union(set(Lk[j])))))
11         return ck_set

```

↑ ck\_generator函數用來將Lk產生所有Ck+1候選集，其中針對所有Lk，利用self join的方式，將兩個長度為k的itemset做聯集以產生一組Ck+1候選itemset。函數output為所有Ck+1的候選集，以利後續建立Hash Tree。

```

1  class HNode:
2
3      def __init__(self):
4          self.children = {}
5          self.isLeaf = True
6          self.bucket = {}

```

↑ hash tree的建立方式參考github其他作者的公開程式碼。建立hash tree節點的class，每個節點皆能儲存子節點、葉節點狀態、及input之itemset

```

1  class HTree:
2
3      def __init__(self, max_leaf_cnt, max_child_cnt):
4          self.root = HNode()
5          self.max_leaf_cnt = max_leaf_cnt
6          self.max_child_cnt = max_child_cnt
7          self.frequent_itemsets = []
8
9      def recur_insert(self, node, itemset, index, cnt):
10         def insert(self, itemset):
11             def add_support(self, itemset):
12                 def dfs(self, node, support_cnt):
13                     def get_frequent_itemsets(self, support_cnt):
14                         def hash(self, val):

```

(因程式碼內容太多，function內容省略)

↑ 建立hash tree的class，這個tree的attribute紀錄著其根節點，最多葉節點及子節點數量、和由這個tree產生的frequent itemsets。

hash tree分為三大部分：

1. 以候選集建立樹狀結構
  - 利用**insert function**以某候選**itemset**建立樹枝，並呼叫**recur\_insert**透過不斷疊帶的方式依照**itemset**中為每個**item**建立相對應的**node**，疊代的過程將會判斷某**node**是否已超過其**max\_leaf**的數量而往下增加子節點
  - 其中也會依照層級將**itemset**中第**k**個**item**利用**hash function**做**hash**後將其視為其子節點的**key**。
2. 利用transaction dataset中長度為k的所有組合將itemset插入已建立完成的hash tree中
  - 利用 **add\_support**的**function**將長度為k的**itemset**以**item**的順序先後插入**hash tree**當中，並依照其條件在節點的部分增加該**itemset**的**support**
3. 利用get\_frequent\_itemsets函數取得該hash tree中有滿足minimum support的itemset集合
  - 判斷該**node**是否為葉節點，並將葉節點中**itemset**以字典的方式儲存其**itemset**的**support**

```

1  def apriori(_transactions,min_support,max_leaf_cnt,max_child_cnt, freq_patterns, fr
2      k=2
3      L1_dict, L1_list = get_L1(_transactions, min_support)
4      ... generate_hash_tree(C2_candidates,max_leaf_cnt,max_child_cnt)
5      k_subsets = get_k_subsets(_transactions,k)
6      for subset in k_subsets:
7          h_tree.add_support(subset)
8      L2 = [items[0] for items in h_tree.get_frequent_itemsets(min_support)]
9      ...
10     while(len(freq_patterns[k-2])>0):
11         Ck_candidates = ck_generator(freq_patterns[k-2], k)
12         ...
13     return

```

( 因程式碼內容太多 部分已省略 )

↑ **apriori**這個**function**主要是統整利用**hash tree**建立**frequent patterns**的順序，首先會從**input**中的**transactions dataset**建立**L1**，再由**L1**利用**ck\_generator**函數建立**C2**，之後再利用**C2**建立**hash tree**後，從**transactions**找出所有長度為2的**itemset**後，利用**tree**中的**functions**找出**L2**，然後以此邏輯繼續疊代後即可找出所有**frequent itemsets**

```

1  def rule_generator(freq_itemset, _candidate_sets, init_num, min_confidence, rules_l
2      ...
3      if init_num == 1 : #第一層
4          _candidate_sets = ...
5          init_num += 1
6          rule_generator(freq_itemset...)
7      else :
8          pruned_subsets = ...
9          if pruned_subsets != []:
10             ...
11             else:
12                 init_num += 1
13                 rule_generator(freq_itemset,...)
14

```

( 部分程式碼已省略 )

↑ 針對每一個**frequent itemset**，先找出其**k-1**的所有組合，再利用**self join**的方式及課堂上講解的技巧，取左邊交集右邊聯集的方式，找出關聯組合，並計算每個組合的**confidence**，若低於**minimum**值則**prune**掉，再利用**pruned**完後的**itemset**繼續以相同邏輯疊代找出該組**itemset**所有符合條件的關聯集合。

```

1
2  min_sup = 10
3  max_leaf_cnt = 20
4  max_child_cnt = 20
5  apriori(_transactions...)
6
7  with open(file_freq, 'w') as file:
8      ...
9      for freq_pattern, support in freq_dict_sorted.items():
10         file_writer.writerow([set(freq_pattern), support])
11  '''
12  Rule Generation
13  '''
14  min_conf = 0.5
15  rules_list = []
16  for _freq_pattern in flattened_patterns:
17      rule_generator(freq_itemset...)
18  ...
19  final_rules = sorted(flattened_rules, key = lambda x: x[2], reverse = True)
20
21  with open(file_rule, 'w') as csv_file:
22      ...
23  print('writing done!')

```

(部分程式碼已省略)

↑ 最後僅分別依照設定條件呼叫**apriori**函數及**rule\_generator**函數得到**frequent itemsets**及**rules**，並將結果分別寫入以**file\_freq**與**file\_rule**為名的檔案