Vivek Khanolkar
ECE49595 - Data Mining
Assignment 2

Candidate list of 3 for gene:



Frequency lists for gene:



Code is below.

```python
# generate candidate set Ck
Ck = []
kLess = len(freq_sets[0]) - 1 #k = 3 then kLess = 1
for x in freq_sets:
    for y in freq_sets:
        # print("x: " , x , "y: " , y)
        if x != y:
            i = map(list, x)
            j = map(list, y)
            if sorted(i[0:kLess]) == sorted(j[0:kLess]) :
                # i = set(x)
                # j = set(y)
                merged = x | y
                # print("merged: ", merged)
                if(merged not in Ck):
                    Ck.append(merged)


# now prune the set
prune = []
for x in Ck:
    xNew = set(x)
    pp = combinations(xNew, kLess + 1)
    for i in pp:
        i = set(i)
        if i not in freq_sets:
            if x in Ck:
                # Ck.remove(x)  #this does not work for some reason (different from removing with prune at end) but how?
                prune.append(x)

for y in prune:
    Ck.remove(y)

Ck = list(map(frozenset, Ck))
# print("FreqLength: ", len(freq_sets))
# print("FREQ: ", freq_sets)
# print("CandidateLength: ", len(Ck))
# print("CANDIDATE: ", Ck)
return Ck
```

```python
def get_freq(dataset, candidates, min_support, verbose=False):
    """

    This function separates the candidates itemsets into frequent itemset and infrequent itemsets based on the min_support,
    and returns all candidate itemsets that meet a minimum support threshold.

    Parameters
    ----------
    dataset : list
        The dataset (a list of transactions) from which to generate candidate
        itemsets.

    candidates : frozenset
        The list of candidate itemsets.

    min_support : float
        The minimum support threshold.

    Returns
    -------
    freq_list : list
        The list of frequent itemsets.

    support_data : dict
        The support data for all candidate itemsets.
    """
    # print("DATASET: ", dataset)
    # print("CAN: ", candidates)
    # print("MIN: ", min_support)
    minSup = min_support * len(dataset) #need 50% of dataset to have item
    freq_list = []
    dic = {}
    for x in dataset:
        for y in candidates:
            if y not in dic:
                dic[y] = 0
            if (y.issubset(x)):
                dic[y] = dic[y] + 1

    for k in dic:
        if (dic[k] >= minSup):
            freq_list.append(k)

    support_data = dic

    return freq_list, support_data
```