

CSE 601: Data Mining and Bioinformatics

Project 1: Dimensionality Reduction & Association Analysis

Part 2: Association Analysis

Team Members:

- 1. Name:** Bhagutharivalan Natarajan Muthukkannu
UB Person No: 50314871
UBIT Name: bhagutha
- 2. Name:** Harish Kannan Venkataramanan
UB Person No: 50316999
UBIT Name: hvenkata
- 3. Name:** Praveen Mohan
UB Person No: 50321225
UBIT Name: pmohan2

Generating Frequent Itemset using Apriori Algorithm

Packages used:

1. *“import itertools”* - To iterate through each frequent itemsets to generate candidate rules.
2. *“import pandas as pd”* - To store and retrieve RULES | HEAD | BODY for obtaining template results.

Functions implemented:

- *freq_items_generation(candidate_items, tr_list, support, length)*

```
1 def freq_items_generation(candidate_items, tr_list, support, length):
2     frequent_dict = dict()
3     for candidate in candidate_items:
4         if length != 1:
5             candidate_set = set(candidate.split(','))
6         else:
7             candidate_set = {candidate}
8         count = 0
9         for row in tr_list:
10            if(candidate_set.issubset(row)):
11                count+=1
12            if (count/len(tr_list))*100 >= support:
13                frequent_dict[candidate] = count
14    return list(frequent_dict.keys())
```

The function `freq_items_generation(candidate_items, tr_list, support, length)` generates a list of frequent itemset.

Inputs:

1. `candidate_items` - A set of candidate frequent itemset for generating frequent itemsets for the corresponding threshold.
2. `tr_list` - Transaction list or list of all rows from the original data to find the support count for candidate frequent itemset.
3. `support` - Minimum Support Threshold in percentage.
4. `length` - Length of the candidate frequent itemset

Output:

1. `list(frequent_dict.keys())` - A list of frequent itemset.

- *candidate_set_generation(items, length)*

```

1 def candidate_set_generation(items, length):
2     candidate_items = set()
3     for a in range(len(items)-1):
4         items_a = set(items[a].split(','))
5         for b in range(a+1, len(items)):
6             items_b = set(items[b].split(','))
7             if (len(items_a | items_b) == length):
8                 temp = sorted(items_a | items_b)
9                 temp_list = ','.join(temp)
10                candidate_items.add(temp_list)
11    return candidate_items

```

The function `candidate_set_generation(items, length)` generates a set of candidate frequent itemset.

Inputs:

1. `items` - A set of frequent itemset for generating next length candidate frequent itemset.
2. `length` - Length of the itemset.

Output:

1. `candidate_items` - A set of candidate frequent itemset.

- *open_file(filename)*

```

1 def open_file(filename):
2     file = open(filename, "r")
3     candidate_items_l1 = set()
4     tr_list = []
5     for line in file:
6         row = line.strip("\n").split("\t")
7         for i in range(len(row)):
8             if i != len(row)-1:
9                 row[i] = "G"+str(i+1)+" "+row[i].upper()
10                candidate_items_l1.add(row[i])
11            tr_list.append(set(row))
12    return(candidate_items_l1, tr_list)

```

The function `open_file(filename)` generates a list of candidate frequent itemset of length 1 and a transaction list or list of all rows from the original data to find the support count for candidate frequent itemset.

Inputs:

1. `filename` - Name of the original gene data file.

Output:

1. `candidate_items_l1` - A list of candidate frequent itemset of length 1.
2. `tr_list` - Transaction list or list of all rows from the original data to find the support count for candidate frequent itemset.

- *def apriori_imp_template(filename, support)*

```

1 def apriori_imp_template(filename, support):
2     candidate_items_l1, tr_list = open_file(filename)
3     length = 1
4     candidate_items = candidate_items_l1
5     freq_items = freq_items_generation(candidate_items, tr_list, support, length)
6     ans = [len(freq_items)]
7     while True:
8         length += 1
9         candidate_items = candidate_set_generation(freq_items, length)
10        freq_items = freq_items_generation(candidate_items, tr_list, support, length)
11        if len(freq_items) == 0:
12            break
13        else:
14            ans.append(len(freq_items))
15    print("Support is set to be "+ str(support)+"%")
16    for i in range(len(ans)):
17        print("number of length-"+str(i+1)+ " frequent itemsets: "+str(ans[i]))
18    print("number of all lengths frequent itemsets: "+str(sum(ans))+"\n\n")

```

The function `apriori_imp_template(filename, support)` generates the template for part 1 of Apriori Algorithm in generating the frequent itemset for the given support.

Inputs:

1. `filename` - Name of the original gene data file.
2. `support` - Minimum Support Threshold in percentage.

Output:

1. The template for part 1 of Apriori Algorithm in generating the frequent itemsets for the given support.

- *apriori_imp_result(filename, support)*

```

1 def apriori_imp_result(filename, support):
2     candidate_items_l1, tr_list = open_file(filename)
3     length = 1
4     candidate_items = candidate_items_l1
5     freq_items = freq_items_generation(candidate_items, tr_list, support, length)
6     result = set(freq_items)
7     while True:
8         length += 1
9         candidate_items = candidate_set_generation(freq_items, length)
10        freq_items = freq_items_generation(candidate_items, tr_list, support, length)
11        if len(freq_items) == 0:
12            break
13        else:
14            result = result | set(freq_items)
15    return result, tr_list

```

Results for given support values [30%, 40%, 50%, 60%, 70%]

```
1 support_values = [30, 40, 50, 60, 70]
2 for val in support_values:
3     apriori_imp_template("association-rule-test-data.txt", val)
```

Support is set to be 30%
number of length-1 frequent itemsets: 196
number of length-2 frequent itemsets: 5340
number of length-3 frequent itemsets: 5287
number of length-4 frequent itemsets: 1518
number of length-5 frequent itemsets: 438
number of length-6 frequent itemsets: 88
number of length-7 frequent itemsets: 11
number of length-8 frequent itemsets: 1
number of all lengths frequent itemsets: 12879

Support is set to be 40%
number of length-1 frequent itemsets: 167
number of length-2 frequent itemsets: 753
number of length-3 frequent itemsets: 149
number of length-4 frequent itemsets: 7
number of length-5 frequent itemsets: 1
number of all lengths frequent itemsets: 1077

Support is set to be 50%
number of length-1 frequent itemsets: 109
number of length-2 frequent itemsets: 63
number of length-3 frequent itemsets: 2
number of all lengths frequent itemsets: 174

Support is set to be 60%
number of length-1 frequent itemsets: 34
number of length-2 frequent itemsets: 2
number of all lengths frequent itemsets: 36

Support is set to be 70%
number of length-1 frequent itemsets: 7
number of all lengths frequent itemsets: 7

Generating Association Rules from frequent itemset

- *freq_count(freq_itemset, tr_list)*

```
1 def freq_count(freq_itemset, tr_list):
2     freq_itemset = set(freq_itemset)
3     count = 0
4     for row in tr_list:
5         if(freq_itemset.issubset(row)):
6             count+=1
7     return count
```

The function `freq_count(freq_itemset, tr_list)` generates the count for the given itemset in the transaction database.

Inputs:

1. `freq_itemset` - A frequent itemset for the given support threshold.
2. `tr_list` - Transaction list or list of all rows from the original data to find the support count for candidate frequent itemset.

Output:

1. `count` = The count for the given frequent itemset in the `tr_list`.

- `association_rules(filename, support, confidence)`

```

1 def association_rules(filename, support, confidence):
2     import itertools
3     import pandas as pd
4     rules = set()
5     rules_list = []
6     result, tr_list = apriori_imp_result(filename, support)
7     result = list(result)
8     for i in range(len(result)):
9         result[i] = result[i].split(",")
10    result = list(sorted(result, key = len, reverse=True))
11    for freq_itemset in result:
12        if len(freq_itemset) != 1:
13            rule_cnt = freq_count(freq_itemset, tr_list)
14            for i in range(len(freq_itemset) - 1, 0, -1):
15                head_list = list(itertools.combinations(freq_itemset, i))
16                for head in head_list:
17                    head_cnt = freq_count(head, tr_list)
18                    rule_confidence = round((rule_cnt/head_cnt)*100, 2)#####
19                    if rule_confidence >= confidence:
20                        freq_itemset = set(freq_itemset)
21                        head = set(head)
22                        body = freq_itemset ^ head
23                        rule_set = set(head|body)
24                        rule = ",".join(head) + " -> " + ",".join(body)
25                        if rule not in rules:
26                            rules.add(rule)
27                            rules_list.append([rule, rule_set, head, body, rule_confidence])
28    rules_df = pd.DataFrame(rules_list, columns= ["RULE_SET", "RULE", "HEAD", "BODY", "Confidence"])
29    return rules_df

```

Generating Association Rules with a minimum Support Threshold of 50% and Confidence Threshold of 70%.

```

1 rules = association_rules("association-rule-test-data.txt", 50, 70)
2 rules

```

	RULE_SET	RULE	HEAD	BODY	Confidence
0	G72_UP,G59_UP -> G82_DOWN	{G72_UP, G82_DOWN, G59_UP}	{G72_UP, G59_UP}	{G82_DOWN}	83.87
1	G82_DOWN,G59_UP -> G72_UP	{G72_UP, G82_DOWN, G59_UP}	{G82_DOWN, G59_UP}	{G72_UP}	91.23
2	G82_DOWN,G72_UP -> G59_UP	{G82_DOWN, G59_UP, G72_UP}	{G82_DOWN, G72_UP}	{G59_UP}	89.66
3	G72_UP -> G82_DOWN,G59_UP	{G82_DOWN, G59_UP, G72_UP}	{G72_UP}	{G82_DOWN, G59_UP}	70.27
4	G82_DOWN -> G72_UP,G59_UP	{G72_UP, G59_UP, G82_DOWN}	{G82_DOWN}	{G72_UP, G59_UP}	76.47
...
112	G88_DOWN -> G24_DOWN	{G88_DOWN, G24_DOWN}	{G88_DOWN}	{G24_DOWN}	70.42
113	G59_UP -> G88_DOWN	{G88_DOWN, G59_UP}	{G59_UP}	{G88_DOWN}	72.37
114	G88_DOWN -> G59_UP	{G88_DOWN, G59_UP}	{G88_DOWN}	{G59_UP}	77.46
115	G88_DOWN -> G38_DOWN	{G88_DOWN, G38_DOWN}	{G88_DOWN}	{G38_DOWN}	70.42
116	G2_DOWN -> G38_DOWN	{G38_DOWN, G2_DOWN}	{G2_DOWN}	{G38_DOWN}	75.76

117 rows x 5 columns

- *asso_rule_template1(a, b, c)*

```

1 def asso_rule_template1(a, b, c):
2     c = set(",".join(c).upper().split(','))
3     result = []
4     for i in range(len(rules)):
5         if b == 'ANY' and len(c & rules.iloc[i][a]) > 0:
6             result.append(rules.iloc[i]["RULE_SET"])
7         elif b == 'NONE' and len(c & rules.iloc[i][a]) == 0:
8             result.append(rules.iloc[i]["RULE_SET"])
9         elif b == 1 and len(c & rules.iloc[i][a]) == 1:
10            result.append(rules.iloc[i]["RULE_SET"])
11    return result, len(result)

```

The function `asso_rule_template1(a, b, c)` generates the results for template 1 for the given query.

Inputs:

1. a - "RULE" | "HEAD" | "BODY"
2. b - "ANY" | "NONE" | 1
3. c - ["Gene", ...]

Output:

1. result - A list of rules for the given query.
2. len(result) - Total number of rules generated for the given query.

Template 1 results

(support = 50%, confidence = 70%)

```

1 template1_query = [{"RULE", "ANY", ['G59_UP']}, {"RULE", "NONE", ['G59_UP']}, {"RULE", 1, ['G59_UP', 'G10_
2 for val in template1_query:
3     result, count = asso_rule_template1(val[0], val[1], val[2])
4     print("The total number of rules generated for the template 1 query "+str(val[0])+" "+str(val[1])+"

```

The total number of rules generated for the template 1 query 'RULE, ANY, ['G59_UP']': 26
The total number of rules generated for the template 1 query 'RULE, NONE, ['G59_UP']': 91
The total number of rules generated for the template 1 query 'RULE, 1, ['G59_UP', 'G10_Down']': 39
The total number of rules generated for the template 1 query 'HEAD, ANY, ['G59_UP']': 9
The total number of rules generated for the template 1 query 'HEAD, NONE, ['G59_UP']': 108
The total number of rules generated for the template 1 query 'HEAD, 1, ['G59_UP', 'G10_Down']': 17
The total number of rules generated for the template 1 query 'BODY, ANY, ['G59_UP']': 17
The total number of rules generated for the template 1 query 'BODY, NONE, ['G59_UP']': 100
The total number of rules generated for the template 1 query 'BODY, 1, ['G59_UP', 'G10_Down']': 24

- *asso_rule_template2(a, b)*

```

1 def asso_rule_template2(a, b):
2     result = []
3     count = 0
4     for i in range(len(rules)):
5         if len(rules.iloc[i][a]) >= b:
6             result.append(rules.iloc[i]["RULE_SET"])
7             count += 1
8     return result, count

```

The function `asso_rule_template2(a, b)` generates the results for template 2 for the given query.

Inputs:

1. `a` - "RULE" | "HEAD" | "BODY"
2. `b` - integer (length)

Output:

1. `result` - A list of rules for the given query.
2. `count` - Total number of rules generated for the given query.

Template 2 results

```

1 template2_query = [{"RULE", 3}, {"HEAD", 2}, {"BODY", 1}]
2 for val in template2_query:
3     result, count = asso_rule_template2(val[0], val[1])
4     print("The total number of rules generated for the template 2 query " + val[0] + ", " + str(val[1]) + ": " + str(count))

```

The total number of rules generated for the template 2 query 'RULE, 3': 9
 The total number of rules generated for the template 2 query 'HEAD, 2': 6
 The total number of rules generated for the template 2 query 'BODY, 1': 117

- *temp_operator(string)*

```

1 def temp_operator(string):
2     if len(string) == 4:
3         string = string.split("or")
4         string.append("or")
5     elif len(string) == 5:
6         string = string.split("and")
7         string.append("and")
8     return string

```

The function `temp_operator(string)` splits the first input for template 3 into respective template value and the corresponding operator.

Inputs:

1. `string` - "1or1", "2or2", "1or2", "1and1", "2and2", "1and2".

Output:

1. `list` - [template number, template number, operator]

- *asso_rule_template3(a,b,c,d,e,f=None,g=None)*

```

1 def asso_rule_template3(a,b,c,d,e,f=None,g=None):
2     a = temp_operator(a)
3     if a[0] == '1' and a[1] == '1':
4         result1, count1 = asso_rule_template1(b, c, d)
5         result2, count2 = asso_rule_template1(e, f, g)
6     elif a[0] == '2' and a[1] == '2':
7         result1, count1 = asso_rule_template2(b, c)
8         result2, count2 = asso_rule_template2(d, e)
9     elif a[0] == '1' and a[1] == '2':
10        result1, count1 = asso_rule_template1(b, c, d)
11        result2, count2 = asso_rule_template2(e, f)
12    elif a[0] == '2' and a[1] == '1':
13        result1, count1 = asso_rule_template2(b, c)
14        result2, count2 = asso_rule_template1(d, e, f)
15
16    if a[2] == "and":
17        final = set(result1) & set(result2)
18        final = set(final)
19        count = len(final)
20        return final, count
21    if a[2] == "or":
22        final = result1 + result2
23        final = set(final)
24        count = len(final)
25        return final, count

```

The function `asso_rule_template2(a, b)` generates the results for template 3 for the given query.

Inputs:

1. a - "1or1", "2or2", "1or2", "1and1", "2and2", "1and2".
2. b - "RULE" | "HEAD" | "BODY"
3. c - "ANY" | "NONE" | 1 (or) integer
4. d - ["Gene", ...] (or) "RULE" | "HEAD" | "BODY"
5. e - "RULE" | "HEAD" | "BODY" (or) "ANY" | "NONE" | 1
6. f - "ANY" | "NONE" | 1 (or) ["Gene", ...] (or) None
7. g - ["Gene", ...] (or) None

Output:

1. final - A set of rules for the given query.
2. count - Total number of rules generated for the given query.

Template 3 results

```

1 template3_query = [['1or1', 'HEAD', 'ANY', ['G10_Down'], 'BODY', 1, ['G59_UP']], ['1and1', 'HEAD', 'ANY',
2 for val in template3_query:
3     if len(val) == 7:
4         result, count = asso_rule_template3(val[0], val[1], val[2], val[3], val[4], val[5], val[6])
5         print("The total number of rules generated for the template 3 query "+str(val)+"": "+ str(count)+"\
6     elif len(val) == 6:
7         result, count = asso_rule_template3(val[0], val[1], val[2], val[3], val[4], val[5])
8         print("The total number of rules generated for the template 3 query "+str(val)+"": "+ str(count)+"\
9     elif len(val) == 5:
10        result, count = asso_rule_template3(val[0], val[1], val[2], val[3], val[4])
11        print("The total number of rules generated for the template 3 query "+str(val)+"": "+ str(count)+"\
12
13

```

The total number of rules generated for the template 3 query ['1or1', 'HEAD', 'ANY', ['G10_Down'], 'BODY', 1, ['G59_UP']]: 24

The total number of rules generated for the template 3 query ['1and1', 'HEAD', 'ANY', ['G10_Down'], 'BODY', 1, ['G59_UP']]: 1

The total number of rules generated for the template 3 query ['1or2', 'HEAD', 'ANY', ['G10_Down'], 'BODY', 2] : 11

The total number of rules generated for the template 3 query ['1and2', 'HEAD', 'ANY', ['G10_Down'], 'BODY', 2] : 0

The total number of rules generated for the template 3 query ['2or2', 'HEAD', 1, 'BODY', 2]: 117

The total number of rules generated for the template 3 query ['2and2', 'HEAD', 1, 'BODY', 2]: 3