

Frequent Pattern Mining: Comparative Analysis

PROJECT REPORT

Mini Project of Data Mining (DSCC 440-2)

MS Data Science

By
Aradhya Mathur

(Implementation project) Using a programming language that you are familiar with, such as C++ or Java, implement three frequent itemset mining algorithms introduced in this chapter: (1) Apriori [AS94b], (2) FP-growth [HPY00], and (3) Eclat [Zak00] (mining using the vertical data format). Compare the performance of each algorithm with various kinds of large data sets. Write a report to analyze the situations (e.g., data size, data distribution, minimal support threshold setting, and pattern density) where one algorithm may perform better than the others, and state why.

Implementation project #1 (not HW#4)

6.7 (1) & (2)*, plus one improvement of your choice for 6.7(1)
Due 10/06

Use the UCI Adult Census Dataset to test your code

<http://archive.ics.uci.edu/ml/datasets/Adult>

Important note: you can use the open-source code as a reference, but you should implement the algorithms independently. The point of the assignment is for you to know how the algorithms are implemented, not just how to run them. It would be easy to detect the latter, e.g. if more than one of you uses the same code.



UNIVERSITY of
ROCHESTER

Fall Semester 2022

INDEX

Sr. No.	Section	Page No
1.	Abstract	3
2.	Problem Statement	3
3.	Objectives	3
4.	Dataset Description	4
5.	Pre-Processing	4-6
6.	Algorithms	
6.1	Apriori	7-8
6.2	Improved Apriori	9-10
6.3	FP Growth	11-14
7.	Library Used	15
8.	Results	15-19
9.	Conclusions	20
10.	References	21

ABSTRACT

Frequent pattern mining searches for recurrent relationships in a certain data set. In frequent mining typically the interesting associations and correlations between itemsets in transactional and relational databases are determined.

Support: Support of 10% means that 10% of all the transactions under analysis show that mobile and mobile cover are purchased together.

Confidence: A confidence of 60% means that 60% of the customers who purchased a mobile and mobile cover bought screen guard.

$\text{Support}(A \rightarrow B) = \text{Sup_count}(A \cup B)$

$\text{Confidence}(A \rightarrow B) = \text{Sup_count}(A \cup B) / \text{Sup_count}(A)$

A strong rule always satisfies both minimum support and minimum confidence.

$\text{Sup_count}(A)$: Total transactions in which A appears.

Closed Itemset: An itemset in which none of its direct supersets have support count same as itemset is called closed itemset.

K- Itemset: The itemset in which there are K items.

To conclude it can be said that an itemset is frequent if its support count is greater than minimum support count set by users or domain experts.

PROBLEM STATEMENT:

(Implementation project) Using a programming language that you are familiar with, such as Python, implement (1) Apriori, (2) FP-growth, and (3) Improved Apriori. Compare the performance of each algorithm with UCI Adult dataset. Analyse why one algorithm may perform better than the others.

OBJECTIVES

Implementing Apriori, Improved Apriori and FP Growth algorithms. Comparing performance of each other based on time taken to implement algorithms at different minimum support.

DATASET DESCRIPTION

Adult dataset is census income dataset which based on census data predicts whether income exceeds \$50K per year.

In this dataset there are 32561 rows and 15 columns.

Attributes in the dataset are:

- 1) Age: continuous values.
- 2) Workclass: Self-emp-not-inc, Federal-gov, Without-pay, etc.
- 3) fnlwgt: continuous values.
- 4) Education: Bachelors, HS-grad, Assoc-acdm, etc.
- 5) Education-num: continuous values.
- 6) Marital-status: Married-civ-spouse, Separated, Married-AF-spouse, etc.
- 7) Occupation: Exec-managerial, Machine-op-inspct, Armed-Forces, etc.
- 8) Relationship: Husband, Unmarried, etc.
- 9) Race: White, Amer-Indian-Eskimo, Black, etc.
- 10) Sex: Female or Male.
- 11) Capital-gain: continuous values.
- 12) Capital-loss: continuous values.
- 13) Hours-per-week: continuous values.
- 14) Native-country: United-States, India, China, France, Holand-Netherlands etc.
- 15) Income: >50K or <=50K.

In this data we can determine the pattern of which kind of people lie in income range >50K and <=50K. For example, after mining we obtain a frequent pattern {Female, United-States, Full-Time, >50K} pattern, we can say that a person who is a Female, lives in United-States and works Full-Time earns >50K.

There are some continuous attributes and rest are categorical. Our primary job is to convert continuous attributes into categorical.

PRE-PROCESSING:

Step 1) Reading data and adding column names

```
df = pd.read_csv('adult.data', sep=",", header = None , na_values = "?")
```

In this we observed there was no header.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	40	United-States	<=50K
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<=50K
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<=50K
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	40	United-States	<=50K
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<=50K

Fig 1: Raw Data

In order to add headers (column names), the following process was done:

```
df.columns = ['age', 'workclass', 'fnlwgt', 'education_num', 'education', 'marital_status',
'occupation', 'relationship', 'race', 'sex', 'capital_gain',
'capital_loss', 'hours_per_week', 'native_country', 'class']
```

	age	workclass	fnlwgt	education	education_num	marital_status	occupation	relationship	race	sex	capital_gain	capital_loss	hours_per_week	native
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	40	Un
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	Un
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	Un
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	40	Un
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	

Fig 2: Added Columns

Step 2) Converting continuous data into categorical data

We use binning method to convert continuous data into categorical data

```
df['Age'] = pd.cut(x=df['age'], bins=[0, 18, 30, 50, 100], labels=['Underage', 'Young',
'Adult', 'Elderly'])
```

Bins	Labels
0-18	Underage
18-30	Young
30-50	Adult
50-100	Elderly

Similarly, for hours_per_week

```
df['Hours_per_Week'] = pd.cut(x=df['hours_per_week'], bins=[0, 20, 40, 100],
labels=['Part-Time', 'Full-Time', 'Overtime'])
```

Bins	Labels
0-20	Part-Time
20-40	Full-Time
40-100	Overtime

Age	Hours_per_Week
Adult	Full-Time
Adult	Part-Time
Adult	Full-Time
Elderly	Full-Time
Young	Full-Time

Fig 3: New Age and Hours_per_Week Column

This is what new Age and Hours_per_Week columns look.

Step 3) Dropping continuous age and hours_per_week column as we created categorical ones.

Dropped fnlwgt and education number because they were random and due to inability to convert into categorical.

Dropped capital_gain and capital_loss as they have large number of 0 values.

```
df = df.drop(['age', 'fnlwgt', 'education_num', 'hours_per_week', 'capital_loss', 'capital_gain'],
axis = 1)
```

Before dropping

education_num	marital_status	occupation	relationship	race	sex	capital_gain	capital_loss	hours_per_week	native_country	class	Age	Hours_per_Week
13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	40	United-States	<=50K	Adult	Full-Time
13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	0	13	United-States	<=50K	Adult	Part-Time
9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<=50K	Adult	Full-Time
7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	0	0	40	United-States	<=50K	Elderly	Full-Time
13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<=50K	Young	Full-Time

Fig 4: Before Dropping Columns

After dropping

```
df = df.drop(['age', 'fnlwgt', 'education_num', 'hours_per_week', 'capital_loss', 'capital_gain'], axis = 1)
df.head()
```

	workclass	education	marital_status	occupation	relationship	race	sex	native_country	class	Age	Hours_per_Week
0	State-gov	Bachelors	Never-married	Adm-clerical	Not-in-family	White	Male	United-States	<=50K	Adult	Full-Time
1	Self-emp-not-inc	Bachelors	Married-civ-spouse	Exec-managerial	Husband	White	Male	United-States	<=50K	Adult	Part-Time
2	Private	HS-grad	Divorced	Handlers-cleaners	Not-in-family	White	Male	United-States	<=50K	Adult	Full-Time
3	Private	11th	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male	United-States	<=50K	Elderly	Full-Time
4	Private	Bachelors	Married-civ-spouse	Prof-specialty	Wife	Black	Female	Cuba	<=50K	Young	Full-Time

Fig 5: After Dropping Columns

Step 4) Removing whitespace

```
df = df.applymap(lambda space: space.strip() if type(space) is str else space)
```

ALGORITHMS

Apriori Algorithm:

```
#We will start the timer for Apriori Algorithm here
start_time = timeit.default_timer()
# Earlier we cleaned the dataframe. Now we will be using dataframe.values.tolist() to load the dataset.
def load_df():
    return df.values.tolist()
#First step is generating Candidate1.
def gen_cand1(itemset):
    CANDIDATE1 = []
    for i in itemset:
        for j in i:
            if not [j] in CANDIDATE1:
                CANDIDATE1.append([j])
    return list(map(frozenset, CANDIDATE1))
#Set in unhashable therefore we use frozenset which is nothing but immutable version of python set object.
itemset = load_df()
CANDIDATE1 = gen_cand1(itemset)
# Scanning Database
def database_scan(Db, Ck, min_sup):
    sup_count = {}
    sup_data = {}
    r_list = []
    Db_length = len(Db)
    for t in Db:
        for i in Ck:
            if i.issubset(t):
                if not i in sup_count: sup_count[i]=1
                else: sup_count[i] += 1

    total_items = int(Db_length)
    for key in sup_count:
        support = sup_count[key]/total_items
        if support >= min_sup:
            r_list.insert(0,key)
        sup_data[key] = support
    return r_list, sup_data

#In this step we use our knowledge of support count.
#If support of an item is greater than min_support, we insert that item and store it.

# Generating Apriori
def generate_apriori(L_k, k):
    Ck = []
    for l in range(len(L_k)):
        for i in range(l+1, len(L_k)):
```

```

        l1 = list(L_k[l])[:k-2]
        l2 = list(L_k[i])[:k-2]
        l1.sort()
        l2.sort()
        if l1==l2:
            Ck.append(L_k[l] | L_k[i])
    return Ck
# Apriori function in which we specify min_sup and obtain frequent itemsets corresponding to it.
def apriori(itemset, min_sup = 0.12):
    CANDIDATE1 = gen_cand1(itemset)
    D = list(map(set, itemset))
    L1, sup_data = database_scan(D, CANDIDATE1, min_sup)
    L = [L1]
    k = 2
    while (len(L[k-2]) > 0):
        Ck = generate_apriori(L[k-2], k)
        L_k, sup = database_scan(D, Ck, min_sup)
        sup_data.update(sup)
        L.append(L_k)
        k += 1
    return L, sup_data
L,sdata = apriori(itemset)
new_list = []
for index in range(len(L)):
    print("\nL{} ".format(index))
    print("Number of patterns={} \n".format(len(L[index])))
    apriori_freq_pattern = [list(x) for x in L[index]]
    print(apriori_freq_pattern)

```

Output:

For min_sup = 0.4

```

L0
Number of patterns=9

[['Private'], ['Husband'], ['Married-civ-spouse'], ['Full-Time'], ['Adult'], ['<=50K'], ['United-States'], ['Male'], ['White']]

L1
Number of patterns=21

[['White', 'Private'], ['Male', 'Private'], ['United-States', 'Private'], ['Private', '<=50K'], ['Full-Time', 'Private'], ['White', 'Married-civ-spouse'], ['Male', 'Married-civ-spouse'], ['United-States', 'Married-civ-spouse'], ['Husband', 'Male'], ['Husband', 'Married-civ-spouse'], ['Male', 'White'], ['United-States', 'White'], ['Male', 'United-States'], ['White', '<=50K'], ['Male', '<=50K'], ['United-States', '<=50K'], ['White', 'Adult'], ['United-States', 'Adult'], ['Full-Time', 'White'], ['Full-Time', 'United-States'], ['Full-Time', '<=50K']]

L2
Number of patterns=13

[['United-States', 'Private', '<=50K'], ['United-States', 'White', 'Private'], ['Male', 'United-States', 'Private'], ['Male', 'White', 'Private'], ['White', 'Private', '<=50K'], ['Husband', 'Male', 'Married-civ-spouse'], ['Full-Time', 'United-States', '<=50K'], ['Full-Time', 'White', '<=50K'], ['Full-Time', 'United-States', 'White'], ['Male', 'United-States', '<=50K'], ['United-States', 'White', '<=50K'], ['Male', 'White', '<=50K'], ['Male', 'United-States', 'White']]

L3
Number of patterns=1

[['United-States', 'White', 'Private', '<=50K']]

L4
Number of patterns=0

[]

```

Fig 6: Output for min_sup = 0.4

Improved Apriori Algorithm

```
#Taking a random sample containing 60% of original data
df = df.sample(frac=0.6)
df.describe()
#We will start the timer for Apriori Algorithm here
start_time = timeit.default_timer()
# Earlier we cleaned the dataframe. Now we will be using dataframe.values.tolist() to load the dataset.
def load_df():
    return df.values.tolist()
#First step is generating Candidate1.
def gen_cand1(itemset):
    CANDIDATE1 = []
    for i in itemset:
        for j in i:
            if not [j] in CANDIDATE1:
                CANDIDATE1.append([j])
    return list(map(frozenset, CANDIDATE1))
#Set in unhashable therefore we use frozenset which is nothing but immutable version of python set object.

itemset = load_df()
CANDIDATE1 = gen_cand1(itemset)

# Scanning Database
def database_scan(Db, Ck, min_sup):
    sup_count = {}
    sup_data = {}
    r_list = []
    Db_length = len(Db)
    for t in Db:
        for i in Ck:
            if i.issubset(t):
                if not i in sup_count: sup_count[i]=1
                else: sup_count[i] += 1

    total_items = int(Db_length)
    for key in sup_count:
        support = sup_count[key]/total_items
        if support >= min_sup:
            r_list.insert(0,key)
            sup_data[key] = support
    return r_list, sup_data

#In this step we use our knowledge of support count.
#If support of an item is greater than min_support, we insert that item and store it.

# Generating Apriori
def generate_apriori(L_k, k):
```

```

Ck = []
for l in range(len(L_k)):
    for i in range(l+1, len(L_k)):
        l1 = list(L_k[l]):k-2]
        l2 = list(L_k[i]):k-2]
        l1.sort()
        l2.sort()
        if l1==l2:
            Ck.append(L_k[l] | L_k[i])
return Ck
# Apriori function in which we specify min_sup and obtain frequent itemsets corresponding to it.
def apriori(itemset, min_sup = 0.12):
    CANDIDATE1 = gen_cand1(itemset)
    D = list(map(set, itemset))
    L1, sup_data = database_scan(D, CANDIDATE1, min_sup)
    L = [L1]
    k = 2
    while (len(L[k-2]) > 0):
        Ck = generate_apriori(L[k-2], k)
        L_k, sup = database_scan(D, Ck, min_sup)
        sup_data.update(sup)
        L.append(L_k)
        k += 1
    return L, sup_data
L,sdata = apriori(itemset)
new_list = []
for index in range(len(L)):
    print("\nL{} ".format(index))
    print("Number of patterns={} \n".format(len(L[index])))
    apriori_freq_pattern = [list(x) for x in L[index]]
    print(apriori_freq_pattern)

```

Output:

For min_sup = 0.4

```

L0
Number of patterns=9
[['<=50K'], ['Full-Time'], ['Adult'], ['United-States'], ['Male'], ['White'], ['Husband'], ['Married-civ-spouse'], ['Private']]

L1
Number of patterns=21
[['Male', '<=50K'], ['Private', '<=50K'], ['White', '<=50K'], ['United-States', '<=50K'], ['Full-Time', '<=50K'], ['Private', 'Full-Time'], ['White', 'Full-Time'], ['United-States', 'Full-Time'], ['Husband', 'Married-civ-spouse'], ['Private', 'White'], ['Married-civ-spouse', 'White'], ['Male', 'Private'], ['Married-civ-spouse', 'Male'], ['Husband', 'Male'], ['Male', 'White'], ['United-States', 'Private'], ['United-States', 'Married-civ-spouse'], ['United-States', 'White'], ['United-States', 'Male'], ['White', 'Adult'], ['United-States', 'Adult']]

L2
Number of patterns=12
[['United-States', 'Male', '<=50K'], ['United-States', 'White', '<=50K'], ['United-States', 'Private', '<=50K'], ['United-States', 'Full-Time', '<=50K'], ['White', 'Full-Time', '<=50K'], ['Private', 'White', '<=50K'], ['United-States', 'White', 'Full-Time'], ['United-States', 'Male', 'White'], ['United-States', 'Male', 'Private'], ['United-States', 'Private', 'White'], ['Male', 'White', 'Private'], ['Male', 'Husband', 'Married-civ-spouse']]

L3
Number of patterns=1
[['United-States', 'White', '<=50K', 'Private']]

L4
Number of patterns=0
[]

```

Fig 7: Output for min_sup = 0.4

FP-Growth Algorithm

```
#Here we start timing
start_time = timeit.default_timer()

#Defining class for Tree
class Tree(object):

    def __init__(self, value, count, parent):
        self.value = value
        self.count = count
        self.parent = parent
        self.link = None
        self.child = []

    #function for getting child
    def getting_child(self, value):
        for n in self.child:
            if n.value == value:
                return n
        return None

    #function for adding child
    def adding_child(self, value):
        a_child = Tree(value, 1, self)
        self.child.append(a_child)
        return a_child

#Defining class for Building FP Growth Tree
class Build_FPGrowth_Tree(object):

    def __init__(self, trans, threshold, r_value, r_count):
        self.frequent = self.get_freq_items(trans, threshold)
        self.headers = self.gen_header(self.frequent)
        self.root = self.gen_fp_tree(trans, r_value, r_count, self.frequent, self.headers)

    #Function for getting frequent items
    def get_freq_items(self, db_trans, sup_threshold):
        freq_item = {}
        for t in db_trans:
            for i in t:
                if i in freq_item:
                    freq_item[i] += 1
                else:
                    freq_item[i] = 1

        for key in list(freq_item.keys()):
            if freq_item[key] < sup_threshold:
                del freq_item[key]
        return freq_item

    def gen_header(self, freq):
        h_table = {}
```

```

    for key in freq.keys():
        h_table[key] = None
    return h_table

#Function for building FP Tree
def gen_fp_tree(self, db_trans, r_value, r_count, freq, heads):
    root_node = Tree(r_value, r_count, None)
    for t in db_trans:
        sort_item = [i for i in t if i in freq]
        sort_item.sort(key=lambda i: freq[i], reverse=True)
        if len(sort_item) > 0:
            #Checking if sorted items are more than 0 and if they are, we append.
            self.node_insert(sort_item, root_node, heads)
    return root_node

#Function for inserting tree
def node_insert(self, items, node, head):
    f = items[0]
    new_child = node.getting_child(f) #Getting Child
    if new_child is not None: #Checking Child
        new_child.count += 1
    else:
        new_child = node.adding_child(f) #Adding Child
        if head[f] is None:
            head[f] = new_child
        else:
            head_list = head[f]
            while head_list.link is not None:
                head_list = head_list.link
            head_list.link = new_child

    rem_items = items[1:] #Recurrive calling
    if len(rem_items) > 0:
        #Checking whether items are present and if they are, we append
        self.node_insert(rem_items, new_child, head)

def tree_path(self, n):
    child_num = len(n.child)
    if child_num > 1:
        return False
    elif child_num == 0:
        return True

#Pattern Mining
def pattern_mining(self, t_hold):
    if self.tree_path(self.root):
        return self.pattern_generation()
    else:
        return self.zpattern(self.subtrees_mining(t_hold))

#Conditional tree
def zpattern(self, f_pattern):

```

```

i = self.root.value
if i is not None:
    new_pattern = {}
    for key in f_pattern.keys():
        new_pattern[tuple((list(key) + [i]))] = f_pattern[key]
    return new_pattern
return f_pattern
#Pattern Generation
def pattern_generation(self):
    fre_pattern = {}
    i = self.frequent.keys() #Merging Index and Values
    if self.root.value is None:
        s_value = []
    else:
        s_value = [self.root.value]
        fre_pattern[tuple(s_value)] = self.root.count
    for j in range(1, len(i)):
        for k in itertools.combinations(i, j):
            ptn = tuple((list(k) + s_value))
            fre_pattern[ptn] = min([self.frequent[f] for f in k])
    return fre_pattern

def subtrees_mining(self, threshold):
    fre_pat = {}
    m_order = sorted(self.frequent.keys(), key=lambda l : self.frequent[l])
    for each_item in m_order:
        cond_tree = []
        head_node = self.headers[each_item]
        tree_suff = []
        while head_node is not None: # When node is not null we append
            tree_suff.append(head_node)
            head_node = head_node.link

        for i in tree_suff:
            freq = i.count
            path_tree = []
            parent = i.parent
            while parent.parent is not None:
                path_tree.append(parent.value)
                parent = parent.parent
            for i in range(freq):
                cond_tree.append(path_tree)
            #Constructing subtree with frequent patterns
            stree = Build_FPGrowth_Tree(cond_tree, threshold, each_item, self.frequent[each_item])
            stree_pat = stree.pattern_mining(threshold)
            # Adding patterns generated in subtree to the main tree
            for freq_pa in stree_pat.keys():
                if freq_pa in fre_pat:
                    fre_pat[freq_pa] += stree_pat[freq_pa]
                else:
                    fre_pat[freq_pa] = stree_pat[freq_pa]

```

```

        return fre_pat
#Getting Frequent patterns
def fp_growth_freq_patterns(data, sup_threshold):
    tree = Build_FPGrowth_Tree(data, sup_threshold, None, None)
    return tree.pattern_mining(sup_threshold)
#Defining minimum support
min_sup = 0.12
x = min_sup*32561
print("((Pattern) , Support Count) are:- ")
fp_freq_itemsets = fp_growth_freq_patterns(df, x)
fpgrowth_freq_itemsets = list(fp_freq_itemsets.items())
end_time = timeit.default_timer()
fpgrowth_freq_itemsets

```

Output: This is not complete output, just a snapshot of first few frequent patterns
For min_sup = 0.4

```

min_sup = 0.4
x = min_sup*32561
print("((Pattern) , Support Count) are:- ")
fp_freq_itemsets = fp_growth_freq_patterns(df, x)
fpgrowth_freq_itemsets = list(fp_freq_itemsets.items())
end_time = timeit.default_timer()
fpgrowth_freq_itemsets

((Pattern) , Support Count) are:-

[('Husband', 'Married-civ-spouse'), 13184),
 ('Husband', 'Male', 'Married-civ-spouse'), 13183),
 ('Husband', 'Male'), 13192),
 ('Male', 'Married-civ-spouse'), 13319),
 ('Married-civ-spouse', 'United-States'), 13368),
 ('Married-civ-spouse', 'White'), 13410),
 ('Adult', 'White'), 13194),
 ('Adult', 'United-States'), 13887),
 ('Full-Time', 'Private'), 14465),
 ('<=50K', 'Full-Time', 'White'), 13200),
 ('<=50K', 'Full-Time', 'United-States'), 14315),
 ('Full-Time', 'White'), 16535),
 ('Full-Time', 'United-States', 'White'), 15063),
 ('Full-Time', 'United-States'), 17734),
 ('Male', 'Private', 'White'), 13123),
 ('Male', 'Private', 'United-States'), 13209),
 ('<=50K', 'Male', 'White'), 13085),
 ('<=50K', 'Male', 'United-States'), 13389),
 ('Male', 'White'), 19174),
 ('Male', 'United-States', 'White'), 17653),
 ('Male', 'United-States'), 19488),
 ('<=50K', 'Private', 'White'), 14872),
 ('<=50K', 'Private', 'United-States', 'White'), 13452),

```

Fig 8: Output for min_sup = 0.4

Libraries Used

Pandas, NumPy and Timeit

Results:

Using the following code to get time taken by each algorithm at different minimum support

```
start_time = timeit.default_timer()
end_time = timeit.default_timer()
total_time = end_time - start_time
total_time
```

Comparison Table:

<i>Min_sup</i>	<i>Time taken by Apriori in secs</i>	<i>Time taken by improved Apriori in secs</i>	<i>Time taken by FP Growth in secs</i>
0.05	19.27	15.34	14.04
0.08	13.47	10.33	9.89
0.12	9.54	8.63	7.21

Comparing all the algorithms for different minimum support:

For Minimum Support = 0.05

Apriori

```
for index in range(len(L)):
    print("\nL{}".format(index))
    print("Number of patterns={}".format(len(L[index])))
    apriori_freq_pattern = [list(x) for x in L[index]]
    print(apriori_freq_pattern)

L7
Number of patterns=9

[['HS-grad', 'Husband', 'United-States', 'White', 'Married-civ-spouse', 'Male', 'Private', '<=50K'], ['Husband', 'United-States', 'Adult', 'White', 'Full-Time', 'Married-civ-spouse', 'Male', '<=50K'], ['Husband', 'United-States', 'White', 'Full-Time', 'Married-civ-spouse', 'Male', '<=50K'], ['Husband', 'United-States', 'Overtime', '>50K', 'White', 'Married-civ-spouse', 'Male', 'Private'], ['Husband', 'United-States', 'Overtime', '>50K', 'Adult', 'White', 'Married-civ-spouse', 'Male', 'Private'], ['Husband', 'United-States', 'Adult', 'White', 'Married-civ-spouse', 'Male', 'Private', '<=50K'], ['Husband', 'United-States', 'Overtime', 'Adult', 'White', 'Married-civ-spouse', 'Male', 'Private'], ['Husband', 'United-States', 'Adult', 'White', 'Full-Time', 'Married-civ-spouse', 'Male', 'Private']]

L8
Number of patterns=0

[]

end_time = timeit.default_timer() #ending timer
total_time = end_time - start_time
total_time

19.27536420000024
```

Fig 8: Apriori Output for min_sup = 0.05

Improved Apriori

```

for index in range(len(L)):
    print("\nL{}".format(index))
    print("Number of patterns={}".format(len(L[index])))
    improved_apriori_freq_pattern = [list(x) for x in L[index]]
    print(improved_apriori_freq_pattern)

L7
Number of patterns=9

[['United-States', 'Husband', 'Married-civ-spouse', 'Full-Time', 'Adult', 'White', '<=50K', 'Male'], ['United-States', 'Husband', 'Male', 'Private', 'Full-Time', 'White', '<=50K', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Private', 'Full-Time', 'Adult', 'White', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Private', 'White', '<=50K', 'Married-civ-spouse', 'HS-grad'], ['United-States', 'Husband', 'Male', 'Private', 'Adult', 'White', '<=50K', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Private', 'Adult', 'White', 'Overtime', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Private', '>50K', 'White', 'Overtime', 'Married-civ-spouse'], ['United-States', 'Husband', 'Married-civ-spouse', '>50K', 'Adult', 'Overtime', 'White', 'Male'], ['United-States', 'Husband', 'Male', 'Private', '>50K', 'Adult', 'White', 'Married-civ-spouse']]

L8
Number of patterns=0

[]

end_time = timeit.default_timer()
total_time = end_time - start_time
total_time

15.345461100000193

```

Fig 9: Improved Apriori Output for $\min_sup = 0.05$

FP Tree

```

min_sup = 0.05
x = min_sup*32561
print("(Pattern) , Support Count) are:- ")
fp_freq_itemsets = fp_growth_freq_patterns(df, x)
fpgrowth_freq_itemsets = list(fp_freq_itemsets.items())
end_time = timeit.default_timer()
fpgrowth_freq_itemsets

((Pattern) , Support Count) are:-

[('Masters',), 1723),
 ('Machine-op-inspct', 'United-States'), 1687),
 ('<=50K', 'Machine-op-inspct'), 1752),
 ('<=50K', 'Machine-op-inspct', 'Private'), 1688),
 ('Machine-op-inspct', 'Private'), 1913),
 ('Local-gov', 'White'), 1720),
 ('Local-gov', 'United-States', 'White'), 1632),
 ('Local-gov', 'United-States'), 1956),
 ('Married-civ-spouse', 'Self-emp-not-inc'), 1680),
 ('<=50K', 'Self-emp-not-inc', 'United-States'), 1654),
 ('<=50K', 'Self-emp-not-inc', 'White'), 1671),
 ('Male', 'Self-emp-not-inc', 'United-States'), 1956),
 ('Male', 'Self-emp-not-inc', 'United-States', 'White'), 1860),
 ('Male', 'Self-emp-not-inc', 'White'), 1985),
 ('Self-emp-not-inc', 'United-States'), 2313),
 ('Self-emp-not-inc', 'United-States', 'White'), 2190),
 ('Self-emp-not-inc', 'White'), 2342),
 ('Part-Time', 'Private', 'United-States'), 1665)

total_time = end_time - start_time
total_time

14.046819700000015

```

Fig 10: FP Growth Output for $\min_sup = 0.05$

For minimum support = 0.08

Apriori

```
for index in range(len(L)):
    print("\nL{}".format(index))
    print("Number of patterns={}".format(len(L[index])))
    apriori_freq_pattern = [list(x) for x in L[index]]
    print(apriori_freq_pattern)

L6
Number of patterns=14

[['Husband', 'White', 'Full-Time', 'Married-civ-spouse', 'Male', 'Private', '<=50K'], ['Husband', 'United-States', 'White', 'Full-Time', 'Married-civ-spouse', 'Male', '<=50K'], ['Husband', 'United-States', 'White', 'Married-civ-spouse', 'Male', 'Private', '<=50K'], ['Husband', 'United-States', 'Overtime', 'White', 'Married-civ-spouse', 'Male', 'Private'], ['Husband', 'United-States', 'Overtime', 'Adult', 'White', 'Married-civ-spouse', 'Male'], ['Young', 'United-States', 'Never-married', 'White', 'Full-Time', 'Private', '<=50K'], ['Husband', 'United-States', 'White', 'Full-Time', 'Married-civ-spouse', 'Male', 'Private'], ['Husband', 'United-States', 'Adult', 'White', 'Full-Time', 'Married-civ-spouse', 'Male'], ['Husband', 'United-States', '>50K', 'White', 'Married-civ-spouse', 'Male', 'Private'], ['Husband', 'United-States', 'Adult', '>50K', 'White', 'Married-civ-spouse', 'Male'], ['Husband', 'United-States', 'Adult', 'White', 'Married-civ-spouse', 'Male', 'Private'], ['Husband', 'United-States', 'Overtime', '>50K', 'White', 'Married-civ-spouse', 'Male'], ['Husband', 'United-States', 'Full-Time', 'Married-civ-spouse', 'Male', 'Private', '<=50K'], ['Husband', 'United-States', 'Adult', 'White', 'Married-civ-spouse', 'Male', '<=50K']]

L7
Number of patterns=0

[]

end_time = timeit.default_timer() #ending timer
total_time = end_time - start_time
total_time

13.4794423000003
```

Fig 11: Apriori Output for min_sup = 0.08

Improved Apriori

```
for index in range(len(L)):
    print("\nL{}".format(index))
    print("Number of patterns={}".format(len(L[index])))
    improved_apriori_freq_pattern = [list(x) for x in L[index]]
    print(improved_apriori_freq_pattern)

L6
Number of patterns=14

[['United-States', 'Husband', 'Male', 'Full-Time', 'White', '<=50K', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Private', 'Full-Time', 'White', '<=50K', 'Married-civ-spouse'], ['Husband', 'Male', 'Private', 'Full-Time', 'White', '<=50K', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Full-Time', 'Adult', 'White', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Adult', 'White', '<=50K', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Private', 'White', '<=50K', 'Married-civ-spouse'], ['United-States', 'Full-Time', 'Young', 'White', 'Never-married', '<=50K', 'Private'], ['United-States', 'Husband', 'Male', 'Private', 'Full-Time', 'White', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', '>50K', 'Adult', 'White', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', '>50K', 'White', 'Overtime', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Adult', 'White', 'Overtime', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Private', 'Overtime', 'White', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Private', '>50K', 'White', 'Married-civ-spouse'], ['United-States', 'Husband', 'Male', 'Private', 'Adult', 'White', 'Married-civ-spouse']]

L7
Number of patterns=0

[]

end_time = timeit.default_timer()
total_time = end_time - start_time
total_time

10.33158370000001
```

Fig 12: Improved Apriori Output for min_sup = 0.08

```
min_sup = 0.08
x = min_sup*32561
print("((Pattern) , Support Count) are:- ")
fp_freq_itemsets = fp_growth_freq_patterns(df, x)
fpgrowth_freq_itemsets = list(fp_freq_itemsets.items())
end_time = timeit.default_timer()
fpgrowth_freq_itemsets
```

```
((Pattern) , Support Count) are:-  
[[('Part-Time', 'United-States'), 2680],  
 (('<=50K', 'Part-Time'), 2733),  
 (('<=50K', 'Black'), 2737),  
 (('Black', 'United-States'), 2832),  
 (('Other-service', 'Private'), 2740),  
 (('<=50K', 'Other-service', 'Private'), 2640),  
 (('Other-service', 'United-States'), 2777),  
 (('<=50K', 'Other-service', 'United-States'), 2671),  
 (('<=50K', 'Other-service'), 3158),  
 (('Female', 'Unmarried'), 2654),  
 (('United-States', 'Unmarried'), 3033),  
 (('<=50K', 'United-States', 'Unmarried'), 2837),  
 (('<=50K', 'Unmarried'), 3228),  
 (('<=50K', 'Sales'), 2667),  
 (('Private', 'Sales'), 2942),  
 (('Private', 'Sales', 'United-States'), 2734),  
 (('Sales', 'White'), 3237),  
 (('Sales', 'United-States', 'White'), 2065)]
```

```
total_time = end_time - start_time
total_time
```

Fig 13: FP Growth Output for $\min_sup = 0.08$

Apriori

```
for index in range(len(L)):
    print("\nL{ } ".format(index))
    print("Number of patterns= { } \n".format(len(L[index])))
    apriori_freq_pattern = [list(i) for i in L[index]]
    print(apriori_freq_pattern)
```

```

L0
Number of patterns=24

[['Craft-repair'], ['Own-child'], ['Some-college'], ['Overtime'], ['>50K'], ['Young'], ['Female'], ['Prof-specialty'], ['Elderly'], ['Divorced'], ['HS-grad'], ['Private'], ['Husband'], ['Exec-managerial'], ['Married-civ-spouse'], ['Full-Time'], ['Adult'], ['<=50K'], ['United-States'], ['Male'], ['White'], ['Not-in-family'], ['Never-married'], ['Bachelors']]

```

L1
Number of patterns=106

```
[['White', 'Some-college'], ['Full-Time', 'Some-college'], ['<=50K', 'Some-college'], ['<=50K', 'Overtime'], ['Young', 'Never-married'], ['Young', 'White'], ['United-States', 'Young'], ['Own-child', 'Never-married'], ['White', 'Own-child'], ['United-States', 'Own-child'], ['<=50K', 'Own-child'], ['Young', 'Male'], ['Male', 'Some-college'], ['United-States', 'Some-college'], ['Private', 'Some-college'], ['Private', 'Never-married'], ['Female', 'Never-married'], ['>50K', 'Adult'], ['>50K', 'Private'], ['Adult', 'Overtime'], ['Private', 'Overtime'], ['Married-civ-spouse', 'HS-grad'], ['Husband', 'HS-grad'], ['Elderly', 'White'], ['>50K', 'White'], ['>50K', 'Male'], ['>50K', 'United-States'], ['>50K', 'Married-civ-spouse'], ['>50K', 'Husband'], ['White', 'Overtime'], ['Male', 'Overtime'], ['United-States', 'Overtime'], ['Married-civ-spouse', 'Overtime'], ['Husband', 'Overtime'], ['White', 'Female'], ['United-States', 'Female'], ['Female', 'Adult'], ['<=50K', 'Female'], ['Female', 'Full-Time']]
```

```
end_time = timeit.default_timer() #ending timer
total_time = end_time - start_time
total time
```

Fig 14: Apriori Output for $\min \text{ sup} = 0.12$

```
for index in range(len(L)):
    print("\nL{} ".format(index))
    print("Number of patterns={} \n".format(len(L[index])))
    improve_apriori_freq_pattern = [list(i) for i in L[index]]
    print(improve_apriori_freq_pattern)
```

ale', 'Adult'], ['Full-Time', 'White', 'Adult'], ['Full-Time', 'Private', 'Adult'], ['Full-Time', 'Married-civ-spouse', 'Adul
t'], ['Full-Time', 'Husband', 'United-States'], ['Full-Time', 'Male', 'Husband'], ['Full-Time', 'White', 'Husband'], ['Full-T
ime', 'Private', 'Husband'], ['Full-Time', 'Married-civ-spouse', 'Husband'], ['Full-Time', 'Adult', 'Husband'], ['Female', 'A
dult', 'United-States'], ['Female', 'Never-married', '<=50K'], ['Female', 'Never-married', 'United-States'], ['<=50K', 'Unit
ed-States', 'Overtime'], ['Male', 'United-States', 'Overtime'], ['White', 'United-States', 'Overtime'], ['<=50K', 'Adult', 'Un
ited-States'], ['Male', 'Adult', 'United-States'], ['White', 'Adult', 'United-States'], ['<=50K', 'Husband', 'United-State
s'], ['Male', 'Husband', 'United-States'], ['White', 'Husband', 'United-States'], ['Male', 'Married-civ-spouse', 'United-Sta
es'], ['<=50K', 'Adult', 'HS-grad'], ['Male', 'White', 'HS-grad'], ['Male', 'Husband', 'HS-grad'], ['Male', 'Married-civ-spo
se', 'HS-grad'], ['Private', 'United-States', 'Overtime'], ['Adult', 'Private', 'United-States'], ['Private', 'United-Stat
e', 'Husband'], ['Married-civ-spouse', 'Husband', 'HS-grad'], ['White', 'Adult', 'HS-grad'], ['White', 'Private', 'HS-grad'],
['Married-civ-spouse', 'United-States', 'Overtime'], ['Married-civ-spouse', 'Adult', 'United-States'], ['Married-civ-spouse',
'Husband', 'United-States'], ['Married-civ-spouse', 'Private', 'United-States'], ['Adult', 'United-States', 'Overtime'], ['Ad
ult', 'United-States', 'Husband'], ['HS-grad', 'Adult', 'United-States'], ['Husband', 'United-States', 'Overtime'], ['HS-gra
d', 'Husband', 'United-States'], ['White', 'Private', 'United-States'], ['White', 'Private', 'Never-married'], ['Full-time',
'White', 'Young'], ['<=50K', 'White', 'Young'], ['Male', 'White', 'Young'], ['<=50K', 'White', 'Own-child'], ['Full-time', 'W
hite', 'Private'], ['Bachelors', 'White', 'United-States'], ['White', 'Young', 'United-States'], ['White', 'Young', 'Never-ma
rried'], ['White', 'Private', 'Young'], ['White', 'United-States', 'Own-child'], ['<=50K', 'Full-Time', 'Married-civ-spous
e'], ['<=50K', 'Married-civ-spouse', 'United-States'], ['Full-Time', 'White', 'HS-grad'], ['<=50K', 'White', 'HS-grad'], ['Fe
male', 'Full-Time', '<=50K'], ['Female', 'Full-Time', 'United-States'], ['Female', '<=50K', 'United-States'], ['Female', 'Full

```
end_time = timeit.default_timer()
total_time = end_time - start_time
total_time
```

8.635148800000024

FP Growth

```
min_sup = 0.12
x = min_sup*32561
print("(Pattern) , Support Count) are:- ")
fp_freq_itemsets = fp_growth_freq_patterns(df, x)
fpgrowth_freq_itemsets = list(fp_freq_itemsets.items())
end_time = timeit.default_timer()
fpgrowth_freq_itemsets

((Pattern) , Support Count) are:-

[ (('Exec-managerial',), 4066),
  (('Craft-repair',), 4099),
  (('Prof-specialty',), 4140),
  (('?',), 4262),
  (('<=50K', 'Divorced'), 3980),
  (('United-States', 'Divorced'), 4162),
  (('<=50K', 'United-States', 'White', 'Own-child'), 3966),
  (('<=50K', 'White', 'Own-child'), 4196),
  (('<=50K', 'United-States', 'Never-married', 'Own-child'), 4134),
  (('<=50K', 'Never-married', 'Own-child'), 4451),
  (('<=50K', 'United-States', 'Own-child'), 4632),
  (('<=50K', 'Own-child'), 5001),
  (('United-States', 'White', 'Bachelors'), 4380),
  (('United-States', 'Bachelors'), 4766),
  (('Married-civ-spouse', 'Elderly'), 4009),
  (('United-States', '<=50K', 'Elderly'), 3941),
  (('White', 'Male', 'Elderly'), 4114),
  (('United-States', 'Male', 'Elderly'), 4159)

total_time = end_time - start_time
total_time

7.21551240000008
```

19

Conclusion

It is evident improved Apriori algorithm is better than the traditional Apriori algorithm as during the comparison we found time required to execute improved algorithm is less than the time required to execute traditional algorithm. Also, it is quite clear that FP Growth is the fastest algorithm out of all.

For the improvement, I used sampling technique and found out that the sample taken gives all the frequent patterns efficiently. Minimum support is lowered because of less number of rows. Only one scan is required as it covered all the frequent patterns as compared to normal Apriori algorithm.

L6 is same for Apriori and Improved Apriori algorithm

```
for index in range(len(L)):
    print("\nL{} ".format(index))
    print("Number of patterns={}".format(len(L[index])))
    apriori_freq_pattern = [list(i) for i in L[index]]
    print(apriori_freq_pattern)

ates', 'White', 'Full-Time'], ['Married-civ-spouse', 'Private', 'Husband', 'Male', 'White', 'Full-Time'], ['Married-civ-spouse', 'Private', 'Husband', 'Male', 'United-States', 'White'], ['Adult', 'Married-civ-spouse', 'Private', 'Husband', 'Male', 'United-States'], ['Adult', 'Married-civ-spouse', 'Private', 'Husband', 'Male', 'White'], ['Adult', 'Married-civ-spouse', 'Private', 'Husband', 'United-States', 'White'], ['Adult', 'Married-civ-spouse', 'Private', 'Male', 'United-States', 'White'], ['Adult', 'Private', 'Husband', 'Male', 'United-States', 'White'], ['>50K', 'Married-civ-spouse', 'Husband', 'Male', 'United-States', 'White'], ['Married-civ-spouse', 'Husband', 'Male', 'Overtime', 'United-States', 'White'], ['Married-civ-spouse', 'Private', 'Husband', 'Male', 'United-States', 'Full-Time'], ['<=50K', 'Married-civ-spouse', 'Private', 'Husband', 'Male', 'United-States'], ['<=50K', 'Private', 'Male', 'United-States', 'White', 'Full-Time'], ['Adult', 'Married-civ-spouse', 'Husband', 'Male', 'United-States', 'White'], ['<=50K', 'Married-civ-spouse', 'Husband', 'Male', 'United-States', 'White']]

L6
Number of patterns=1

[['Adult', 'Married-civ-spouse', 'Private', 'Husband', 'Male', 'United-States', 'White']]

L7
Number of patterns=0

[]

end_time = timeit.default_timer() #ending timer
total_time = end_time - start_time
total_time

9.549866300000076
```

Fig 17: Apriori Output for $\min_sup = 0.12$

```
for index in range(len(L)):
    print("\nL{} ".format(index))
    print("Number of patterns={}".format(len(L[index])))
    improve_apriori_freq_pattern = [list(i) for i in L[index]]
    print(improve_apriori_freq_pattern)

e', 'Married-civ-spouse', 'White', 'Male', 'Husband', 'United-States'], ['<=50K', 'Full-Time', 'Married-civ-spouse', 'White', 'Male', 'Husband'], ['Married-civ-spouse', 'White', 'Male', 'Husband', 'Private', 'United-States'], ['Married-civ-spouse', 'White', 'Adult', 'Husband', 'Private', 'United-States'], ['Married-civ-spouse', 'White', 'Adult', 'Male', 'Private', 'United-States'], ['Married-civ-spouse', 'White', 'Adult', 'Male', 'Husband', 'United-States'], ['Married-civ-spouse', 'White', 'Overtime', 'Male', 'Husband', 'United-States'], ['<=50K', 'Married-civ-spouse', 'Male', 'Husband', 'Private', 'United-States'], ['Married-civ-spouse', 'Adult', 'Male', 'Husband', 'Private', 'United-States'], ['White', 'Adult', 'Male', 'Husband', 'Private', 'United-States'], ['<=50K', 'Married-civ-spouse', 'White', 'Male', 'Husband', 'United-States'], ['<=50K', 'Full-Time', 'White', 'Male', 'Private', 'United-States'], ['Married-civ-spouse', 'White', 'Adult', 'Male', 'Husband', 'Private'], ['<=50K', 'Married-civ-spouse', 'White', 'Male', 'Husband', 'Private']]

L6
Number of patterns=1

[['Married-civ-spouse', 'White', 'Adult', 'Male', 'Husband', 'Private', 'United-States']]

L7
Number of patterns=0

[]

end_time = timeit.default_timer()
total_time = end_time - start_time
total_time

8.635148800000024
```

Fig 18: Improved Apriori Output for $\min_sup = 0.12$

Reference

- [1] Yu Cheng, Ying Xiong. Research and Improvement of Apriori Algorithm for Association Rules. 2010 2nd International Workshop on Intelligent Systems and Applications
- [2] Jiawei Han, Micheline Kamber, Jian Pei. Data Mining Concepts and Techniques, 248-264.
- [3] https://notebook.community/aleph314/K2/Data%20Mining/Association%20Rules/assoc_mining_problems
- [4] Reynaldo John Tristan Mahinay Jr., Franz Stewart Dizon, Stephen Kyle Farinas and Harry Pardo. Learning of High Dengue Incidence with Clustering and FP-Growth Algorithm using WHO Historical Data. *3rd IEEE International Conference on Agents (ICA 2018)*
- [5] Jiao Yabing. Research of an Improved Apriori Algorithm in Data Mining Association Rules. International Journal of Computer and Communication Engineering, Vol. 2, No. 1, January 2013
- [6] <https://towardsdatascience.com/understand-and-build-fp-growth-algorithm-in-python-d8b989bab342>
- [7] <https://www.geeksforgeeks.org/ml-frequent-pattern-growth-algorithm/>