

## Assignment 4 - Helper Functions

### Samuel Steiner -- Assignment 4 --- Association Rules

We begin by including the functions to generate frequent itemsets (via the Apriori algorithm) and resulting association rules:

In [ ]:

```
# (c) 2016 Everaldo Aguiar & Reid Johnson
#
# Modified from:
# Marcel Caraciolo (https://gist.github.com/marcelcaraciolo/1423287)
#
# Functions to compute and extract association rules from a given frequent itemset
# generated by the Apriori algorithm.
#
# The Apriori algorithm is defined by Agrawal and Srikant in:
# Fast algorithms for mining association rules
# Proc. 20th int. conf. very large data bases, VLDB. Vol. 1215. 1994
import csv
import numpy as np

def load_dataset(filename):
    '''Loads an example of market basket transactions from a provided csv file.

    Returns: A list (database) of lists (transactions). Each element of a transaction is
    an item.
    '''

    with open(filename, 'r') as dest_f:
        data_iter = csv.reader(dest_f, delimiter=',', quotechar='"')
        data = [data for data in data_iter]
        data_array = np.asarray(data)

    return data_array

def apriori(dataset, min_support=0.5, verbose=False):
    """Implements the Apriori algorithm.

    The Apriori algorithm will iteratively generate new candidate
    k-itemsets using the frequent (k-1)-itemsets found in the previous
    iteration.

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

    min_support : float
        The minimum support threshold. Defaults to 0.5.

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

    support_data : dict
        The support data for all candidate itemsets.

    References
    -----
    .. [1] R. Agrawal, R. Srikant, "Fast Algorithms for Mining Association
        Rules", 1994.

    """
    C1 = create_candidates(dataset)
    D = list(map(set, dataset))
    F1, support_data = support_prune(D, C1, min_support, verbose=False) # prune candidate 1-itemsets
    F = [F1] # list of frequent itemsets; initialized to frequent 1-itemsets
    k = 2 # the itemset cardinality
    while (len(F[k - 2]) > 0):
        Ck = apriori_gen(F[k-2], k) # generate candidate itemsets
        Fk, supK = support_prune(D, Ck, min_support) # prune candidate itemsets
        support_data.update(supK) # update the support counts to reflect pruning
        F.append(Fk) # add the pruned candidate itemsets to the list of frequent itemsets
        k += 1

    if verbose:
        # Print a list of all the frequent itemsets.
        for kset in F:
            for item in kset:
                print("\n \
                    + "{" \
                    + "".join(str(i) + ", " for i in iter(item)).rstrip(', ') \
                    + "}" \
                    + ": sup = " + str(round(support_data[item], 3)))

    return F, support_data

def create_candidates(dataset, verbose=False):
    """Creates a list of candidate 1-itemsets from a list of transactions.

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

    Returns
    -----
    C1 : list
        The list of candidate 1-itemsets.
    """
```

Parameters

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

Returns

The list of candidate itemsets (c1) passed as a frozenset (a set that is immutable and hashable).

```
"""
c1 = [] # list of all items in the database of transactions
for transaction in dataset:
    for item in transaction:
        if not [item] in c1:
            c1.append([item])
c1.sort()

if verbose:
    # Print a list of all the candidate items.
    print(""" \
+ "{" \
+ "".join(str(i[0]) + ", " for i in iter(c1)).rstrip(', ') \
+ "}")

# Map c1 to a frozenset because it will be the key of a dictionary.
return list(map(frozenset, c1))
```

**def** support\_prune(dataset, candidates, min\_support, verbose=False):  
"""Returns all candidate itemsets that meet a minimum support threshold.

By the apriori principle, if an itemset is frequent, then all of its subsets must also be frequent. As a result, we can perform support-based pruning to systematically control the exponential growth of candidate itemsets. Thus, itemsets that do not meet the minimum support level are pruned from the input list of itemsets (dataset).

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

retlist : list  
The list of frequent itemsets.

support\_data : dict  
The support data for all candidate itemsets.

```
"""
sscnt = {} # set for support counts
for tid in dataset:
    for can in candidates:
        if can.issubset(tid):
            sscnt.setdefault(can, 0)
            sscnt[can] += 1

num_items = float(len(dataset)) # total number of transactions in the dataset
retlist = [] # array for unpruned itemsets
support_data = {} # set for support data for corresponding itemsets
for key in sscnt:
    # Calculate the support of itemset key.
    support = sscnt[key] / num_items
    if support >= min_support:
        retlist.insert(0, key)
        support_data[key] = support

# Print a list of the pruned itemsets.
if verbose:
    for kset in retlist:
        for item in kset:
            print("{ " + str(item) + "}")
    print("""
for key in sscnt:
    print(""" \
+ "{" \
+ "".join([str(i) + ", " for i in iter(key)]).rstrip(', ') \
+ "}" \
+ ": sup = " + str(support_data[key]))
```

```
return retlist, support_data
```

```
def apriori_gen(freq_sets, k):
    """Generates candidate itemsets (via the Fk-1 x Fk-1 method).

    This operation generates new candidate k-itemsets based on the frequent
    (k-1)-itemsets found in the previous iteration. The candidate generation
    procedure merges a pair of frequent (k-1)-itemsets only if their first k-2
    items are identical.

    Parameters
    -----
    freq_sets : list
        The list of frequent (k-1)-itemsets.

    k : integer
        The cardinality of the current itemsets being evaluated.

    Returns
    -----
    retlist : list
        The list of merged frequent itemsets.
    """
    retList = [] # list of merged frequent itemsets
    lenLk = len(freq_sets) # number of frequent itemsets
    for i in range(lenLk):
        for j in range(i+1, lenLk):
            a=list(freq_sets[i])
            b=list(freq_sets[j])
            a.sort()
            b.sort()
            F1 = a[:k-2] # first k-2 items of freq_sets[i]
            F2 = b[:k-2] # first k-2 items of freq_sets[j]

            if F1 == F2: # if the first k-2 items are identical
                # Merge the frequent itemsets.
                retList.append(freq_sets[i] | freq_sets[j])

    return retList

def rules_from_conseq(freq_set, H, support_data, rules, min_confidence=0.5, verbose=False):
    """Generates a set of candidate rules.

    Parameters
    -----
    freq_set : frozenset
        The complete list of frequent itemsets.

    H : list
        A list of frequent itemsets (of a particular length).

    support_data : dict
        The support data for all candidate itemsets.

    rules : list
        A potentially incomplete set of candidate rules above the minimum
        confidence threshold.

    min_confidence : float
        The minimum confidence threshold. Defaults to 0.5.
    """
    m = len(H[0])
    if m == 1:
        Hmp1 = calc_confidence(freq_set, H, support_data, rules, min_confidence, verbose)
    if (len(freq_set) > (m+1)):
        Hmp1 = apriori_gen(H, m+1) # generate candidate itemsets
        Hmp1 = calc_confidence(freq_set, Hmp1, support_data, rules, min_confidence, verbose)
        if len(Hmp1) > 1:
            # If there are candidate rules above the minimum confidence
            # threshold, recurse on the list of these candidate rules.
            rules_from_conseq(freq_set, Hmp1, support_data, rules, min_confidence, verbose)

def calc_confidence(freq_set, H, support_data, rules, min_confidence=0.5, verbose=False):
    """Evaluates the generated rules.

    One measurement for quantifying the goodness of association rules is
    confidence. The confidence for a rule 'P implies H' (P -> H) is defined as
    the support for P and H divided by the support for P
    (support (P|H) / support(P)), where the | symbol denotes the set union
    (thus P|H means all the items in set P or in set H).

    To calculate the confidence, we iterate through the frequent itemsets and
    associated support data. For each frequent itemset, we divide the support
    of the itemset by the support of the antecedent (left-hand-side of the
    rule).

    Parameters
```

```

-----
freq_set : frozenset
    The complete list of frequent itemsets.

H : list
    A list of frequent itemsets (of a particular length).

min_support : float
    The minimum support threshold.

rules : list
    A potentially incomplete set of candidate rules above the minimum
    confidence threshold.

min_confidence : float
    The minimum confidence threshold. Defaults to 0.5.

Returns
-----
pruned_H : list
    The list of candidate rules above the minimum confidence threshold.
    """
pruned_H = [] # list of candidate rules above the minimum confidence threshold
for consequent in H: # iterate over the frequent itemsets
    conf = support_data[freq_set] / support_data[freq_set - consequent]
    if conf >= min_confidence:
        rules.append((freq_set - consequent, consequent, conf))
        pruned_H.append(consequent)

    if verbose:
        print(""" \
            + "{" \
            + "".join([str(i) + ", " for i in iter(freq_set - consequent)]).rstrip(', ') \
            + "}" \
            + " ----> " \
            + "{" \
            + "".join([str(i) + ", " for i in iter(consequent)]).rstrip(', ') \
            + "}" \
            + ": conf = " + str(round(conf, 3)) \
            + ", sup = " + str(round(support_data[freq_set], 3)))

return pruned_H

def generate_rules(F, support_data, min_confidence=0.5, verbose=True):
    """Generates a set of candidate rules from a list of frequent itemsets.

    For each frequent itemset, we calculate the confidence of using a
    particular item as the rule consequent (right-hand-side of the rule). By
    testing and merging the remaining rules, we recursively create a list of
    pruned rules.

    Parameters
    -----
    F : list
        A list of frequent itemsets.

    support_data : dict
        The corresponding support data for the frequent itemsets (L).

    min_confidence : float
        The minimum confidence threshold. Defaults to 0.5.

    Returns
    -----
    rules : list
        The list of candidate rules above the minimum confidence threshold.
        """
    rules = []
    for i in range(1, len(F)):
        for freq_set in F[i]:
            H1 = [frozenset([itemset]) for itemset in freq_set]
            if (i > 1):
                rules_from_conseq(freq_set, H1, support_data, rules, min_confidence, verbose)
            else:
                calc_confidence(freq_set, H1, support_data, rules, min_confidence, verbose)

    return rules

```

## FP Growth code from notes

In [ ]:

```
# (c) 2014 Reid Johnson
#
# Modified from:
# Eric Naeseth <eric@naeseth.com>
# (https://github.com/enaeseth/python-fp-growth/blob/master/fp_growth.py)
#
# A Python implementation of the FP-growth algorithm.

from collections import defaultdict, namedtuple
from itertools import imap

__author__ = 'Eric Naeseth <eric@naeseth.com>'
__copyright__ = 'Copyright © 2009 Eric Naeseth'
__license__ = 'MIT License'

def fpgrowth(dataset, min_support=0.5, include_support=True, verbose=False):
    """Implements the FP-growth algorithm.

    The `dataset` parameter can be any iterable of iterables of items.
    `min_support` should be an integer specifying the minimum number of
    occurrences of an itemset for it to be accepted.

    Each item must be hashable (i.e., it must be valid as a member of a
    dictionary or a set).

    If `include_support` is true, yield (itemset, support) pairs instead of
    just the itemsets.

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

    min_support : float
        The minimum support threshold. Defaults to 0.5.

    include_support : bool
        Include support in output (default=False).

    References
    -----
    .. [1] J. Han, J. Pei, Y. Yin, "Mining Frequent Patterns without Candidate
        Generation," 2000.

    """

    F = []
    support_data = {}
    for k,v in find_frequent_itemsets(dataset, min_support=min_support, include_support=include_support, verbose=verbose):
        F.append(frozenset(k))
        support_data[frozenset(k)] = v

    # Create one array with subarrays that hold all transactions of equal length.
    def bucket_list(nested_list, sort=True):
        bucket = defaultdict(list)
        for sublist in nested_list:
            bucket[len(sublist)].append(sublist)
        return [v for k,v in sorted(bucket.items())] if sort else bucket.values()

    F = bucket_list(F)

    return F, support_data

def find_frequent_itemsets(dataset, min_support, include_support=False, verbose=False):
    """
    Find frequent itemsets in the given transactions using FP-growth. This
    function returns a generator instead of an eagerly-populated list of items.

    The `dataset` parameter can be any iterable of iterables of items.
    `min_support` should be an integer specifying the minimum number of
    occurrences of an itemset for it to be accepted.

    Each item must be hashable (i.e., it must be valid as a member of a
    dictionary or a set).

    If `include_support` is true, yield (itemset, support) pairs instead of
    just the itemsets.

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

    min_support : float
    """
```

```
min_support : float
    The minimum support threshold. Defaults to 0.5.
```

```
include_support : bool
    Include support in output (default=False).
```

```
"""
```

```
items = defaultdict(lambda: 0) # mapping from items to their supports
processed_transactions = []
```

```
# Load the passed-in transactions and count the support that individual
# items have.
```

```
for transaction in dataset:
    processed = []
    for item in transaction:
        items[item] += 1
        processed.append(item)
    processed_transactions.append(processed)
```

```
# Remove infrequent items from the item support dictionary.
items = dict((item, support) for item, support in items.items()
             if support >= min_support)
```

```
# Build our FP-tree. Before any transactions can be added to the tree, they
# must be stripped of infrequent items and their surviving items must be
# sorted in decreasing order of frequency.
```

```
def clean_transaction(transaction):
    #transaction = filter(lambda v: v in items, transaction)
    transaction.sort(key=lambda v: items[v], reverse=True)
    return transaction
```

```
master = FPTree()
for transaction in map(clean_transaction, processed_transactions):
    master.add(transaction)
```

```
support_data = {}
```

```
def find_with_suffix(tree, suffix):
    for item, nodes in tree.items():
        support = float(sum(n.count for n in nodes)) / len(dataset)
        if support >= min_support and item not in suffix:
            # New winner!
            found_set = [item] + suffix
            support_data[frozenset(found_set)] = support
            yield (found_set, support) if include_support else found_set

            # Build a conditional tree and recursively search for frequent
            # itemsets within it.
            cond_tree = conditional_tree_from_paths(tree.prefix_paths(item),
                                                    min_support)
            for s in find_with_suffix(cond_tree, found_set):
                yield s # pass along the good news to our caller
```

```
if verbose:
    # Print a list of all the frequent itemsets.
    for itemset, support in find_with_suffix(master, []):
        print(""" \
              + "{" \
              + "".join(str(i) + ", " for i in iter(itemset)).rstrip(', ') \
              + "}" \
              + ": sup = " + str(round(support_data[frozenset(itemset)], 3))
```

```
# Search for frequent itemsets, and yield the results we find.
```

```
for itemset in find_with_suffix(master, []):
    yield itemset
```

```
class FPTree(object):
```

```
"""
```

```
An FP tree.
```

```
This object may only store transaction items that are hashable (i.e., all
items must be valid as dictionary keys or set members).
```

```
"""
```

```
Route = namedtuple('Route', 'head tail')
```

```
def __init__(self):
    # The root node of the tree.
    self._root = FPNode(self, None, None)
```

```
# A dictionary mapping items to the head and tail of a path of
# "neighbors" that will hit every node containing that item.
self._routes = {}
```

```
@property
```

```
def root(self):
    """The root node of the tree."""
    return self._root
```

```

def add(self, transaction):
    """
    Adds a transaction to the tree.
    """

    point = self._root

    for item in transaction:
        next_point = point.search(item)
        if next_point:
            # There is already a node in this tree for the current
            # transaction item; reuse it.
            next_point.increment()
        else:
            # Create a new point and add it as a child of the point we're
            # currently looking at.
            next_point = FPNode(self, item)
            point.add(next_point)

            # Update the route of nodes that contain this item to include
            # our new node.
            self._update_route(next_point)

        point = next_point

def _update_route(self, point):
    """Add the given node to the route through all nodes for its item."""
    assert self is point.tree

    try:
        route = self._routes[point.item]
        route[1].neighbor = point # route[1] is the tail
        self._routes[point.item] = self.Route(route[0], point)
    except KeyError:
        # First node for this item; start a new route.
        self._routes[point.item] = self.Route(point, point)

def items(self):
    """
    Generate one 2-tuples for each item represented in the tree. The first
    element of the tuple is the item itself, and the second element is a
    generator that will yield the nodes in the tree that belong to the item.
    """
    for item in self._routes.keys():
        yield (item, self.nodes(item))

def nodes(self, item):
    """
    Generates the sequence of nodes that contain the given item.
    """

    try:
        node = self._routes[item][0]
    except KeyError:
        return

    while node:
        yield node
        node = node.neighbor

def prefix_paths(self, item):
    """Generates the prefix paths that end with the given item."""

    def collect_path(node):
        path = []
        while node and not node.root:
            path.append(node)
            node = node.parent
        path.reverse()
        return path

    return (collect_path(node) for node in self.nodes(item))

def inspect(self):
    print("Tree:")
    self.root.inspect(1)

    print("")
    print("Routes:")
    for item, nodes in self.items():
        print("  %r" % item)
        for node in nodes:
            print("    %r" % node)

def removed(self, node):

```



```

        """Called when node` is removed from the tree; performs cleanup."""
        head, tail = self._routes[node.item]
        if node is head:
            if node is tail or not node.neighbor:
                # It was the sole node.
                del self._routes[node.item]
            else:
                self._routes[node.item] = self.Route(node.neighbor, tail)
        else:
            for n in self.nodes(node.item):
                if n.neighbor is node:
                    n.neighbor = node.neighbor # skip over
                    if node is tail:
                        self._routes[node.item] = self.Route(head, n)
                    break

def conditional_tree_from_paths(paths, min_support):
    """Builds a conditional FP-tree from the given prefix paths."""
    tree = FPTree()
    condition_item = None
    items = set()

    # Import the nodes in the paths into the new tree. Only the counts of the
    # leaf nodes matter; the remaining counts will be reconstructed from the
    # leaf counts.
    for path in paths:
        if condition_item is None:
            condition_item = path[-1].item

        point = tree.root
        for node in path:
            next_point = point.search(node.item)
            if not next_point:
                # Add a new node to the tree.
                items.add(node.item)
                count = node.count if node.item == condition_item else 0
                next_point = FPNode(tree, node.item, count)
                point.add(next_point)
                tree._update_route(next_point)
            point = next_point

    assert condition_item is not None

    # Calculate the counts of the non-leaf nodes.
    for path in tree.prefix_paths(condition_item):
        count = path[-1].count
        for node in reversed(path[:-1]):
            node._count += count

    # Eliminate the nodes for any items that are no longer frequent.
    for item in items:
        support = sum(n.count for n in tree.nodes(item))
        if support < min_support:
            # Doesn't make the cut anymore
            for node in tree.nodes(item):
                if node.parent is not None:
                    node.parent.remove(node)

    # Finally, remove the nodes corresponding to the item for which this
    # conditional tree was generated.
    for node in tree.nodes(condition_item):
        if node.parent is not None: # the node might already be an orphan
            node.parent.remove(node)

    return tree

class FPNode(object):
    """A node in an FP tree."""

    def __init__(self, tree, item, count=1):
        self._tree = tree
        self._item = item
        self._count = count
        self._parent = None
        self._children = {}
        self._neighbor = None

    def add(self, child):
        """Adds the given FPNode `child` as a child of this node."""

        if not isinstance(child, FPNode):
            raise TypeError("Can only add other FPNodes as children")

        if not child.item in self._children:
            self._children[child.item] = child
            child._parent = self

```

```

child.parent = self

def search(self, item):
    """
    Checks to see if this node contains a child node for the given item.
    If so, that node is returned; otherwise, `None` is returned.
    """

    try:
        return self._children[item]
    except KeyError:
        return None

def remove(self, child):
    try:
        if self._children[child.item] is child:
            del self._children[child.item]
            child.parent = None
            self._tree._removed(child)
            for sub_child in child.children:
                try:
                    # Merger case: we already have a child for that item, so
                    # add the sub-child's count to our child's count.
                    self._children[sub_child.item]._count += sub_child.count
                    sub_child.parent = None # it's an orphan now
                except KeyError:
                    # Turns out we don't actually have a child, so just add
                    # the sub-child as our own child.
                    self.add(sub_child)
            child._children = {}
        else:
            raise ValueError("that node is not a child of this node")
    except KeyError:
        raise ValueError("that node is not a child of this node")

def __contains__(self, item):
    return item in self._children

@property
def tree(self):
    """The tree in which this node appears."""
    return self._tree

@property
def item(self):
    """The item contained in this node."""
    return self._item

@property
def count(self):
    """The count associated with this node's item."""
    return self._count

def increment(self):
    """Increments the count associated with this node's item."""
    if self._count is None:
        raise ValueError("Root nodes have no associated count.")
    self._count += 1

@property
def root(self):
    """True if this node is the root of a tree; false if otherwise."""
    return self._item is None and self._count is None

@property
def leaf(self):
    """True if this node is a leaf in the tree; false if otherwise."""
    return len(self._children) == 0

def parent():
    doc = "The node's parent."
    def fget(self):
        return self._parent
    def fset(self, value):
        if value is not None and not isinstance(value, FPNode):
            raise TypeError("A node must have an FPNode as a parent.")
        if value and value.tree is not self.tree:
            raise ValueError("Cannot have a parent from another tree.")
        self._parent = value
    return locals()
parent = property(**parent())

def neighbor():
    doc = """
    The node's neighbor; the one with the same value that is "to the right"
    of it in the tree.
    """

```

```

def fget(self):
    return self._neighbor
def fset(self, value):
    if value is not None and not isinstance(value, FPNode):
        raise TypeError("A node must have an FPNode as a neighbor.")
    if value and value.tree is not self.tree:
        raise ValueError("Cannot have a neighbor from another tree.")
    self._neighbor = value
    return locals()
neighbor = property(**neighbor())

@property
def children(self):
    """The nodes that are children of this node."""
    return tuple(self._children.values())

def inspect(self, depth=0):
    print((' ' * depth) + repr(self))
    for child in self.children:
        child.inspect(depth + 1)

def __repr__(self):
    if self.root:
        return "<%s (root)>" % type(self).__name__
    return "<%s %r (%r)>" % (type(self).__name__, self.item, self.count)

```

To load our dataset of grocery transactions, use the command below

```

In [ ]: dataset = load_dataset('grocery.csv')
        D = list(map(set, dataset))

```

/Users/sam/opt/anaconda3/lib/python3.9/site-packages/numpy/core/\_asarray.py:102: VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype=object' when creating the ndarray.  
 return array(a, dtype=object, copy=False, order=order)

*dataset* is now a ndarray containing each of the 9835 transactions

```

In [ ]: type(dataset)

```

```

Out[ ]: numpy.ndarray

```

```

In [ ]: dataset.shape

```

```

Out[ ]: (9835,)

```

```

In [ ]: dataset[0]

```

```

Out[ ]: ['citrus fruit', 'semi-finished bread', 'margarine', 'ready soups']

```

```

In [ ]: dataset[1]

```

```

Out[ ]: ['tropical fruit', 'yogurt', 'coffee']

```

*D* Contains that dataset in a set format (which excludes duplicated items and sorts them)

```

In [ ]: type(D[0])

```

```

Out[ ]: set

```

```

In [ ]: D[0]

```

```

Out[ ]: {'citrus fruit', 'margarine', 'ready soups', 'semi-finished bread'}

```

Complete the assignment below by making use of the provided funtions.

You may use the notebook file attached with lesson 3 as a reference

Part 1

```
In [ ]: import numpy as np
X = [i for i in np.arange(0.02, 0.055, 0.005)]
y = .3
```

```
In [ ]: for x in X:
print(f"SUPPORT = {x} or greater")
F_ap, support_data_ap = apriori(dataset, min_support=x)
H_ap = generate_rules(F_ap, support_data_ap, min_confidence=y)
```

```
SUPPORT = 0.02 or greater
{yogurt} ----> {other vegetables}: conf = 0.311, sup = 0.043
{pip fruit} ----> {other vegetables}: conf = 0.345, sup = 0.026
{citrus fruit} ----> {other vegetables}: conf = 0.349, sup = 0.029
{fruit/vegetable juice} ----> {whole milk}: conf = 0.368, sup = 0.027
{frankfurter} ----> {whole milk}: conf = 0.348, sup = 0.021
{newspapers} ----> {whole milk}: conf = 0.343, sup = 0.027
{margarine} ----> {whole milk}: conf = 0.413, sup = 0.024
{pip fruit} ----> {whole milk}: conf = 0.398, sup = 0.03
{rolls/buns} ----> {whole milk}: conf = 0.308, sup = 0.057
{beef} ----> {whole milk}: conf = 0.405, sup = 0.021
{sausage} ----> {whole milk}: conf = 0.318, sup = 0.03
{frozen vegetables} ----> {whole milk}: conf = 0.425, sup = 0.02
{domestic eggs} ----> {other vegetables}: conf = 0.351, sup = 0.022
{butter} ----> {other vegetables}: conf = 0.361, sup = 0.02
{pastry} ----> {whole milk}: conf = 0.374, sup = 0.033
{brown bread} ----> {whole milk}: conf = 0.389, sup = 0.025
{domestic eggs} ----> {whole milk}: conf = 0.473, sup = 0.03
{pork} ----> {whole milk}: conf = 0.384, sup = 0.022
{pork} ----> {other vegetables}: conf = 0.376, sup = 0.022
{whipped/sour cream} ----> {whole milk}: conf = 0.45, sup = 0.032
{whipped/sour cream} ----> {other vegetables}: conf = 0.403, sup = 0.029
{root vegetables} ----> {whole milk}: conf = 0.449, sup = 0.049
{sausage} ----> {rolls/buns}: conf = 0.326, sup = 0.031
{root vegetables} ----> {other vegetables}: conf = 0.435, sup = 0.047
{citrus fruit} ----> {whole milk}: conf = 0.369, sup = 0.031
{tropical fruit} ----> {whole milk}: conf = 0.403, sup = 0.042
{bottled water} ----> {whole milk}: conf = 0.311, sup = 0.034
{curd} ----> {whole milk}: conf = 0.49, sup = 0.026
{tropical fruit} ----> {other vegetables}: conf = 0.342, sup = 0.036
{yogurt} ----> {whole milk}: conf = 0.402, sup = 0.056
{butter} ----> {whole milk}: conf = 0.497, sup = 0.028
{other vegetables} ----> {whole milk}: conf = 0.387, sup = 0.075
{yogurt, whole milk} ----> {other vegetables}: conf = 0.397, sup = 0.022
{yogurt, other vegetables} ----> {whole milk}: conf = 0.513, sup = 0.022
{whole milk, root vegetables} ----> {other vegetables}: conf = 0.474, sup = 0.023
{other vegetables, root vegetables} ----> {whole milk}: conf = 0.489, sup = 0.023
{other vegetables, whole milk} ----> {root vegetables}: conf = 0.31, sup = 0.023
SUPPORT = 0.025 or greater
{yogurt} ----> {other vegetables}: conf = 0.311, sup = 0.043
{pip fruit} ----> {other vegetables}: conf = 0.345, sup = 0.026
{citrus fruit} ----> {other vegetables}: conf = 0.349, sup = 0.029
{fruit/vegetable juice} ----> {whole milk}: conf = 0.368, sup = 0.027
{newspapers} ----> {whole milk}: conf = 0.343, sup = 0.027
{pip fruit} ----> {whole milk}: conf = 0.398, sup = 0.03
{rolls/buns} ----> {whole milk}: conf = 0.308, sup = 0.057
{sausage} ----> {whole milk}: conf = 0.318, sup = 0.03
{pastry} ----> {whole milk}: conf = 0.374, sup = 0.033
{brown bread} ----> {whole milk}: conf = 0.389, sup = 0.025
{domestic eggs} ----> {whole milk}: conf = 0.473, sup = 0.03
{whipped/sour cream} ----> {whole milk}: conf = 0.45, sup = 0.032
{whipped/sour cream} ----> {other vegetables}: conf = 0.403, sup = 0.029
{root vegetables} ----> {whole milk}: conf = 0.449, sup = 0.049
{sausage} ----> {rolls/buns}: conf = 0.326, sup = 0.031
{root vegetables} ----> {other vegetables}: conf = 0.435, sup = 0.047
{citrus fruit} ----> {whole milk}: conf = 0.369, sup = 0.031
{tropical fruit} ----> {whole milk}: conf = 0.403, sup = 0.042
{bottled water} ----> {whole milk}: conf = 0.311, sup = 0.034
{curd} ----> {whole milk}: conf = 0.49, sup = 0.026
{tropical fruit} ----> {other vegetables}: conf = 0.342, sup = 0.036
{yogurt} ----> {whole milk}: conf = 0.402, sup = 0.056
{butter} ----> {whole milk}: conf = 0.497, sup = 0.028
{other vegetables} ----> {whole milk}: conf = 0.387, sup = 0.075
SUPPORT = 0.030000000000000002 or greater
{yogurt} ----> {other vegetables}: conf = 0.311, sup = 0.043
{pip fruit} ----> {whole milk}: conf = 0.398, sup = 0.03
{rolls/buns} ----> {whole milk}: conf = 0.308, sup = 0.057
{pastry} ----> {whole milk}: conf = 0.374, sup = 0.033
{whipped/sour cream} ----> {whole milk}: conf = 0.45, sup = 0.032
{root vegetables} ----> {whole milk}: conf = 0.449, sup = 0.049
{sausage} ----> {rolls/buns}: conf = 0.326, sup = 0.031
{root vegetables} ----> {other vegetables}: conf = 0.435, sup = 0.047
{citrus fruit} ----> {whole milk}: conf = 0.369, sup = 0.031
{tropical fruit} ----> {whole milk}: conf = 0.403, sup = 0.042
{bottled water} ----> {whole milk}: conf = 0.311, sup = 0.034
```

```

{tropical fruit} ----> {other vegetables}:  conf = 0.342, sup = 0.036
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
SUPPORT = 0.035 or greater
{yogurt} ----> {other vegetables}:  conf = 0.311, sup = 0.043
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{tropical fruit} ----> {whole milk}:  conf = 0.403, sup = 0.042
{tropical fruit} ----> {other vegetables}:  conf = 0.342, sup = 0.036
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
SUPPORT = 0.04000000000000001 or greater
{yogurt} ----> {other vegetables}:  conf = 0.311, sup = 0.043
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{tropical fruit} ----> {whole milk}:  conf = 0.403, sup = 0.042
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
SUPPORT = 0.045000000000000005 or greater
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
SUPPORT = 0.05 or greater
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
SUPPORT = 0.055000000000000001 or greater
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075

```

Using the provided code for apriori makes this step simple, once the code was loaded I made it simple to experiment with the support and confidence variables by using a tuple to describe them. From there I experimented with values close to those in the directions for this assignment for part 3, .02 and .3 respectively. I found that there was a lot of information to gather from the data and since the data had ten thousand entries the seemingly low support values still had a large count. 1% of 10000 is 100 transactions. Looking at 5% or 500 transactions with the same confidence value as part 3 generated only 3 rules.

## Part 2

```

In [ ]: for x in X:
        print(f"SUPPORT = {x} or greater")
        F_fp, support_data_fp = fpgrowth(dataset, min_support=x)
        H_fp = generate_rules(F_fp, support_data_fp, min_confidence=y)

```

```

SUPPORT = 0.02 or greater
{citrus fruit} ----> {whole milk}:  conf = 0.369, sup = 0.031
{citrus fruit} ----> {other vegetables}:  conf = 0.349, sup = 0.029
{margarine} ----> {whole milk}:  conf = 0.413, sup = 0.024
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{yogurt} ----> {other vegetables}:  conf = 0.311, sup = 0.043
{tropical fruit} ----> {other vegetables}:  conf = 0.342, sup = 0.036
{tropical fruit} ----> {whole milk}:  conf = 0.403, sup = 0.042
{pip fruit} ----> {whole milk}:  conf = 0.398, sup = 0.03
{pip fruit} ----> {other vegetables}:  conf = 0.345, sup = 0.026
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
{butter} ----> {whole milk}:  conf = 0.497, sup = 0.028
{butter} ----> {other vegetables}:  conf = 0.361, sup = 0.02
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{bottled water} ----> {whole milk}:  conf = 0.311, sup = 0.034
{curd} ----> {whole milk}:  conf = 0.49, sup = 0.026
{beef} ----> {whole milk}:  conf = 0.405, sup = 0.021
{frankfurter} ----> {whole milk}:  conf = 0.348, sup = 0.021
{newspapers} ----> {whole milk}:  conf = 0.343, sup = 0.027
{fruit/vegetable juice} ----> {whole milk}:  conf = 0.368, sup = 0.027
{pastry} ----> {whole milk}:  conf = 0.374, sup = 0.033
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
{sausage} ----> {rolls/buns}:  conf = 0.326, sup = 0.031
{sausage} ----> {whole milk}:  conf = 0.318, sup = 0.03
{brown bread} ----> {whole milk}:  conf = 0.389, sup = 0.025
{whipped/sour cream} ----> {whole milk}:  conf = 0.45, sup = 0.032
{whipped/sour cream} ----> {other vegetables}:  conf = 0.403, sup = 0.029
{pork} ----> {whole milk}:  conf = 0.384, sup = 0.022
{pork} ----> {other vegetables}:  conf = 0.376, sup = 0.022
{domestic eggs} ----> {whole milk}:  conf = 0.473, sup = 0.03
{domestic eggs} ----> {other vegetables}:  conf = 0.351, sup = 0.022
{frozen vegetables} ----> {whole milk}:  conf = 0.425, sup = 0.02
{yogurt, whole milk} ----> {other vegetables}:  conf = 0.397, sup = 0.022
{yogurt, other vegetables} ----> {whole milk}:  conf = 0.513, sup = 0.022
{whole milk, root vegetables} ----> {other vegetables}:  conf = 0.474, sup = 0.023

```

```

{other vegetables, root vegetables} ----> {whole milk}:  conf = 0.489, sup = 0.023
{other vegetables, whole milk} ----> {root vegetables}:  conf = 0.31, sup = 0.023
SUPPORT = 0.025 or greater
{citrus fruit} ----> {whole milk}:  conf = 0.369, sup = 0.031
{citrus fruit} ----> {other vegetables}:  conf = 0.349, sup = 0.029
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{yogurt} ----> {other vegetables}:  conf = 0.311, sup = 0.043
{tropical fruit} ----> {other vegetables}:  conf = 0.342, sup = 0.036
{tropical fruit} ----> {whole milk}:  conf = 0.403, sup = 0.042
{pip fruit} ----> {whole milk}:  conf = 0.398, sup = 0.03
{pip fruit} ----> {other vegetables}:  conf = 0.345, sup = 0.026
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
{butter} ----> {whole milk}:  conf = 0.497, sup = 0.028
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{bottled water} ----> {whole milk}:  conf = 0.311, sup = 0.034
{curd} ----> {whole milk}:  conf = 0.49, sup = 0.026
{newspapers} ----> {whole milk}:  conf = 0.343, sup = 0.027
{fruit/vegetable juice} ----> {whole milk}:  conf = 0.368, sup = 0.027
{pastry} ----> {whole milk}:  conf = 0.374, sup = 0.033
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
{sausage} ----> {rolls/buns}:  conf = 0.326, sup = 0.031
{sausage} ----> {whole milk}:  conf = 0.318, sup = 0.03
{brown bread} ----> {whole milk}:  conf = 0.389, sup = 0.025
{whipped/sour cream} ----> {whole milk}:  conf = 0.45, sup = 0.032
{whipped/sour cream} ----> {other vegetables}:  conf = 0.403, sup = 0.029
{domestic eggs} ----> {whole milk}:  conf = 0.473, sup = 0.03
SUPPORT = 0.030000000000000002 or greater
{citrus fruit} ----> {whole milk}:  conf = 0.369, sup = 0.031
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{yogurt} ----> {other vegetables}:  conf = 0.311, sup = 0.043
{tropical fruit} ----> {other vegetables}:  conf = 0.342, sup = 0.036
{tropical fruit} ----> {whole milk}:  conf = 0.403, sup = 0.042
{pip fruit} ----> {whole milk}:  conf = 0.398, sup = 0.03
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{bottled water} ----> {whole milk}:  conf = 0.311, sup = 0.034
{pastry} ----> {whole milk}:  conf = 0.374, sup = 0.033
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
{sausage} ----> {rolls/buns}:  conf = 0.326, sup = 0.031
{whipped/sour cream} ----> {whole milk}:  conf = 0.45, sup = 0.032
SUPPORT = 0.035 or greater
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{yogurt} ----> {other vegetables}:  conf = 0.311, sup = 0.043
{tropical fruit} ----> {other vegetables}:  conf = 0.342, sup = 0.036
{tropical fruit} ----> {whole milk}:  conf = 0.403, sup = 0.042
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
SUPPORT = 0.040000000000000001 or greater
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{yogurt} ----> {other vegetables}:  conf = 0.311, sup = 0.043
{tropical fruit} ----> {whole milk}:  conf = 0.403, sup = 0.042
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
SUPPORT = 0.045000000000000005 or greater
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
SUPPORT = 0.05 or greater
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
SUPPORT = 0.055000000000000001 or greater
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057

```

The provided code again made this section relatively easy! The big thing I wanted to see here was if the rules that were generated were the same, to ensure the implementation was correct. This was the case. With the experimentation around this done it was easy to move onto the next step.

## Part 3

```
In [ ]: F_final, support_data_final = fpgrowth(dataset, min_support=.02)
```

```
In [ ]: rules = generate_rules(F_final, support_data_final, min_confidence=.3)
```

```
{citrus fruit} ----> {whole milk}:  conf = 0.369, sup = 0.031
```



```

{citrus fruit} ----> {other vegetables}:  conf = 0.349, sup = 0.029
{margarine} ----> {whole milk}:  conf = 0.413, sup = 0.024
{yogurt} ----> {whole milk}:  conf = 0.402, sup = 0.056
{yogurt} ----> {other vegetables}:  conf = 0.311, sup = 0.043
{tropical fruit} ----> {other vegetables}:  conf = 0.342, sup = 0.036
{tropical fruit} ----> {whole milk}:  conf = 0.403, sup = 0.042
{pip fruit} ----> {whole milk}:  conf = 0.398, sup = 0.03
{pip fruit} ----> {other vegetables}:  conf = 0.345, sup = 0.026
{other vegetables} ----> {whole milk}:  conf = 0.387, sup = 0.075
{butter} ----> {whole milk}:  conf = 0.497, sup = 0.028
{butter} ----> {other vegetables}:  conf = 0.361, sup = 0.02
{rolls/buns} ----> {whole milk}:  conf = 0.308, sup = 0.057
{bottled water} ----> {whole milk}:  conf = 0.311, sup = 0.034
{curd} ----> {whole milk}:  conf = 0.49, sup = 0.026
{beef} ----> {whole milk}:  conf = 0.405, sup = 0.021
{frankfurter} ----> {whole milk}:  conf = 0.348, sup = 0.021
{newspapers} ----> {whole milk}:  conf = 0.343, sup = 0.027
{fruit/vegetable juice} ----> {whole milk}:  conf = 0.368, sup = 0.027
{pastry} ----> {whole milk}:  conf = 0.374, sup = 0.033
{root vegetables} ----> {other vegetables}:  conf = 0.435, sup = 0.047
{root vegetables} ----> {whole milk}:  conf = 0.449, sup = 0.049
{sausage} ----> {rolls/buns}:  conf = 0.326, sup = 0.031
{sausage} ----> {whole milk}:  conf = 0.318, sup = 0.03
{brown bread} ----> {whole milk}:  conf = 0.389, sup = 0.025
{whipped/sour cream} ----> {whole milk}:  conf = 0.45, sup = 0.032
{whipped/sour cream} ----> {other vegetables}:  conf = 0.403, sup = 0.029
{pork} ----> {whole milk}:  conf = 0.384, sup = 0.022
{pork} ----> {other vegetables}:  conf = 0.376, sup = 0.022
{domestic eggs} ----> {whole milk}:  conf = 0.473, sup = 0.03
{domestic eggs} ----> {other vegetables}:  conf = 0.351, sup = 0.022
{frozen vegetables} ----> {whole milk}:  conf = 0.425, sup = 0.02
{yogurt, whole milk} ----> {other vegetables}:  conf = 0.397, sup = 0.022
{yogurt, other vegetables} ----> {whole milk}:  conf = 0.513, sup = 0.022
{whole milk, root vegetables} ----> {other vegetables}:  conf = 0.474, sup = 0.023
{other vegetables, root vegetables} ----> {whole milk}:  conf = 0.489, sup = 0.023
{other vegetables, whole milk} ----> {root vegetables}:  conf = 0.31, sup = 0.023

```

```

In [ ]: import pandas as pd

df = pd.DataFrame()

```

```

In [ ]: df['rule'] = [f"{'', '.join([i for i in _[0]])} ----> [{}].join([i for i in _[1]])}" for _ in rules]
df['confidence'] = [_[2] for _ in rules]
df['support'] = [support_data_final[frozenset().union(_[1],[0])] for _ in rules]
df['interest factor'] = df['confidence']/[support_data_final[_[1]] for _ in rules]
# df['i factor'] = [support_data_final[frozenset().union(_[1],[0])] / (support_data_final[_[1]] * support_data_final[_[0]])]
# fully coded intrest factor the lift is equivalent in this case

```

```

In [ ]: conf_sort = df.sort_values(by=['confidence'], ascending=False)
conf_sort.head(5)

```

```

Out [ ]:

```

	rule	confidence	support	interest factor
33	[yogurt, other vegetables] ----> [whole milk]	0.512881	0.022267	2.007235
10	[butter] ----> [whole milk]	0.497248	0.027555	1.946053
14	[curd] ----> [whole milk]	0.490458	0.026131	1.919481
35	[other vegetables, root vegetables] ----> [whol...	0.489270	0.023183	1.914833
34	[whole milk, root vegetables] ----> [other vege...	0.474012	0.023183	2.449770

```

In [ ]: supprt_sort = df.sort_values(by=['support'], ascending=False)
supprt_sort.head(5)

```

```

Out [ ]:

```

	rule	confidence	support	interest factor
9	[other vegetables] ----> [whole milk]	0.386758	0.074835	1.513634
12	[rolls/buns] ----> [whole milk]	0.307905	0.056634	1.205032
3	[yogurt] ----> [whole milk]	0.401603	0.056024	1.571735
21	[root vegetables] ----> [whole milk]	0.448694	0.048907	1.756031
20	[root vegetables] ----> [other vegetables]	0.434701	0.047382	2.246605

```

In [ ]: int_sort = df.sort_values(by=['interest factor'], ascending=False)
int_sort.head(5)

```

```

Out [ ]:

```

	rule	confidence	support	interest factor
--	------	------------	---------	-----------------

	rule	confidence	support	interest factor
36	[other vegetables, whole milk] ---> [root vege...	0.309783	0.023183	2.842082
34	[whole milk, root vegetables] ---> [other vege...	0.474012	0.023183	2.449770
20	[root vegetables] ---> [other vegetables]	0.434701	0.047382	2.246605
26	[whipped/sour cream] ---> [other vegetables]	0.402837	0.028876	2.081924

For part 3 to find the intrest factor I was able to substitute lift which was cheaper to calculate as all the values necessary were already calculated and needed less lookups. The reason this is possible is because the transactional data is binary, either the product is purchased or it is not purchased. Moving on to the 5 rules which were generated by sorting the rules by confidence, support, or intrest factor.

- First off when it comes to support, I found it unsurprising that these rules seem to be 1 item to 1 item rules as the number of transactions that contain a pair of items will always be more than a set of 3 items as or more.
- Second, intrestingly there was no rule which overlapped on all 3 tables. There was overlap for support and intrest factor as well as confidence and intrest factor but no support and confidence. Its not supprising as intrest factor is a composite variable made from the confidence and support.

The overlapping rules are as followed:

rule	table 1	table 2
[whole milk, root vegetables] ---> [other vegetables]	confidence	interest factor
[root vegetables] ---> [other vegetables]	support	interest factor