In [1]: import pandas as pd import numpy as np 1. Data Preprocessing 1.1 Raw Data In [2]: raw\_supermarket = pd.read\_excel("AS5 supermarket.xlsx") 1.2 Procedure 1.2.1 Replace In [3]: supermarket = raw\_supermarket.replace({"?":0,"high\n":"high","low\n":"low"}) 1.2.2 Add column names In [4]: names = pd.read\_excel("AS5 supermarket\_attribute\_name.xlsx",header=None)[1] names = [col.strip() for col in names] supermarket.columns = names 1.2.3 Drop columns drop department columns: department col = [col for col in supermarket.columns if col.startswith('department')] In [5]: supermarket = supermarket.drop(columns=department\_col) drop non-frequent columns: supermarket = supermarket.drop(columns=supermarket.columns[supermarket.sum()==0]) drop high / low class column: In [7]: | supermarket = supermarket.drop(columns=['total']) 1.3 Data Info supermarket.info() In [8]: <class 'pandas.core.frame.DataFrame'> RangeIndex: 4627 entries, 0 to 4626 Data columns (total 100 columns): Column Non-Null Count Dtype grocery misc 4627 non-null int64 1 baby needs 4627 non-null int64 bread and cake 4627 non-null int64 int64 baking needs 4627 non-null juice-sat-cord-ms 4627 non-null int64 tea 4627 non-null int64 biscuits 4627 non-null int64 canned fish-meat 4627 non-null int64 canned fruit 4627 non-null int64 canned vegetables 4627 non-null int64 breakfast food 4627 non-null int64 11 cigs-tobacco pkts 4627 non-null int64 12 cigarette cartons 4627 non-null int64 13 cleaners-polishers 4627 non-null int64 coffee 4627 non-null int64 sauces-gravy-pkle 4627 non-null int64 16 confectionary 4627 non-null int64 puddings-deserts 4627 non-null int64 dishcloths-scour 4627 non-null int64 4627 non-null deod-disinfectant int64 frozen foods 4627 non-null int64 21 razor blades int64 4627 non-null 22 fuels-garden aids 4627 non-null int64 4627 non-null spices int64 jams-spreads 4627 non-null int64 25 insecticides 4627 non-null int64 pet foods 4627 non-null int64 27 laundry needs int64 4627 non-null party snack foods 4627 non-null int64 tissues-paper prd 4627 non-null int64 int64 wrapping 4627 non-null 4627 non-null 31 dried vegetables int64 pkt-canned soup int64 4627 non-null soft drinks int64 4627 non-null health food other 4627 non-null int64 beverages hot 4627 non-null int64 health&beauty misc 4627 non-null int64 37 deodorants-soap 4627 non-null int64 mens toiletries int64 4627 non-null 39 medicines 4627 non-null int64 40 haircare 4627 non-null int64 dental needs 4627 non-null int64 lotions-creams 4627 non-null int64 sanitary pads 4627 non-null int64 43 cough-cold-pain 4627 non-null int64 meat misc int64 4627 non-null 46 cheese 4627 non-null int64 4627 non-null chickens int64 47 4627 non-null milk-cream int64 cold-meats 4627 non-null int64 49 deli gourmet 4627 non-null int64 margarine 51 4627 non-null int64 52 salads 4627 non-null int64 small goods 4627 non-null int64 dairy foods 4627 non-null int64 fruit drinks 4627 non-null int64 delicatessen misc 4627 non-null int64 int64 57 beef 4627 non-null 4627 non-null int64 58 hogget lamb 4627 non-null int64 4627 non-null int64 pet food 61 4627 non-null int64 pork poultry 4627 non-null int64 veal 4627 non-null int64 gourmet meat 4627 non-null int64 produce misc 4627 non-null int64 fruit int64 66 4627 non-null plants 67 4627 non-null int64 4627 non-null int64 potatoes vegetables 4627 non-null int64 70 variety misc 4627 non-null int64 71 brushware 4627 non-null int64 4627 non-null electrical int64 haberdashery 4627 non-null int64 kitchen int64 4627 non-null manchester int64 4627 non-null pantyhose 4627 non-null int64 plasticware 4627 non-null int64 int64 stationary 4627 non-null prepared meals 4627 non-null int64 preserving needs 4627 non-null int64 condiments int64 4627 non-null cooking oils 4627 non-null int64 bake off products int64 4627 non-null small goods2 4627 non-null int64 85 offal 4627 non-null int64 4627 non-null int64 86 mutton 87 trim pork 4627 non-null int64 88 trim lamb 4627 non-null int64 89 imported cheese 4627 non-null int64 90 casks white wine 4627 non-null int64 91 casks red wine 4627 non-null int64 92 750ml white nz 4627 non-null int64 93 750ml red nz 4627 non-null int64 4627 non-null 94 750ml white imp int64 95 750ml red imp 4627 non-null int64 96 sparkling nz 4627 non-null int64 4627 non-null 97 sparkling imp int64 98 port and sherry int64 4627 non-null 99 non host support 4627 non-null int64 dtypes: int64(100) memory usage: 3.5 MB 2. Data Mining 2.1 Apriori Algorithm 2.1.1 Apriori class In [9]: class Apriori: def \_\_init\_\_(self, min\_sup = 0.2): self.min sup = min sup # min support \*Generating frequent 1-itemsets\* in: data<pandas.core.frame.DataFrame> out: C<list>[<tuple>], D<dict> def frequent singleton(self, data): sup\_count = data.sum() sup\_count = sup\_count[sup\_count>=self.min\_sup\_count] sup count.index = [(idx,) for idx in sup count.index] C = sup count.index D = sup\_count.to\_dict() return C, D \*Generating candidates\* in: l<list>[<tuple>], k<int> out: <list>[<tuple>] def \_apriori\_gen(self, l, k): # 1:frequent (k-1)-itemsets C = set()for i in range(len(1)-1): for j in range(i+1,len(l)): c = tuple(sorted(set(l[i]).union(set(l[j])))) if (k == 2) or (len(c) == k) and not (self.\_has\_infrequent\_subset(c, l, k-1)): C.add(c) return list(C) 1 1 1 \*Pruning\* in: c<tuple>, l<list>[<tuple>], k<int> out: <boolean> def \_has\_infrequent\_subset(self, c, l, k): # c: candidate k-itemset, L:frequent (k-1)-itemsets S = self.\_ksubsets(set(c), k) for s in S: if tuple(sorted(s)) not in 1: return True return False \*Generating all subsets of c with k elements\* in: c<set>, k<int> out: <set> def \_ksubsets(self, c, k): import itertools return set(itertools.combinations(c, k)) \*Mining all frequent itemsets\* in: data<pandas.core.frame.DataFrame> out: L<list>[<list>[<tuple>]] def mine(self, data): import time start\_time = time.time() self.min\_sup\_count = int(self.min\_sup\*data.shape[0])+1 C, D= self.\_frequent\_singleton(data) k = 1while len(C) > 0: C = self.\_apriori\_gen(C, k+1) for i in range(data.shape[0]): transaction = data.loc[i] s = set(transaction[transaction>0].index) for c in C: if set(c).issubset(s): D[c] = D.get(c,0)+1C = [c for c in C if D.get(c,0) >= self.min\_sup\_count] k += 1D = {key:val for key,val in D.items() if val >= self.min sup count} elapsed\_time = time.time() - start\_time print('Time elapsed:', elapsed time) return D \*Showing top 10 most frequent 2,3,4-itemsets\* in: D<dict> def display(self, D): for i in range(2,5): print('Top 10 most frequent ', i, '-itemset:', sep='') L = [item for item in D.items() if len(item[0])==i] L = sorted(L, key = lambda item: float(item[1]), reverse = True) for j in range(10): print(j+1,': ', L[j], sep='') 2.1.2 Pattern mining In [10]: apriori = Apriori(min sup=0.2) D = apriori.mine(supermarket) Time elapsed: 32.38124108314514 2.1.3 Frequent patterns In [11]: apriori.display(D) Top 10 most frequent 2-itemset: 1: (('bread and cake', 'milk-cream'), 2337) 2: (('bread and cake', 'fruit'), 2325) 3: (('bread and cake', 'vegetables'), 2298) 4: (('fruit', 'vegetables'), 2207) 5: (('baking needs', 'bread and cake'), 2191) 6: (('bread and cake', 'frozen foods'), 2129) 7: (('biscuits', 'bread and cake'), 2083) 8: (('fruit', 'milk-cream'), 2038) 9: (('milk-cream', 'vegetables'), 2025) 10: (('baking needs', 'vegetables'), 1949) Top 10 most frequent 3-itemset: 1: (('bread and cake', 'fruit', 'vegetables'), 1791) 2: (('bread and cake', 'fruit', 'milk-cream'), 1684) 3: (('bread and cake', 'milk-cream', 'vegetables'), 1658) 4: (('baking needs', 'bread and cake', 'vegetables'), 1586) 5: (('baking needs', 'bread and cake', 'milk-cream'), 1580) 6: (('fruit', 'milk-cream', 'vegetables'), 1571) 7: (('baking needs', 'bread and cake', 'fruit'), 1564) 8: (('bread and cake', 'frozen foods', 'vegetables'), 1548) 9: (('bread and cake', 'frozen foods', 'fruit'), 1548) 10: (('biscuits', 'bread and cake', 'fruit'), 1541) Top 10 most frequent 4-itemset: 1: (('bread and cake', 'fruit', 'milk-cream', 'vegetables'), 1311) 2: (('baking needs', 'bread and cake', 'fruit', 'vegetables'), 1255) 3: (('bread and cake', 'frozen foods', 'fruit', 'vegetables'), 1242) 4: (('biscuits', 'bread and cake', 'fruit', 'vegetables'), 1216) 5: (('baking needs', 'bread and cake', 'milk-cream', 'vegetables'), 1169) 6: (('baking needs', 'bread and cake', 'fruit', 'milk-cream'), 1161) 7: (('biscuits', 'bread and cake', 'frozen foods', 'fruit'), 1143) 8: (('bread and cake', 'frozen foods', 'fruit', 'milk-cream'), 1138) 9: (('bread and cake', 'frozen foods', 'milk-cream', 'vegetables'), 1131) 10: (('biscuits', 'bread and cake', 'fruit', 'milk-cream'), 1122) 2.2 FP-growth Algorithm 2.2.1 Header class In [12]: class Node: def \_\_init\_\_(self, item, sup): self.item = item self.sup = sup self.adjacent = None self.tail = None 2.2.2 Tree class In [13]: class treeNode: def \_\_init\_\_(self, item, count=1): self.item = item self.count = count self.parent = None self.next = None self.child = set() def single\_path(self): temp = self while len(temp.child)>0: if len(temp.child)>1: return False temp = tuple(temp.child)[0] return True def insert\_tree(self, header, transaction, i, count = 1): if i >= len(transaction): return head = **None** for node in header: if node.item == transaction[i]: head = nodebreak if head == None: return for c in self.child: if c.item == transaction[i]: c.count += count c.insert\_tree(header,transaction, i+1, count) return c = self.\_\_class\_\_(transaction[i],count) c.parent = self self.child.add(c) if head.adjacent == None: head.tail = head.adjacent = c else: head.tail.next = c head.tail = head.tail.next c.insert tree(header, transaction, i+1, count) return 2.2.3 FP-tree class In [14]: class FP\_tree: def \_\_init\_\_(self,header, tree): self.header = header self.tree = tree \*Generating conditional pattern base\* in: item<str> out: pattern base<dict> def \_cond\_pattern\_base(self, item): pattern base = dict() for node in self.header: if node.item == item: break head = node.adjacent while True: temp = headl = list()while temp.parent.item != '': l.append(temp.parent.item) temp = temp.parent if len(1) > 0: l.reverse() pattern base[tuple(1)] = head.count if head.next == None: break head = head.next return pattern\_base \*Generating conditional FP-tree\* in: item<str>, min\_sup\_count<int> out: <FP-tree> 1 1 1 def cond\_FP\_tree(self, item, min\_sup\_count): cpb = self.\_cond\_pattern\_base(item) root = treeNode('',0) header = dict() for itemset in cpb: for item in itemset: header[item] = header.get(item,0) + cpb[itemset] header = [Node(key,val) for key,val in header.items() if val >= min\_sup\_count] for itemset in cpb: root.insert\_tree(header, itemset, 0, cpb[itemset]) return FP\_tree(header, root) 2.2.4 FP-growth class In [15]: class FP growth: def \_\_init\_\_(self, min\_sup = 0.2): self.min sup = min sup \*Generating frequent 1-itemsets\* in: data<pandas.core.frame.DataFrame> out: header<list>[<Node>], D<dict> 1 1 1 def \_frequent\_item(self, data): sup count = data.sum() sup\_count = sup\_count[sup\_count >= self.min\_sup\_count] sup count = sup count.sort values(ascending=False) header = [Node(item, sup\_count[item]) for item in sup\_count.index] sup\_count.index = [(idx,) for idx in sup\_count.index] D = sup\_count.to\_dict() return header, D \*Constructing tree\* in: header<list>[<Node>], data<pandas.core.frame.DataFrame> out: root<treeNode> 1 1 1 def \_construct\_tree(self, header, data): root = treeNode('',0) for i in range(data.shape[0]): transaction = data.loc[i] transaction = [node.item for node in header if transaction[node.item] > 0] root.insert\_tree(header, transaction, 0) return root \*Implementing FP-growth\* in: fp\_tree<FP-tree>, pattern<list> out: D<dict> def \_FP\_growth(self, fp\_tree, pattern): import itertools D = dict()if fp tree.tree.single path(): temp = fp tree.tree l = list()d = dict()while len(temp.child) > 0: c = tuple(temp.child)[0] l.append(c.item) d[c.item] = c.count temp = cfor i in range(0,len(1)): combination = list(itertools.combinations(1,i+1)) for itemset in combination: sup\_count = d[itemset[-1]] itemset = list(itemset)+ pattern D[tuple(itemset)] = sup\_count else: for node in fp\_tree.header: cfpt = fp\_tree.cond\_FP\_tree(node.item, self.min\_sup\_count) D[tuple([node.item]+pattern)] = node.sup if len(cfpt.tree.child) > 0: D = {\*\*D, \*\*self. FP growth(cfpt, [node.item]+pattern)} return D \*Mining all frequent itemsets\* in: data<pandas.core.frame.DataFrame> out: <dict> def mine(self,data): import time start time = time.time() self.min\_sup\_count = int(self.min\_sup\*data.shape[0])+1 header, self.D = self. frequent item(data) # D may be unnecessary tree = self. construct\_tree(header,data) fp\_tree = FP\_tree(header, tree) self.D = {\*\*self.D, \*\*self.\_FP\_growth(fp\_tree,[])} elapsed time = time.time() - start time print('Time elapsed:', elapsed\_time) return self.D \*Showing top 10 most frequent 2,3,4-itemsets\* in: D<dict> 1.1.1 def display(self,D): for i in range(2,5): print('Top 10 most frequent ', i, '-itemset:',sep='') L = [(tuple(sorted(item[0])),item[1]) for item in D.items() if len(item[0])==i] L = sorted(L, key = lambda item: float(item[1]), reverse = True) for j in range(10): print(j+1,': ', L[j], sep='') 2.2.5 Pattern mining In [16]: fp\_growth = FP\_growth(min\_sup=0.2) D = fp growth.mine(supermarket) Time elapsed: 5.522256135940552 2.2.6 Frequent patterns In [17]: | fp\_growth.display(D) Top 10 most frequent 2-itemset: 1: (('bread and cake', 'milk-cream'), 2337) 2: (('bread and cake', 'fruit'), 2325) 3: (('bread and cake', 'vegetables'), 2298) 4: (('fruit', 'vegetables'), 2207) 5: (('baking needs', 'bread and cake'), 2191) 6: (('bread and cake', 'frozen foods'), 2129) 7: (('biscuits', 'bread and cake'), 2083) 8: (('fruit', 'milk-cream'), 2038) 9: (('milk-cream', 'vegetables'), 2025) 10: (('baking needs', 'vegetables'), 1949) Top 10 most frequent 3-itemset: 1: (('bread and cake', 'fruit', 'vegetables'), 1791) 2: (('bread and cake', 'fruit', 'milk-cream'), 1684) 3: (('bread and cake', 'milk-cream', 'vegetables'), 1658) 4: (('baking needs', 'bread and cake', 'vegetables'), 1586) 5: (('baking needs', 'bread and cake', 'milk-cream'), 1580) 6: (('fruit', 'milk-cream', 'vegetables'), 1571) 7: (('baking needs', 'bread and cake', 'fruit'), 1564) 8: (('bread and cake', 'frozen foods', 'fruit'), 1548) 9: (('bread and cake', 'frozen foods', 'vegetables'), 1548) 10: (('biscuits', 'bread and cake', 'fruit'), 1541) Top 10 most frequent 4-itemset: 1: (('bread and cake', 'fruit', 'milk-cream', 'vegetables'), 1311) 2: (('baking needs', 'bread and cake', 'fruit', 'vegetables'), 1255) 3: (('bread and cake', 'frozen foods', 'fruit', 'vegetables'), 1242) 4: (('biscuits', 'bread and cake', 'fruit', 'vegetables'), 1216) 5: (('baking needs', 'bread and cake', 'milk-cream', 'vegetables'), 1169) 6: (('baking needs', 'bread and cake', 'fruit', 'milk-cream'), 1161)

7: (('biscuits', 'bread and cake', 'frozen foods', 'fruit'), 1143)
8: (('bread and cake', 'frozen foods', 'fruit', 'milk-cream'), 1138)

10: (('biscuits', 'bread and cake', 'fruit', 'milk-cream'), 1122)

3. Comparison

The end

9: (('bread and cake', 'frozen foods', 'milk-cream', 'vegetables'), 1131)

under min sup =0.2, FP-growth algorithm (~6s) runs almost five times faster than Apriori algorithm (>30s).

Different from the multiple scans of data in Apriori, FP-growth only scans the data twice and boosts the speed.

**CSC4008** Assignment 5

@arthor: JIANG, Jingxin 117020119

0. Macros