

数据科学基础第二次实践

1. 考查内容

2. 实验说明

4. 代码设计与实现

4.1 数据处理

4.2 Dummy Apriori

4.3 Apriori1

4.4 Apriori2

4.5 Apriori3

4.6 FP_Growth算法

4.7 calculate_association_rules

5. 实验结果与分析

5.1 使用 apriori 算法在 Groceries数据集上挖掘频繁 3-项集(支持度 0.01)

5.2 使用 apriori 算法在 Groceries数据集上挖掘关联规则(支持度 0.01, 置信度 0.5)

5.3 使用 FP-Growth 算法在 Groceries 数据集上挖掘的一些关联规则(置信度 0.5)

5.4 对比 dummy apriori、仅使用第一种剪枝策略 advanced apriori、同时 使用第一种和第二种剪枝策略 a...

5.5 使用命令行接收参数(数据集、支持率、挖掘频繁项集的大小)

5.6 寻找 2 个关联规则, 比较强的关联规则, 进行市场分析

5.7 实现第四部分第三种剪枝策略, 并将三种剪枝策略一起使用算法运行时间一起画在算法时 间对比图表中

5.8

apriori1 的内存使用情况

apriori2 的内存使用情况

apriori3 的内存使用情况

fpgrowth 的内存使用情况

6. 运行文件说明

1. 考查内容

频繁模式与关联规则挖掘, 考察对 `apriori` 与 `fpgrowth` 算法的掌握程度。

2. 实验说明

1. 在给定数据集上(10,000 量级的数据)上使用关联规则挖掘算法。
2. 通过改变不同等级的支持度和置信度，比较 `apriori`、`fpgrowth` 和 `baseline` 算法(说明见)的性能。
3. 试图根据 `apriori` 或 `fpgrowth` 算法揭露出的物品集相关性，发现一些有趣的关联规则。

4. 代码设计与实现

4.1 数据处理

将给定的csv文件处理成字典列表形式，每一个list代表是一次事务。

```
1 def load_data(path='./data/Groceries.csv'):
2     r"""
3     load info from origin data
4     generate the itemset and c_1.
5     """
6     itemset = []
7     goods = {}
8     with open(path, 'r') as f:
9         for i, line in tqdm(enumerate(f.readlines())):
10             if i == 0:
11                 continue
12             words = line.replace('{', '').replace('}', '').replac
13 e(
14             '', '').strip().split(',')[1:]
15             item = {}
16             if len(words) > 0:
17                 for word in words:
18                     word = word.replace(' ', '').replace('/', '')
19                     goods[word] = 1
20                     item[word] = 1
21                 itemset.append(item)
22             c_1 = []
23             for word in list(goods.keys()):
24                 c_1.append(set([word]))
25     return itemset, c_1
```

4.2 Dummy Apriori

没有进行剪枝，即用暴力搜索的方法根据 k 项频繁项集生成候选 $k+1$ 项集，对事务表不做任何处理。

下面对该算法进行一个简单的 `review`

Step1. 生成候选集

```
1 def generate_c_k_plus_1(l_ks):
2     """ generate C_{k+1} 候选集
3     l_ks: a set of c_k
4     """
5     c_k_plus_1s = {}
6     for i in tqdm(range(0, len(l_ks)),
7                   desc=f'Generate the c_k_plus_1 {len(l_ks[0])+1}
dataset', mininterval=0.1):
8         for j in range(i + 1, len(l_ks)):
9             # 去重复
10            item = sorted(l_ks[i].union(l_ks[j]))
11            c_k_plus_1s[tuple(item)] = 1
12        logger.info(f'get the c_{k+1} {len(l_ks[0])+1}:{len(c_k_plus_1s
)}}')
13    return [set(c_k_1) for c_k_1 in list(c_k_plus_1s.keys())]
```

Step2. 根据候选集过滤出符合支持度的频繁集

支持度的计算方法

```
1 def support_rate(itemset, c_k, num):
2     """ support rate
3     支持度
4     itemset: 数据集
5     c_k, the k-th 候选集
6     """
7     count = 0.0
8     record_nums = [0 for _ in range(len(itemset))]
9     for i, item in enumerate(itemset):
```

```

10         hit = 0
11         for c in c_k:
12             if c in item:
13                 hit += 1
14         if len(c_k) == hit:
15             count += 1
16             record_nums[i] = 1
17     return count / num, record_nums

```

根据候选集过滤出符合支持度的频繁集

```

1 def generate_l_k(itemset, c_ks, support, num, counter, trick=None
  ):
2     """ Generate l_k from c_k
3     trick: trick or not
4     Return:
5         候选集、支持度、新的事务表（可能会使用trick）
6     """
7     l_k = []
8     sp_s = []
9     k = len(c_ks[0])
10    desc = f'Generate l_k k={k}'
11    all_record_nums = [0 for _ in range(len(itemset))]
12    for c_k in tqdm(c_ks, desc=desc, mininterval=0.1):
13        sp, record_nums = support_rate(itemset, c_k, num)
14        for i in range(len(record_nums)):
15            all_record_nums[i] += record_nums[i]
16        if sp > support:
17            l_k.append(c_k)
18            sp_s.append(sp)
19
20    # Advanced aprorio 2 filter
21    if trick is None:
22        logger.info('dummy algorithm')
23    elif trick == 'trick_2':
24        itemset = trick_2(all_record_nums, itemset, k)
25    elif trick == 'trick_3':
26        itemset = trick_3(counter, k, itemset)

```

```

27
28     logger.info(f' get l_k len:{len(l_k)} ')
29     return l_k, sp_s, itemset

```

Step3. 找出所有的频繁集，计算其关联规则，找到符合要求的关联规则

```

1  args = opt()
2  all_lks, all_sps, interval = main(args)
3  k_set = get_k_set(all_lks, k=args.k)
4  logger.info(
5      f"len(all_lks):{len(all_lks)}, len(k_set):{len(k_set)}, k_
    set:{k_set}")
6  rules = get_association_rules(
7      all_lks, all_sps, args.confidence, model=args.model)
8  print(f"len(rules):{len(rules)}, rules:{rules}")

```

4.3 Apriori1

主要作用为：减小生成候选项集规模。如果一个 $k+1$ 候选项是频繁的，那么生成它的 \square 频繁项集中必包含其 $k+1$ 个子集。例如，如果 $\{t1, t2, t3, t5\}$ 是频繁的，那么它将由 $\{t1, t2, t3\}$ 和 $\{t1, t2, t5\}$ 生成，否则它就不会被生成。这么做要求生成候选项集的频繁项集的项内是有序的，各项之间也是有序的。

```

1  def advanced_apriori1(itemset, c_ks, support, num):
2      all_lks = []
3      all_sps = []
4      counter = count_freq(itemset)
5
6      def apriori(itemset, c_k, support):
7          l_ks, sp_s, itemset = generate_l_k(itemset, c_k, support,
            num, counter)
8          all_lks.extend(l_ks)
9          all_sps.extend(sp_s)
10         if len(l_ks) > 1:
11             ck_plus_1s = generate_c_k_plus_1(l_ks)
12             # 类间排序

```

```

13         ck_plus_1s = sorted(ck_plus_1s)
14         apriori(itemset, ck_plus_1s, support)
15
16     apriori(itemset, c_ks, support)
17
18     return all_lks, all_sps

```

4.4 Apriori2

主要作用是：减小事务表(数据记录表)规模。

如果一个 $k+1$ 项候选项能与事务表中一条数据记录匹配，那么其 $k+1$ 个子集也必能与事务表中该条记录匹配。因此在事务表中匹配 k 候选项集的时候，统计每条记录被匹配到的次数，如果少于 $k+1$ 次，那么将该条记录从事务表中移除，因为它绝不可能与下一轮的任一 $k+1$ 项候选项匹配。这样每次迭代都减小了事务表的规模，从而减小扫描事务表的时间消耗。

```

1
2 def trick_2(all_record_nums, itemset, k):
3     """
4     trick_2: 减小事务表(数据记录表)规模。
5
6     """
7     dataset = []
8     for record_num, item in zip(all_record_nums, itemset):
9         if record_num > k:
10             dataset.append(item)
11     logger.info(f"before: {len(itemset)}, after: {len(dataset)}")
12     return dataset

```

4.5 Apriori3

主要作用：减少事务表中元组的项

如果事务表中元组的某一项能包含在一个 $k+1$ 频繁项中，那么该项必出现在这个 $k+1$ 频繁项的 k 个 k 项子集中。所以在扫描事务表统计 k 项候选集的出现频次时，如果事务表中任一元组的某一项未被匹配中 k 次，那么该项将在筛选出 k 频繁项集后从元组中除去，从而减少统计 $k+1$ 项候选集频次时与事务表中元组的匹配次数

```

1 # 首先统计一次全部元素的频率
2 def count_freq(itemset):
3     """ Count freq of item in itemset
4     itemset: dataset
5     """
6     counter = {}
7     for item in tqdm(itemset):
8         for record in item:
9             if record in counter:
10                 counter[record] += 1
11             else:
12                 counter[record] = 1
13
14     return counter
15
16
17 def trick_3(counter, k, itemset):
18     """
19     trick_3 : 减少事务表中元组的项。
20
21     Parameters
22     -----
23     counter : [dict]
24         每一个元素的词频
25     k : [int]
26         候选集元素的个数
27     itemset : [list]
28         事务表
29     """
30     dataset = []
31     for item in itemset:
32         elements = copy.deepcopy(item)
33         for element in item.keys():
34             if counter[element] <= k:
35                 del elements[element]
36         if len(elements) > 0:
37             dataset.append(elements)
38     return dataset

```

4.6 FP_Growth算法

这个通过调包实现，比较简单。使用时为了方便调用，做了一下封装。

```
1 def get_association_rules(all_lks, all_sps, confidence, model='fp'):  
2     if model == 'fpgrowth':  
3         return generate_association_rules(all_lks, confidence)  
4     elif model == 'apriori':  
5         return calculate_association_rules(all_lks, all_sps, confidence)  
6     else:  
7         raise(  
8             f"Not support algorithm {model}, Please choose from ['fpgrowth', 'apriori']")
```

4.7 calculate_association_rules

计算其置信度，找出其关联规则。这个是自己实现的一个查找关联规则的方法，但是在与pyfpgrowth包挖掘关联规则的结果对比中，发现结果存在不一致的问题。

```
1 def calculate_association_rules(all_lks, all_sps, confidence):  
2     counter = {}  
3  
4     # Step1. 将 list to map  
5     for lk, sp in tqdm(zip(all_lks, all_sps)):  
6         counter[tuple(sorted(lk))] = sp  
7     rules = []  
8     filters = {}  
9     full_rules = []  
10    for lk in tqdm(all_lks):  
11        # print(type(lk))  
12        for i in range(1, len(lk)):  
13            combines = list(combinations(lk, i))  
14            for a in combines:  
15                a_s = tuple(sorted(a))  
16                b = list(set(lk).difference(set(a)))  
17                b_s = tuple(sorted(b))
```



```

18         key = tuple(sorted(lk))
19         max_info = None
20         rule = None
21         max_confidence = 0
22         if a_s in counter and (counter[key] / counter[a_s
    ]) > confidence and (counter[key] / counter[a_s]) > max_confidenc
    e:
23             info = f'{a_s}:{b_s}, {(counter[key] / counte
    r[a_s])}'
24             max_confidence = (counter[key] / counter[a_s]
    )
25             if info not in filters:
26                 max_info = info
27                 filters[info] = 1
28                 rule = (a_s, b_s, max_confidence)
29             if b_s in counter and (counter[key] / counter[b_s
    ]) > confidence and (counter[key] / counter[b_s]) > max_confidenc
    e:
30                 info = f'{b_s}:{a_s}, {(counter[key] / count
    er[b_s])}'
31                 max_confidence = (counter[key] / counter[b_s]
    )
32                 if info not in filters:
33                     max_info = info
34                     filters[info] = 1
35                     rule = (a_s, b_s, max_confidence)
36             if max_info is not None:
37                 rules.append(rule)
38                 full_rules.append(lk)
39         # print(full_rules)
40         return rules

```

5. 实验结果与分析

下面要求均已完成

5.1 使用 apriori 算法在 Groceries数据集上挖掘频繁 3-项集(支持度 0.01)

使用 Apriori2算法得到的结果如下:

```
1 len(all_lks):333, len(k_set):32, k_set:[{'wholemilk', 'yogurt', 'citrusfruit'}, {'othervegetables', 'wholemilk', 'citrusfruit'}, {'othervegetables', 'rootvegetables', 'citrusfruit'}, {'wholemilk', 'yogurt', 'tropicalfruit'}, {'othervegetables', 'yogurt', 'tropicalfruit'}, {'othervegetables', 'wholemilk', 'tropicalfruit'}, {'wholemilk', 'tropicalfruit', 'rollsbuns'}, {'rootvegetables', 'wholemilk', 'tropicalfruit'}, {'othervegetables', 'rootvegetables', 'tropicalfruit'}, {'othervegetables', 'wholemilk', 'yogurt'}, {'wholemilk', 'yogurt', 'rollsbuns'}, {'curd', 'wholemilk', 'yogurt'}, {'yogurt', 'wholemilk', 'soda'}, {'rootvegetables', 'wholemilk', 'yogurt'}, {'whippedsourcream', 'yogurt', 'wholemilk'}, {'othervegetables', 'yogurt', 'rollsbuns'}, {'othervegetables', 'rootvegetables', 'yogurt'}, {'othervegetables', 'whippedsourcream', 'yogurt'}, {'othervegetables', 'wholemilk', 'pipfruit'}, {'othervegetables', 'butter', 'wholemilk'}, {'othervegetables', 'wholemilk', 'rollsbuns'}, {'othervegetables', 'wholemilk', 'bottledwater'}, {'othervegetables', 'wholemilk', 'soda'}, {'othervegetables', 'fruitvegetablejuice', 'wholemilk'}, {'othervegetables', 'wholemilk', 'pastry'}, {'othervegetables', 'rootvegetables', 'wholemilk'}, {'othervegetables', 'sausage', 'wholemilk'}, {'othervegetables', 'wholemilk', 'pork'}, {'othervegetables', 'whippedsourcream', 'wholemilk'}, {'othervegetables', 'wholemilk', 'domestic eggs'}, {'rootvegetables', 'wholemilk', 'rollsbuns'}, {'othervegetables', 'rootvegetables', 'rollsbuns'}]
```

5.2 使用 apriori 算法在 Groceries数据集上挖掘关联规则(支持度 0.01, 置信度 0.5)

使用 Apriori2算法得到的结果如下

```
1 len(rules):14, rules:[(('othervegetables',), ('citrusfruit', 'rootvegetables'), 0.5862068965517241), (('wholemilk',), ('tropicalfruit', 'yogurt'), 0.5173611111111111), (('wholemilk',), ('rootvegetables', 'tropicalfruit'), 0.570048309178744), (('othervegetables',), ('rootvegetables', 'tropicalfruit'), 0.5845410628019324), (('wholemilk',), ('othervegetables', 'yogurt'), 0.5128805620608898), (('wh
```

```

('wholemilk',), ('curd', 'yogurt'), 0.5823529411764706), (('wholemilk',), ('rootvegetables', 'yogurt'), 0.562992125984252), (('wholemilk',), ('whippedsourcream', 'yogurt'), 0.5245098039215685), (('wholemilk',), ('othervegetables', 'pipfruit'), 0.5175097276264592), (('wholemilk',), ('butter', 'othervegetables'), 0.5736040609137055), (('wholemilk',), ('othervegetables', 'whippedsourcream'), 0.507042535211268), (('wholemilk',), ('domesticeggs', 'othervegetables'), 0.5525114155251142), (('wholemilk',), ('rollsbuns', 'rootvegetables'), 0.5230125523012552), (('othervegetables',), ('rollsbuns', 'rootvegetables'), 0.502092050209205)]

```

5.3 使用 FP-Growth 算法在 Groceries 数据集上挖掘的一些关联规则(置信度 0.5)

```

1 len(rules):7, rules:{('curd', 'yogurt'): (('wholemilk',), 0.5823529411764706), ('butter', 'othervegetables'): (('wholemilk',), 0.5736040609137056), ('domesticeggs', 'othervegetables'): (('wholemilk',), 0.5525114155251142), ('othervegetables', 'whippedsourcream'): (('wholemilk',), 0.5070422535211268), ('othervegetables', 'pipfruit'): (('wholemilk',), 0.5175097276264592), ('citrusfruit', 'rootvegetables'): (('othervegetables',), 0.5862068965517241), ('othervegetables', 'yogurt'): (('wholemilk',), 0.5128805620608899)}

```

5.4 对比 dummy apriori、仅使用第一种剪枝策略 advanced apriori、同时使用第一种和第二种剪枝策略 advanced apriori 的时间损耗

```

1 finish algorithm with 110.9324722290039s # dummy apriori
2 finish algorithm with 112.84710025787354s # apriori1
3 finish algorithm with 71.55795574188232s # apriori2
4 finish algorithm with 106.52765893936157s # apriori3
5 finish algorithm with 0.8581380844116211s # FP-Growth

```

5.5 使用命令行接收参数(数据集、支持率、挖掘频繁项集的大小)

使用方法, 在命令行中使用:

```
1 python main.py -support=0.01 -confidence=0.5 -k=3 -path='' -ty='dummy'
```

各个参数的含义如下:

```
1 parser = argparse.ArgumentParser(  
2     description='The association rule mining')  
3 # 支持率  
4 parser.add_argument('-support', type=float,  
5                     default=0.01, help='The support rate')  
6 # 置信度  
7 parser.add_argument('-confidence', type=float,  
8                     default=0.5, help='The confidence rate')  
9 # 挖掘频繁项集的大小  
10 parser.add_argument('-k', type=int, default=3,  
11                    help='The size of k_item_freq')  
12 # 数据集路径  
13 parser.add_argument('-path', type=str,  
14                    default='./data/Groceries.csv',  
15                    help='The path of dataset')  
16 # 算法选择  
17 parser.add_argument('-ty', type=str,  
18                    default='dummy',  
19                    help='The type of algorithm, must be one  
    of [`dummy`, `apriori1`, `apriori2`, `apriori3`, `fpgrowth`])
```

5.6 寻找 2 个关联规则, 比较强的关联规则, 进行市场分析

('citrusfruit', 'rootvegetables'): (('othervegetables',), 0.5862068965517241)

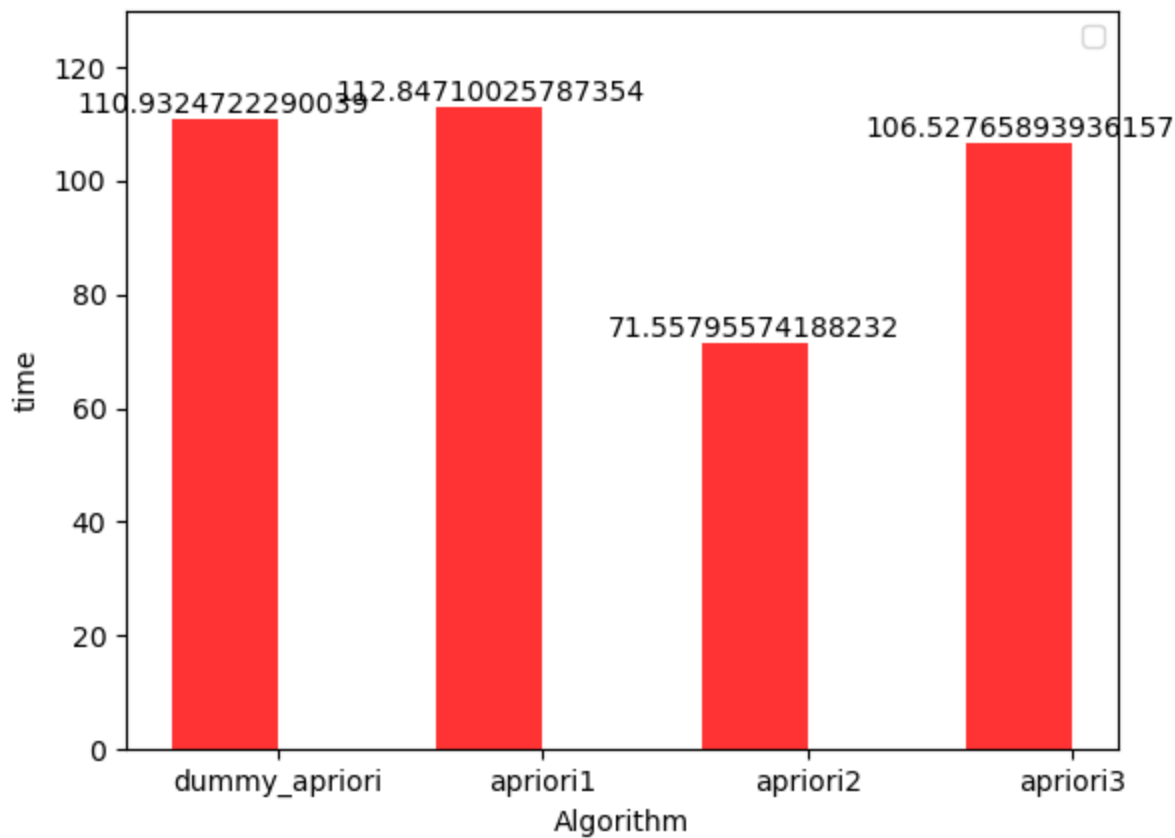
对于素食主义者或者节食的人, 可能在看水果和蔬菜时, 也会考虑其他品种的蔬菜

('curd', 'yogurt'): (('wholemilk',), 0.5823529411764706)

'凝乳', '酸奶', '全脂牛奶', 这三者都是乳制品, 在购买纯牛奶时可以会考虑酸奶或者凝乳。

5.7 实现第四部分第三种剪枝策略，并将三种剪枝策略一起使用算法运行时间一起画在算法时间对比图表中

apriori3 算法即是第三种剪枝方案。这四种方案（baseline与三种剪枝）的算法运行时间对比图如下：



从数据集的表现来看，advanced_2 apriori 性能是最好的，advanced_1 apriori性能与 advanced_3 apriori性能相近。

5.8

dummy apriori 的内存使用情况

Line #	Mem usage	Increment	Occurences	Line Contents
130	39.012 MiB	39.012 MiB	1	@profile

apriori1 的内存使用情况

Line #	Mem usage	Increment	Occurrences	Line Contents
130	39.062 MiB	39.062 MiB	1	@profile
131				def opt():
132	39.062 MiB	0.000 MiB	1	ccscs = csccs

apriori2 的内存使用情况

Line #	Mem usage	Increment	Occurrences	Line Contents
130	39.027 MiB	39.027 MiB	1	@profile
131				def opt():

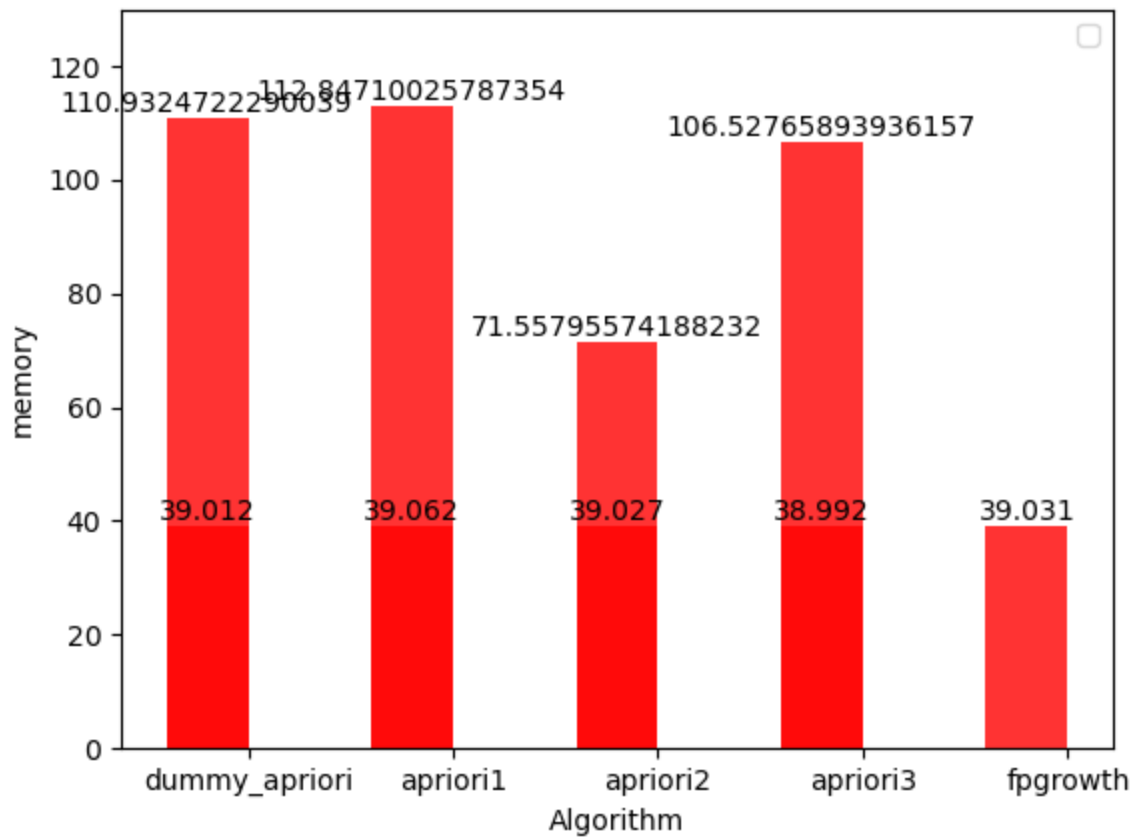
apriori3 的内存使用情况

Line #	Mem usage	Increment	Occurrences	Line Contents
130	38.992 MiB	38.992 MiB	1	@profile

fpgrowth 的内存使用情况

Line #	Mem usage	Increment	Occurrences	Line Contents
130	39.031 MiB	39.031 MiB	1	@profile
131				def opt():

好像没看出啥差别，可能我实现的方法不对



6. 运行文件说明

main.py 为主函数，运行该文件是挖掘频繁集和关联规则，使用方法见5.6说明

draw.py 是对一些运行数据的图形绘制

out.log 是运行时的一些日志