数据挖掘互评作业二: 频繁模式与关联规则挖掘

计算机学院 JohnsonGUO

代码仓库地址: https://github.com/johnsongwx/data-mining-coursework

数据挖掘互评作业二: 频繁模式与关联规则挖掘

作业题目选择

网页浏览行为关联规则挖掘

- 1. 数据预处理
- 2. 频繁模式挖掘
- 3. 关联规则挖掘, 计算其支持度和置信度
- 4. 规则评价, 结果评估
- 5. 结果可视化

作业题目选择

网页浏览行为关联规则挖掘

任务: 分析用户在网站上的浏览行为数据, 挖掘潜在的跳转规律, 为网站提供优化导航结构的建议。

数据集: UCI ML Repository - Anonymous Microsoft Web Data

关联规则挖掘:

- 1. 数据预处理: 清洗数据, 处理缺失值, 提取用户浏览记录。
- 2. 数据探索性分析: 分析最常被访问的页面、页面访问量分布等。
- 3. 关联规则挖掘: 使用Apriori算法或FP-growth算法,根据用户浏览记录计算频繁项集和关联规则。
- 4. 结果评估: 计算关联规则的支持度、置信度和提升度,得出强关联规则。
- 5. 结果分析与应用: 分析得到的关联规则, 为网站提供导航结构优化建议, 以提升用户体验。

网页浏览行为关联规则挖掘

本作业结构:按照作业要求的流程,首先给出带有备注信息的代码,和必要的文字解读,随后给出运行结果。

1. 数据预处理

导入所需的库,包括 pandas 用于数据处理,TransactionEncoder 和 apriori 用于频繁模式挖掘, association rules 用于生成关联规则,以及 matplotlib.pyplot 和 seaborn 用于数据可视化。

- 1 import pandas as pd
- 2 from mlxtend.preprocessing import TransactionEncoder
- 3 from mlxtend.frequent_patterns import apriori, association rules
- 4 import matplotlib.pyplot as plt
- 5 import seaborn as sns

根据每一条数据第一个字段,进行分别处理。

"A"表示attribute,需要将后续内容加入属性字典中;

"C"表示一个case的起始,需要继续读入后续的0个以上的vote;

"V"表示一个case中的一个vote。

```
1
        with open("anonymous-msweb.data", "r", encoding="utf-8") as f:
 2
            lines = f.readlines()
 3
            # attr lines = []
 4
 5
            attr_dic = {}
            case_dic = {}
 6
 7
            vote list = []
            # case vote list = []
 8
9
            valid attr = ["A", "C", "V"]
10
            for i, line in enumerate(lines):
11
                 content = line.strip().split(",")
12
13
14
                attr = content[0]
                 if attr not in valid attr:
15
16
                     continue
17
                 if attr == "A":
18
19
                     attr_dic[content[1]] = [content[3], content[4]]
                 elif attr == "C":
20
                     # c要连着处理之后的v
21
2.2
                     cur_vote = []
23
                     num = i+1
24
                     if num >= len(lines):
25
                         break
26
                     while lines[num].strip().split(",")[0] == 'V':
                         cur_vote.append(lines[num].strip().split(",")[1])
2.7
                         num += 1
28
29
                         if num >= len(lines):
                             break
30
31
                     vote_list.append(cur_vote)
32
                     case_dic[content[1]] = cur_vote[1:]
33
                 else:
34
                     continue
```

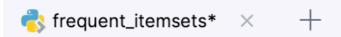
2. 频繁模式挖掘

在开展频繁模式挖掘时,我们首先需要对数据集进行转换。

在使用apriori函数时,输入的数据集需要进行一些预处理才能正确执行。具体来说,需要将数据集转换为适合 apriori函数的形式,可以使用TransactionEncoder类来完成转换。

```
1
       te = TransactionEncoder()
       # 将每一行case_vote数据转换为一个[1, 285]的向量,其中285是在全部数据中出现的不同attr的总
2
   数,如果这行case_vote数据中有某一个attr,那么对应位置为True,否则为False
       # te_ary总共为[32711,285], 32711为case_dic的总数, 也就是所有case的总数
3
       te_ary = te.fit_transform(vote_list)
4
5
       df = pd.DataFrame(te_ary, columns=te.columns_)
6
7
       # 使用apriori函数进行频繁模式挖掘
8
       frequent_itemsets = apriori(df, min_support=0.05, use_colnames=True)
9
10
       print(frequent_itemsets)
```

运行结果如下图所示。



	∨ support	† itemsets
3	0.33126	frozenset({'1008'})
9	0.28685	frozenset({'1034'})
2	0.25872	frozenset({'1004'})
6	0.16294	frozenset({'1018'})
18	0.16080	frozenset({'1008', '1034'})
5	0.15616	frozenset({'1017'})
4	0.14148	frozenset({'1009'})
0	0.13607	frozenset({'1001'})
8	0.09844	frozenset({'1026'})
1	0.09073	frozenset({'1003'})
15	0.07792	frozenset({'1009', '1008'})
17	0.07306	frozenset({'1018', '1008'})
7	0.06490	frozenset({'1025'})
16	0.06123	frozenset({'1017', '1008'})
13	0.06044	frozenset({'1004', '1008'})
12	0.05943	frozenset({'1001', '1018'})
11	0.05521	frozenset({'1003', '1001'})
10	0.05475	frozenset({'1035'})
11	0.05220	frozonco+((11004) 110241)

requent_itemsets

3. 关联规则挖掘, 计算其支持度和置信度

```
# 导出关联规则

rules = association_rules(frequent_itemsets, metric="confidence",
min_threshold=0.1)

print(rules)

# 计算支持度和置信度
rules['support'] = rules['support'].round(4)
rules['confidence'] = rules['confidence'].round(4)
```

运行结果如下图所示。

frozense frozense frozense frozense frozense frozense frozense	et({'10 et({'10 et({'101 et({'10	frozenset({'101 frozenset({'100 frozenset({'100	0.13607 0.13607 0.16294	0.13607 0.09073 0.16294 0.13607	0.05520 0.05520 0.05940	0.60850 0.40580 0.43680	0.60837 0.40567 0.43654	0.04286	1.53011		6193.96933 6193.96933
frozense frozense frozense frozense	et({'10 et({'101 et({'10	frozenset({'101 frozenset({'100 frozenset({'100	0.13607 0.16294	0.16294	0.05940					0.89866	6193.96933
frozense frozense frozense	et({'101 et({'10	frozenset({'100 frozenset({'100	0.16294			0.43680	0.42654				
frozense frozense	et({'10	frozenset({'100		0.13607			0.43034	0.03726	1.48614	0.72567	2827.55741
frozense			0.05070		0.05940	0.36470	0.36455	0.03726	1.35994	0.74896	2827.55741
	et({'10		0.25672	0.33126	0.06040	0.23360	0.23346	-0.02527	0.87257	-0.36060	493.01635
frozense		frozenset({'100	0.33126	0.25872	0.06040	0.18240	0.18233	-0.02527	0.90671	-0.38467	493.01635
	et({'10	frozenset({'103	0.25872	0.28685	0.05330	0.20600	0.20601	-0.02093	0.89813	-0.34633	364.64573
frozense	et({'10	frozenset({'100	0.28685	0.25872	0.05330	0.18580	0.18581	-0.02093	0.91040	-0.35514	364.64573
frozense	et({'10	frozenset({'100	0.14148	0.33126	0.07790	0.55080	0.55060	0.03106	1.48865	0.46423	1170.68463
frozense	et({'10	frozenset({'100	0.33126	0.14148	0.07790	0.23520	0.23516	0.03106	1.12259	0.59598	1170.68463
frozense	et({'101	frozenset({'100	0.15616	0.33126	0.06120	0.39210	0.39192	0.00950	1.10013	0.18394	100.52096
frozense	et({'10	frozenset({'101	0.33126	0.15616	0.06120	0.18480	0.18475	0.00950	1.03520	0.23211	100.52096
frozense	et({'101	frozenset({'100	0.16294	0.33126	0.07310	0.44840	0.44863	0.01909	1.21237	0.31209	395.90282
frozense	et({'10	frozenset({'101	0.33126	0.16294	0.07310	0.22060	0.22067	0.01909	1.07392	0.39065	395.90282
frozense	et({'10	frozenset({'103	0.33126	0.28685	0.16080	0.48540	0.48541	0.06578	1.38589	0.61172	3123.18220
frozense	et({'10	frozenset({'100	0.28685	0.33126	0.16080	0.56060	0.56058	0.06578	1.52189	0.57362	3123.18220

4. 规则评价,结果评估

对规则进行评价,使用Lift、卡方指标;对挖掘结果进行分析。

```
# 评价规则: 使用Lift和卡方指标
 1
 2
        rules['lift'] = rules['support'] / rules['antecedent support']
 3
        rules['chi_square'] = (df.shape[0] * (rules['support'] - rules['antecedent
    support'] * rules['consequent support']) ** 2) \
 4
                             / (rules['antecedent support'] * rules['consequent
    support'] * (1 - rules['antecedent support']) * (1 - rules['consequent support']))
5
        # 对挖掘结果进行分析
 6
 7
        sorted_rules = rules.sort_values(by='lift', ascending=False)
8
        print("关联规则挖掘结果:")
9
10
        print(sorted rules)
11
```

结果如下图所示。



5. 结果可视化

```
# 可视化频繁项集支持度
 1
 2
        plt.figure(figsize=(10, 6))
 3
        plt.subplot(1, 2, 1)
 4
        sns.barplot(x='support', y='itemsets',
    data=frequent itemsets.sort values(by='support', ascending=False),
 5
                    orient='h')
 6
        plt.title('Frequent Itemsets - Support')
 7
        plt.xlabel('Support')
 8
        plt.ylabel('Itemsets')
 9
10
        plt.tight layout()
11
        plt.show()
12
13
        # 可视化关联规则评估指标
14
        plt.figure(figsize=(8, 6))
15
        sns.scatterplot(x='lift', y='confidence', data=rules)
16
        plt.title('Association Rules - Lift vs. Confidence')
17
        plt.xlabel('Lift')
        plt.ylabel('Confidence')
18
19
        plt.show()
```

结果如下图所示:

