

传统文本分析实现网购评论的对象分类

In [1]:

```
1 import re
2 import jieba
3 import pandas as pd
4 import numpy as np
5 from sklearn.model_selection import train_test_split
6 from sklearn.svm import SVC
7 from sklearn.metrics import f1_score, accuracy_score, recall_score
```

数据读取与预处理

In [2]:

```
1 df = pd.read_csv("online_shopping_10_cats.csv")[:60000]
2 df.head()
```

Out[2]:

	cat	label	review
0	书籍	1	做父母一定要有刘墉这样的心态，不断地学习，不断地进步，不断地给自己补充新鲜血液，让自己保持...
1	书籍	1	作者真有英国人严谨的风格，提出观点、进行论述论证，尽管本人对物理学了解不深，但是仍然能感受到...
2	书籍	1	作者长篇大论借用详细报告数据处理工作和计算结果支持其新观点。为什么荷兰曾经县有欧洲最高的生产...
3	书籍	1	作者在战几时之前用了 " 拥抱 " 令人叫绝。日本如果没有战败，就会有会有美军的占领，没胡官僚主义的延...
4	书籍	1	作者在少年时即喜阅读，能看出他精读了无数经典，因而他有一个庞大的内心世界。他的作品最难能可贵...

In [3]:

```
1 #使用re正则提取中文并用jieba分词提取词语语料
2 extract_chinese = re.compile(r'[\u4e00-\u9fa5]+')
3 chinese_corpus_raw = df['review'].tolist()
4 chinese_corpus_raw
5 df['chinese_corpus']=[jieba.lcut("".join(extract_chinese.findall(str(corpus)))) for corpus in
6 df.head()
```

Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\23176\AppData\Local\Temp\jieba.cache
Loading model cost 0.606 seconds.
Prefix dict has been built successfully.

Out[3]:

cat	label	review	chinese_corpus
0	书籍	1 做父母一定要有刘墉这样的心态，不断地学习，不断地进步，不断地给自己补充新鲜血液，让自己保持...	[做, 父母, 一定, 要, 有, 刘墉, 这样, 的, 心态, 不断, 地, 学习, 不断...
1	书籍	1 作者真有英国人严谨的风格，提出观点、进行论述论证，尽管本人对物理学了解不深，但是仍然能感受到...	[作者, 真有, 英国人, 严谨, 的, 风格, 提出, 观点, 进行, 论述, 论证, 尽...
2	书籍	1 作者长篇大论借用详细报告数据处理工作和计算结果支持其新观点。为什么荷兰曾经县有欧洲最高的生产...	[作者, 长篇大论, 借用, 详细, 报告, 数 据处理, 工作, 和, 计算结果, 支持, ...
3	书籍	1 作者在战几时之前用了 " 拥抱 " 令人叫绝。日本如果没有战败，就会有美军的占领，没胡官僚主义的延...	[作者, 在, 战, 几时, 之前, 用, 了, 拥 抱, 令人, 叫绝, 日本, 如果, 没...
4	书籍	1 作者在少年时即喜阅读，能看出他精读了无数经典，因而他有一个庞大的内心世界。他的作品最难能可贵...	[作者, 在, 少年, 时即, 喜, 阅读, 能, 看出, 他, 精读, 了, 无数, 经典...

In [9]:

```
1 #将每条评论分词后整合到一个列表中，将每个词用空格隔开放入一个列表中
2 words_list = []
3 corpus = []
4 for corpu in df['chinese_corpus'].tolist():
5     words_list.append(corpu)
6     corpus.append(' '.join(corpu))
7 words_list[0]
```

Out[9]:

```
['做',
 '父母',
 '一定',
 '要',
 '有',
 '刘墉',
 '这样',
 '的',
 '心态',
 '不断',
 '地',
 '学习',
 '不断',
 '地',
 '进步',
 '不断',
 '地',
 '给']
```

In [5]:

```
1 #每个词用空格分割
2 corpus
```

Out[5]:

['做父母一定要有刘墉这样的心态不断地学习不断地进步不断地给自己补充新鲜血液让自己保持一颗年轻的心我想这是他能很好的和孩子沟通的一个重要因素读刘墉的文章总能让我看到一个快乐的平易近人的父亲他始终站在和孩子同样的高度给孩子创造着一个充满爱和自由的生活环境很喜欢刘墉在字里行间流露出的做父母的那种小狡黠让人总是忍俊不禁父母和子女之间有时候也是一种战斗武力争斗过于低级了智力较量才更有趣味所以做父母的得加把劲了老思想老观念注定会一败涂地生命不息学习不止家庭教育真的是乐在其中',
'作者真有英国人严谨的风格提出观点进行论述论证尽管本人对物理学了解不深但是仍然能感受到真理的火花整本书的结构颇有特点从当时本书写于八十年代流行的计算机话题引入再用数学物理学宇宙学做必要的铺垫这些内容占据了大部分篇幅最后回到关键问题电脑能不能代替人脑和现在流行的观点相反作者认为人的某种洞察是不能被算法模拟的也许作者想说人的灵魂是无可取代的',
'作者长篇大论借用详细报告数据处理工作和计算结果支持其新观点为什么荷兰曾经具有欧洲最高的生产率为什么在文化上有着深刻纽带关系的中国和日本却在经济发展上有着极大的差异为什么英国的北美殖民地造就了经济强大的美国而西班牙的北美殖民却造就了落后的墨西哥很有价值但不包

In [6]:

```
1 #构建类别与编号的转换字典，并将类别转成编号
2 class2idx = {'书籍':0, '平板':1, '手机':2, '水果':3, '洗发水':4, '热水器':5, '蒙牛':6, '衣服':7,
3 idx2class = {idx:class_ for class_, idx in class2idx.items()}
4 class_idx = [class2idx[class_] for class_ in df['cat'].values]
5 class2idx
```

Out[6]:

```
{'书籍': 0,
'平板': 1,
'手机': 2,
'水果': 3,
'洗发水': 4,
'热水器': 5,
'蒙牛': 6,
'衣服': 7,
'计算机': 8,
'酒店': 9}
```

TF-IDF和Bag-of-words实现向量化

In [7]:

```
1 #sklearn实现TF-IDF
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 tfidf_vec=TfidfVectorizer(max_features=1000)
4 tfidf_matrix=tfidf_vec.fit_transform(corpus)
```

In [8]:

```
1 #sklearn实现bagofword
2 from sklearn.feature_extraction.text import CountVectorizer
3 vectorizer = CountVectorizer(max_features=1000)
4 tfidf_matrix1 = vectorizer.fit_transform(corpus)
```

In [9]:

```
1 #划分训练集验证集
2 train_x1, valid_x1, train_y1, valid_y1 = train_test_split(tfidf_matrix, class_idx, random_state=1)
```

In [10]:

```
1 #划分训练集验证集
2 train_x2, valid_x2, train_y2, valid_y2 = train_test_split(tfidf_matrix1, class_idx, random_state=1)
```

使用SVM实现网购评论的对象分类

In [11]:

```
1 #定义SVM分类器并进行分类, 评价指标设为准确率, 召回率和F1值
2 def SVM(train_x, train_y, valid_x, valid_y):
3     svm = SVC(max_iter=500)
4     svm.fit(train_x, train_y)
5     val_pred = svm.predict(valid_x)
6     f1 = f1_score(valid_y, val_pred, average='macro')
7     Accuracy_score = accuracy_score(valid_y, val_pred)
8     Recall_score = recall_score(valid_y, val_pred, average='macro')
9     print(f'Accuracy_score: {Accuracy_score}')
10    print(f'Recall_score: {Recall_score}')
11    print(f'f1_score: {f1}')
```

In [12]:

```
1 #使用SVM分类器并进行分类(TF-IDF)
2 SVM(train_x1, train_y1, valid_x1, valid_y1)
```

D:\Anaconda3\envs\bd\lib\site-packages\sklearn\svm_base.py:289: ConvergenceWarning:
Solver terminated early (max_iter=500). Consider pre-processing your data with StandardScaler or MinMaxScaler.

ConvergenceWarning,

Accuracy_score:0.7315833333333334

Recall_score:0.7237494564062985

f1_score:0.7507279936063987

In [13]:

```
1 #使用SVM分类器并进行分类(Bag-of-words)
2 SVM(train_x2, train_y2, valid_x2, valid_y2)
```

D:\Anaconda3\envs\bd\lib\site-packages\sklearn\svm_base.py:289: ConvergenceWarning:
Solver terminated early (max_iter=500). Consider pre-processing your data with StandardScaler or MinMaxScaler.

ConvergenceWarning,

Accuracy_score:0.6163333333333333

Recall_score:0.6450487458510076

f1_score:0.6550754410041175