



南京工业大学
NANJING TECH
UNIVERSITY

用户数据采集与关联分析

(结课作业)

吴志祥

18205185639

1030624832@qq.com





南京工业大学
NANJING TECH
UNIVERSITY

01 文本数据分析

吴志祥

18205185639

cnwzx2012@njtech.edu.cn



```
from snowlp import SnowNLP
import matplotlib.pyplot as plt
import numpy as np

# ----- 1. 准备待分析文本 -----
# 示例：新能源汽车用户评论文本（包含积极、消极、中性评价）
text = """
这款新能源汽车的续航真的超出预期，充一次电500公里完全没问题，太满意了！
内饰做工精致，智能驾驶功能也很好用，操作简单流畅，推荐大家购买。
价格有点贵，比同级别燃油车贵了3万多，性价比一般。
刹车有点软，试驾时感觉制动距离偏长，安全性让人担心。
空间表现中规中矩，后排坐三个成年人刚好，不算拥挤也不算宽敞。
售后服务态度很好，有问必答很及时，维修保养费用也合理。
续航里程严重，实际只能跑350公里左右，和宣传的差距太大了。
外观设计很时尚，回头率很高，身边朋友都问是什么车型。
充电速度有点慢，快充需要1.5小时才能充满，不太方便长途出行。
整体来说还不错，虽然有小缺点，但在新能源车里算是表现优秀的了。
"""

# ----- 2. 文本分割与情感计算 -----
# 按句子分割文本（基于中文标点）
sentences = [s.strip() for s in text.replace("\n", "").split(".") if s.strip()]

# 逐句计算情感得分（0-1），并分类
emotion_scores = []
positive_sentences = [] # 积极（得分≥0.6）
neutral_sentences = [] # 中性（0.4 < 得分 < 0.6）
negative_sentences = [] # 消极（得分≤0.4）

for sent in sentences:
    score = SnowNLP(sent).sentiments
    emotion_scores.append(score)
    if score >= 0.6:
        positive_sentences.append(sent)
    elif score <= 0.4:
        negative_sentences.append(sent)
    else:
        neutral_sentences.append(sent)

# 统计各类情感占比
total = len(sentences)
pos_ratio = len(positive_sentences) / total * 100
neu_ratio = len(neutral_sentences) / total * 100
neg_ratio = len(negative_sentences) / total * 100

# ----- 3. 生成可视化图表 -----
plt.rcParams['font.sans-serif'] = ['SimHei'] # 解决中文显示问题
plt.rcParams['axes.unicode_minus'] = False # 解决负号显示问题

# 创建2x1子图（情感分布直方图 + 情感占比饼图）
fig, (ax1, ax2) = plt.subplots(nrows=2, ncols=1, figsize=(12, 10))

# 子图1：情感得分分布直方图
ax1.hist(emotion_scores, bins=10, color='#1f77b4', alpha=0.7, edgecolor='black')
ax1.axvline(x=0.4, color='orange', linestyle='--', linewidth=2, label='消极-中性分界 (0.4)')
ax1.axvline(x=0.6, color='red', linestyle='--', linewidth=2, label='中性-积极分界 (0.6)')
ax1.set_xlabel('情感得分 (0=消极, 1=积极)', fontsize=12)
ax1.set_ylabel('句子数量', fontsize=12)
ax1.set_title('新能源汽车用户评论情感得分分布', fontsize=14, fontweight='bold')
ax1.legend()
ax1.grid(alpha=0.3)

# 子图2：情感占比饼图
labels = ['积极评价', '中性评价', '消极评价']
sizes = [pos_ratio, neu_ratio, neg_ratio]
colors = ['#2ca02c', '#ff7f0e', '#d62728']
explode = (0.05, 0, 0) # 突出积极评价

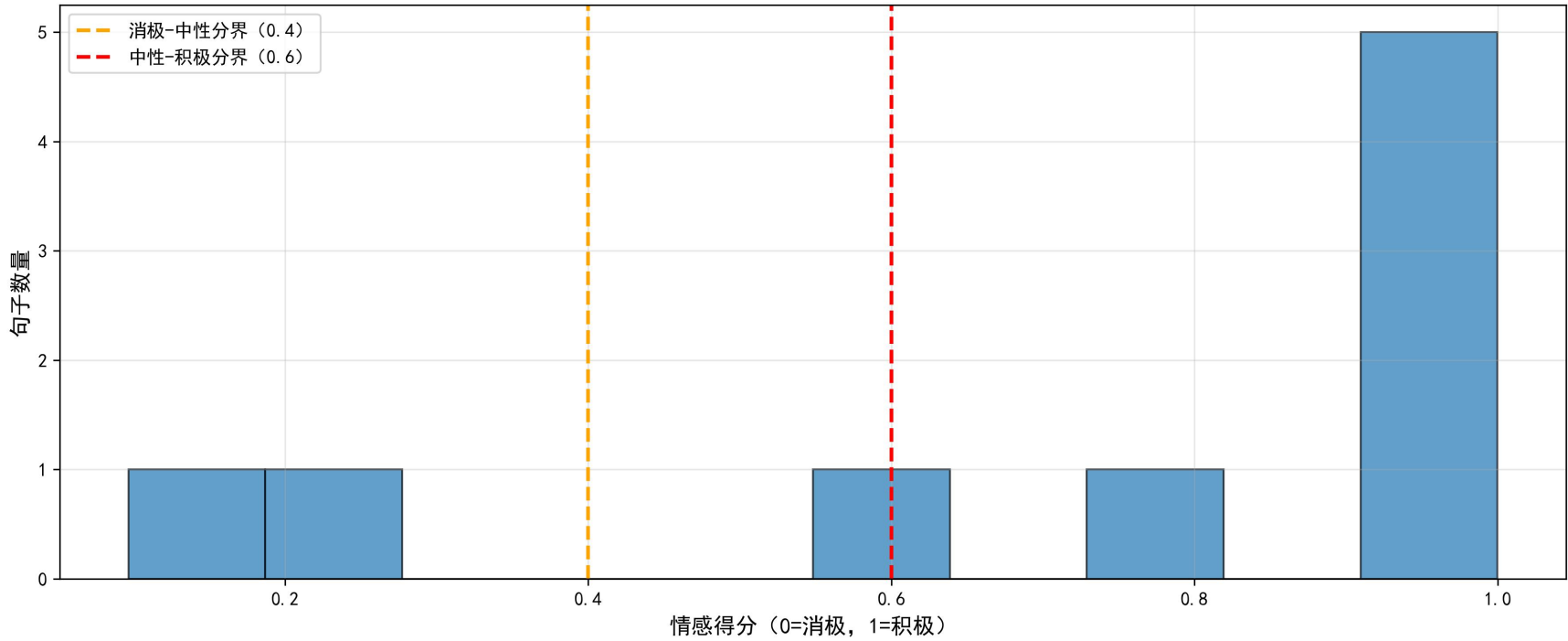
wedges, texts, autotexts = ax2.pie(
    sizes, explode=explode, labels=labels, colors=colors, autopct='%1.1f%%',
    shadow=True, startangle=90, textprops={'fontsize': 11}
)
ax2.set_title('新能源汽车用户评论情感占比', fontsize=14, fontweight='bold')

# 调整子图间距
plt.tight_layout()

# 保存图片（高清图式）
plt.savefig('emomotion_analysis_result.png', dpi=300, bbox_inches='tight')
plt.show()

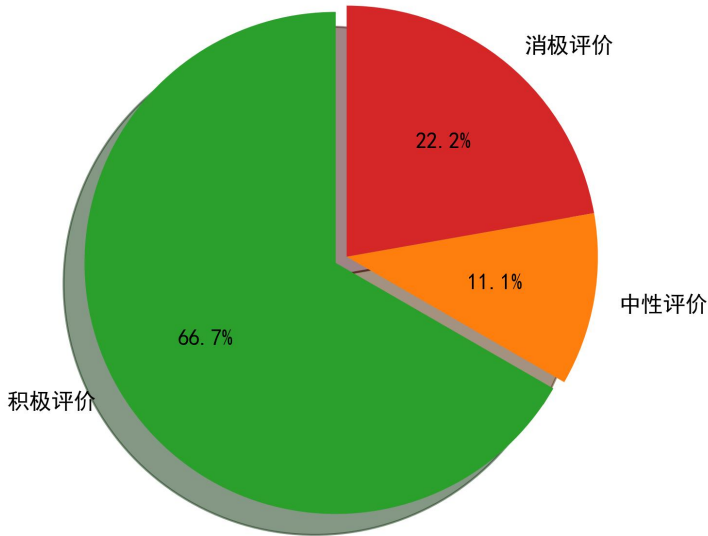
# ----- 4. 输出分析结果 -----
print("=== 情感分析结果汇总 ===")
print(f"总评论句子数: {total} 句")
print(f"积极评价: {len(positive_sentences)} 句 ({pos_ratio:.1f}%)")
print(f"中性评价: {len(neutral_sentences)} 句 ({neu_ratio:.1f}%)")
print(f"消极评价: {len(negative_sentences)} 句 ({neg_ratio:.1f}%)")
print("\n==== 典型句子示例 ====")
print("积极评价示例:", positive_sentences[0] if positive_sentences else "无")
print("消极评价示例:", negative_sentences[0] if negative_sentences else "无")
```

新能源汽车用户评论情感得分分布



情感得分 (0=消极, 1=积极)

新能源汽车用户评论情感占比



第四讲 情感分析

2. 完成sentiment_analysis_4份代码。做截图，并简要做

alysis_4, 1。

```
In [x]: !pip install snowlp
!pip install -U textblob
!python -m textblob.download_corpora

英文

In [x]: # 英文调用TextBlob NLP处理包

Features
• Noun phrase extraction
• Part-of-speech tagging
• Sentiment analysis
• Classification (Naive Bayes, Decision Tree)
• Language translation and detection powered by Google Translate
• Tokenization (splitting text into words and sentences)
• Word and phrase frequencies
• Parsing
• n-grams
• Word inflection (pluralization and singularization) and lemmatization
• Spelling correction
• Add new models or languages through extensions
• WordNet integration

In [x]: text = "I am happy today. I feel sad today."

In [x]: from textblob import TextBlob
blob = TextBlob(text)

In [x]: blob

In [x]: # 很好不动的打印出来了?
# 实际上已经把文本分成了句子了, 看一看
blob.sentences

In [x]: blob.sentences[0].sentiment

In [x]: # 上面的结果什么意思呢?
# 情感值在 0.5 左右的点, 说明一下, 情感极性的变化范围是[-1, 1], -1代表完全负面, 1代表完全正面。
# 我表达的是我很高兴, 那么这个结果是对的

In [x]: blob.sentences[1].sentiment

In [x]: # 整段文本的情感呢?
blob.sentiment

In [x]: # 你可能会觉得没有道理, 怎么一句“高兴”, 一句“沮丧”, 合起来最后会得到正向结果呢?
# 首先不同极性的词, 在数量上是有区别的, 我们应该可以找到一个“沮丧”更为负面的词汇, 而且这也符合逻辑, 谁会这么“天上一脚, 地下一脚”手痒地根据自
« ..... »

中文

In [x]: # 中文分析, 用的是SnowNLP包

Features
• 中文分词 (Character-Based Generative Model)
• 词性标注 (Trie 3-gram 隐马尔可夫模型)
• 情感分析 (现在训练数据主要是采集东西时的评价, 所以对其他的一些可能效果不是很好, 待解决)
• 文本分类 (Naive Bayes)
• 转换成拼音 (Trie树实现的最大匹配)
• 繁体转简体 (Trie树实现的最大匹配)
• 提取文本关键词 (TextRank算法)
• 提取文本摘要 (TextRank算法)
• 词, idf
• Tokenization (分割成句子)
• 文本相似 (BM25)
• 支持python3 (感谢erning)

In [x]: text_cn = u"我今天很快乐, 我今天很愤怒。"

In [x]: # 注意在引号前面我们加了一个字母u, 它很重要, 因为它提示python, “这一段我们输入的文本编码格式是Unicode, 别搞错了哦”。至于文本编码格式的章节, 看
« ..... »

In [x]: from snownlp import SnowNLP

In [x]: senti_cn = SnowNLP(text_cn)

In [x]: # 看看snownlp包的句子能力
for sentence in senti_cn.sentences:
    print(sentence)

In [x]: senti_cn_1 = SnowNLP(senti_cn.sentences[0])

In [x]: # 一个细节上的问题, 英文是s.sentiment, 中文是s.sentiments, 弄了一个p
# 另外, 在句法上和英文的也有不同, 比如直接用语句/ senti_cn.sentences[0].sentiments是会报错的
senti_cn_1.sentiments

In [x]: senti_cn_2 = SnowNLP(senti_cn.sentences[1])

In [x]: senti_cn_2.sentiments

这里你肯定发现了问题——“愤怒”这个词表达了如此强烈的负面情绪, 为何得分依然还是正的?

这是因为SnowNLP和textblob的计分方法不同。SnowNLP的情感分析取值, 表达的是“这句话代表正面情感的概率”。也就是说, 对“我今天很愤怒”一句,
SnowNLP认为, 它表达正面情感的概率很低很低。

这样解释就是ok了

In [x]: senti_cn.sentiments

In [x]: # 整个句子, 那就就有问题了

作业 1
注: 设计课程的时候, 可以把作业单独拿出来, 放在另外一个文档中
```


第四讲 情感分析

2. 完成sentiment_analysis
4份代码。做截图，并简要做

做截图，并简要做

情感分析-高级-时间序列

工作步骤

- 1. 数据清洗和预处理 (Data Cleaning and Preprocessing)
- 2. 文本特征提取 (Text Feature Extraction)
- 3. 模型训练和评估 (Model Training and Evaluation)
- 4. 模型部署和监控 (Model Deployment and Monitoring)

准备工作，之前已经安装好了 Jupyter Notebook，这里再次确认一下！

```
!pip install jupyter
```

```
!pip install pandas numpy matplotlib seaborn nltk
```

```
!pip install word2vec
```

```
!pip install gensim
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

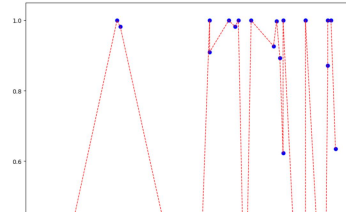
```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```

```
!pip install nltk
```



第四讲 情感分析

2. 完成sentiment_analysis_4份代码。做截图，并简要做

analysis_4, 斤。

在deepseek api平台申请API-key；然后交一点钱，很便宜的。

In []: # https://platform.deepseek.com/transactions

deepseek

选择模型

API keys

充值

账单

账单

充值账单

提现账单

提现账单

提现账单

提现账单

提现账单

提现账单

提现账单

提现账单

提现账单

提现账单

提现账单

In [1]:

```
import requests
import json

# Deepseek API 端点
url = "https://api.deepseek.com/v1/chat/completions"

# 替换为您的 Deepseek API 密钥
API_KEY = "sk-4ef9000e1e6437484c3ff2df4d3a10" # 直接复制过来

# 请求头，包含 API 密钥和请求类型
headers = {
    "Authorization": f"Bearer {API_KEY}",
    "Content-Type": "application/json"
}

# 患者描述文本
text = """
我今年35岁，退休后感到了生活失去了重心，开始出现失眠、头痛和无力乏力的症状。
"但是，这感觉非常奇怪，在衣服的边缘像针扎一样疼痛。"
我常常感到心慌、胸闷，胃部沉重得像压了一块石头。
"对尖锐声音变得特别敏感，电话铃声都会让我紧张。"
"多次到医院检查，结果都显示没有器质性病变。"
"我觉得不能出门，不想与人交流，整天把自己关在屋里，郁郁寡欢，感觉生活毫无意义。"
"""

# 构造提示词，要求模型提取物理和情感实体
prompt = """
请从以下患者描述中提取出具体的身体部位、症状以及对应的情感状态。
并以 JSON 格式返回，格式如下：
{"实体": [{"部位": "...", "症状": "...", "情感": "..."}]}
"""

# 请求体，包含模型参数和提示词
data = {
    "model": "deepseek-chat",
    "messages": [
        {"role": "user", "content": prompt}
    ],
    "temperature": 0.5
}

try:
    # 发送 POST 请求
    response = requests.post(url, headers=headers, data=json.dumps(data))

    # 检查响应状态码
    if response.status_code == 200:
        # 解析 JSON 响应
        result = response.json()
        # 提取模型生成的内容
        generated_text = result["choices"][0]["message"]["content"]
        print("提取情感实体抽取结果:")
        print(generated_text)
    else:
        # 处理错误响应
        print(f"请求失败，状态码: {response.status_code}")
        print(f"错误信息: {response.text}")

except requests.exceptions.RequestException as e:
    # 处理网络请求异常
    print(f"网络请求失败: {e}")
except json.JSONDecodeError as e:
    # 处理 JSON 解析异常
    print(f"JSON 解析失败: {e}")
except Exception as e:
    # 处理其他异常
    print(f"发生未知错误: {e}")
```

In [2]:

```
"""
提取情感实体抽取结果
"""
{
  "实体": [
    {
      "部位": "全身",
      "症状": "疲乏无力",
      "情感": "生活失去重心"
    },
    {
      "部位": "皮肤",
      "症状": "异常敏感，衣服边缘像针扎一样疼痛",
      "情感": "无"
    },
    {
      "部位": "心脏",
      "症状": "心慌",
      "情感": "无"
    },
    {
      "部位": "胸部",
      "症状": "胸闷",
      "情感": "无"
    }
  ]
}
```

In [3]:

```
"""
提取情感实体抽取结果
"""
{
  "实体": [
    {
      "部位": "全身",
      "症状": "疲乏无力",
      "情感": "生活失去重心"
    },
    {
      "部位": "皮肤",
      "症状": "异常敏感，衣服边缘像针扎一样疼痛",
      "情感": "无"
    },
    {
      "部位": "心脏",
      "症状": "心慌",
      "情感": "无"
    },
    {
      "部位": "胸部",
      "症状": "胸闷",
      "情感": "无"
    }
  ]
}
```

经济与管理学院

1/4/2026

第四讲 情感分析

2. 完成sentiment_analysis_4份代码。做截图，并简要做结分析。

