



期末汇报

张亦佳

处理文本特征

- 合并所有文本特征

```
demo["all_text"] = demo["房屋优势"] + " " + demo["核心卖点"] + " "
```

- 使用jieba进行分词并绘制词云图

```
tokens_list = []
for text in demo["all_text"]:
    if isinstance(text, str):
        words = jieba.lcut(text)
        tokens = [word for word in words if len(word) > 1] # 只保留长
        tokens_list.append(tokens)
    else:
        tokens_list.append([]) # 处理空文本
```

```
1 wc=WordCloud(font_path=font_path,
2               background_color='#FFEEEE', # 背景色 (浅红色)
3               colormap='hsv' # 词语颜色方案 (彩虹色)
4               )
5 wc.fit_words(dict(all_tokens))
6
7 fig=plt.figure(figsize=(10,5))
8 plt.imshow(wc)
9 plt.axis('off')
10 plt.show()
```

- 应用Word2Vec训练出词向量

```
1 model = Word2Vec(
2     sentences=tokens_list,
3     vector_size=20, #
4     window=3, # 前后5个词作为上下文
5     min_count=3, # 整体中出现次数少于5的词删掉
6     workers=4, # 并行计算加速训练过程
7     epochs=5 # 模型看5遍所有数据
8 )
9 # 可以参数优化吗?
10 # 语义相近的词在向量空间中相近
```

处理文本特征

- 用训练好的Word2Vec得到每条数据的文本的词向量

```
2 doc_vectors = []
3
4 # 遍历每个房屋的分词结果
5 for tokens in tokens_list:
6     # 初始化一个50维的零向量
7     vector = np.zeros(50)
8     count = 0 # 记录文档中有多少在Word2Vec模型中用到的有效词
9
10    for word in tokens:
11        # 检查词是否在Word2Vec模型的词汇表中
12        if word in model.wv:
13            # 累加该词的向量到文档向量
14            vector += model.wv[word]
15
16            count += 1 # 统计文档中有多少个词被包含在词向量模型中
17
18    # 计算平均向量 保证不同长度的文本有可比性
19    if count > 0:
20        vector /= count # 向量除以词数
21
22    # 将当前文档向量添加到结果列表
23    doc_vectors.append(vector)
24
```

- 随机森林模型预测