

第一讲

1-2

THULAC: 一个高效的中文词

Q

⌂ ... □ ×

THULAC：一个高效的中文词法分析工具包

欢迎使用THULAC中文分词工具包demo系统

周末去公园散步时，刚好赶上樱花开得最盛的时节。粉白的花瓣被风一吹，簌簌落在肩头，空气里都是淡淡的甜香。不远处有小朋友追着泡泡跑，笑声脆生生的，连坐在长椅上晒太阳的老爷爷，都忍不住跟着笑眯了眼。原来最治愈的时光，从来都藏在这些寻常的小美好里呀。

周末去公园散步时，刚好赶上樱花开得最盛的时节。粉白的花瓣被风一吹，簌簌落在肩头，空气里都是淡淡的甜香。不远处有小朋友追着泡泡跑，笑声脆生生的，连坐在长椅上晒太阳的老爷爷，都忍不住跟着笑眯了眼。原来最治愈的时光，从来都藏在这些寻常的小美好里呀。

词性解释

n/名词 np/人名 ns/地名 ni/机构名 nz/其它专名
m/数词 q/量词 mq/数量词 t/时间词 ft/方位词 s/处所词
v/动词 vm/能愿动词 vd/趋向动词 a/形容词 d/副词
h/前接成分 k/后接成分 i/习语 j/简称
z/代词 c/连词 p/介词 u/助词 y/语气助词
e/叹词 o/拟声词 g/语素 w/标点 x/其它

版权所有：清华大学自然语言处理与社会人文计算实验室
Copyright: Natural Language Processing and Computational Social Science Lab, Tsinghua University

```
Python 3.13.5 | packaged by Anaconda, Inc. | (main, Jun 12 2025, 16:37:03) [MSC v.1929 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> print("hello world 谢雨")
hello world 谢雨
>>> |
```

001

```
seg_list = jieba.cut("南京大学生都爱南京市长江大桥")
```

```
print('*'.join(seg_list))
```

南京*大学生*都*爱*南京市*长江大桥

```
jieba.load_userdict('dict.txt')
```

```
seg_list = jieba.cut("吴志祥是南京工业最好大学青年教师，经济与管理学院的，他对那种二次元小魔仙是无感的，这怎么行？python看上去还有点人性")
```

```
print('%'.join(seg_list))
```

吴志祥%是%南京工业最好大学%青年教师%，%经济%与%管理%学院%的%，%他%对%那种%二次元小魔仙%是%无感%的%，%这%怎么%行%？%python%看上去%还%有点人性

```
for seg in seg_list:
    if seg not in stopwords:
        final += seg+'*'
```

```
print (final) # 做实体抽取的时候，停用词表很管用
```

使用*词表*之后*效果*就*好看*很多*什么*之类*词*就*不见*

质量,不大,好

```
ss = jieba.cut('质量不大好')
```

```
print(", ".join(ss))
```

质量,不大好

生成 代码

```
s1 = SnowNLP(u"吴志祥是南京工业大学青年教师，他对那种二次元小魔仙是无感的，这怎么行？")
```

```
print(", ".join(s1.words)) # 因为snownlp擅长处理英文
```

吴志祥是南京工业大学青年教师，他对那种二次元小魔仙是无感的，这怎么行？

002

```
[ '的', '了', '是', '啊', '、', '，', '。', '，', '停用' ]

seg_list_huang = jieba.cut('黄旭华,1926年3月12日出生于广东省汕尾市,原籍广东省揭阳市。1949年毕业于上海交通大学。历任北京海军核潜艇研究室副总工程师、中船重工集团公司核潜艇总体研究设计所研究员

final = ''

for seg in seg_list_huang:
    if seg not in stopwords:
        final+= seg+'/'

print(final)
```

黄旭华/1926/年/3/月/12/日出/生于/广东省/汕尾市/原籍/广东省/揭阳市/1949/年/毕业/于/上海交通大学/历任/北京/海军/核潜艇/研究室/副/总工程师/中船重工集团公司/核潜艇/总体/研究/设计所/研究员/名誉/所长/

003

```
▷ print("特定词汇词频统计结果：")
  for word in target_words:
    print(f"'{word}': {word_counts[word]}次")
Python

... 特定词汇词频统计结果：
'数字化': 2次
'智能化': 3次
'安全': 2次

# 输出所有词汇的词频（按频率降序）
print("\n所有词汇词频统计（前20个）：")
for word, count in word_counts.most_common(20):
    print(f"'{word}': {count}次")
Python

... 所有词汇词频统计（前20个）：
',': 13次
'': 9次
'管理': 5次
'、': 5次
'与': 4次
'“': 3次
'企业': 3次
'”': 3次
'体系': 3次
'智能化': 3次
'经营': 3次
'打造': 3次
'能力': 3次
'供应链': 3次
'建设': 2次
```

004

```
else:
    print(f"请求失败, 状态码: {response.status_code}")
    print(response.text)
```

Python

... 提取到的实体和专业术语:

```
```json
{
 "理论": [
 "肿瘤免疫微环境",
 "T细胞耗竭",
 "免疫编辑理论",
 "免疫抑制"
],
 "方法": [
 "单细胞RNA测序",
 "scRNA-seq",
 "细胞亚群聚类",
 "轨迹分析",
 "pseudotime推断",
 "细胞间通讯网络"
],
 "工具": [
 "Seurat",
 "Monocle3",
 "CellChat"
],
 "疾病与样本": [
 "非小细胞肺癌",
 "肿瘤样本"
]
}
... "个体化免疫治疗"
```

# 4

该研究围绕基于关键词的学术文本聚类展开，核心是探究聚类集成方法对学术文本分类性能的影响，以及关键词抽取方法、关键词个数这两个因素的作用。

聚类集成方法可显著提升学术文本聚类性能。研究对比了以 K-means 为基础的聚类集成、以增量聚类为基础的聚类集成，与单一的 K-means、增量聚类算法，发现聚类集成方法的 F1 值普遍更高，且经 T 检验验证，性能差异具有显著性；同时，在不同关键词抽取方法和关键词个数下，聚类集成方法的稳定性也更好。

关键词个数与聚类集成性能呈正相关。随着关键词个数从 5 个增加到 60 个，学术文本聚类的性能逐步提升，能更充分地体现文本主旨；且在关键词个数较少时，聚类集成方法的优势更突出，可缓解关键词数量不足对分类的影响。

# 第二讲

## 1

关键词总数大幅增长，研究维度持续丰富

环保相关主题文献的关键词总数十年间增长显著，2024 年较 2014 年平均增长超 2 倍

2014 年前后，“污染治理”“PM2.5”“污水处置”等传统治污类关键词占比达 40%-50%，是绝对核心。

早期关键词多集中于单一治理场景，关联性较弱；2024 年关键词形成多维度聚类网络，“协同发展”“生态产品”“绿色经济”等跨领域关键词突现强度最高，成为连接不同研究方向的核心纽带。

关键词关联度提升，反映环保研究从“单点突破”转向“系统协同”的整体趋势。



# 2

tf1\_full\_text\_sanguo.ipynb • tf1\_full\_text\_sanguo-checkpoint.ipynb • 001-word\_cut\_基本分词.ipynb • sentiment\_analysis\_2\_timeline-checkpoint.ipynb • tf1\_full\_text\_huangxuhua.ipynb • sentiment\_an ...

D:\> 88888 > 用户数据 > 用户数据 > 第2讲 感知世界: 词频统计与分析 > tf1\_full\_text\_sanguo.ipynb > M4 本程序是基于jieba分词的全文本词频统计

生成 + 代码 + Markdown 全部运行 ...

articleSet

Python

...

{ '聲大起',  
'正行',  
'可達',  
'角',  
'孫堅同',  
'可請汝兄',  
'堅挺',  
'飲食漸',  
'小',  
'一過',  
'知操',  
'名堅',  
'多疑',  
'求救',  
'卓急',  
'神色',  
'慢',  
'紀律',  
'是',  
'嘗作吏',  
'三軍掩',  
'許吾兒',  
'若引兵',  
'中點',  
'提劍出',  
...  
'路經',  
'顧惜',  
'肅見布',  
'甲',  
...}

tf1\_full\_text\_sanguo.ipynb • tf2\_the\_three\_kingdoms.ipynb X • tf1\_full\_text\_sanguo-checkpoint.ipynb • 001-word\_cut\_基本分词.ipynb • sentiment\_analysis\_2\_timeline-checkpoint.ipynb • tf1\_full\_text\_h ...

D:\> 88888 > 用户数据 > 用户数据 > 第2讲 感知世界: 词频统计与分析 > tf2\_the\_three\_kingdoms.ipynb > M4 你没有发现下面的代码很神奇吗? data\_txt.count一下, 就可以统计词! 实际上是字符串匹配的过程! > name\_dict

生成 + 代码 + Markdown 全部运行 ...

for name in names:  
 name\_dict[name]=data\_txt.count(name)

Python

[20]

name\_dict

Python

[21]

{ '諸葛亮': 149,  
'關羽': 8,  
'劉備': 291,  
'曹操': 907,  
'孫權': 315,  
'張飛': 347,  
'呂布': 332,  
'周瑜': 235,  
'趙雲': 301,  
'龐統': 80,  
'司馬懿': 272,  
'黃忠': 179,  
'馬超': 212}

# 定义 画图 函数

Python

[21]

def make\_chinese\_plot\_ready():  
 from matplotlib import rcParams  
 rcParams['font.family'] = 'Heiti TC' # mac笔记本电脑直接替换字体  
 #rcParams['font.sans-serif'] = ['FangSong'] # 或者直接使用电脑有的字体 FangSong

Python

[23]

# 3

tf1\_full\_text\_sanguo.ipynb • tf2\_the\_three\_kingdoms.ipynb • **tf1\_full\_text\_huangxuhua.ipynb** D:\88888\... X • tf1\_full\_text\_sanguo-checkpoint.ipynb • 001-word\_cut\_基本分词.ipynb • sentiment\_analysis\_...

D: > 88888 > 用户数据 > 用户数据 > 第2讲 感知世界: 词频统计与分析 > tf1\_full\_text\_huangxuhua.ipynb > M4 本程序是基于jieba分词的全文本词频统计 > M4 全文本词频统计的步骤 > M4 集合 (数据结构的一种) > # 输出词频的前N个

生成 + 代码 + Markdown | 全部运行 ... 选择内核

```
print(articlelist[i])
```

[17] Python

```
... ('黄旭华', 53)
('核潜艇', 32)
('采集', 29)
('学术', 22)
('资料', 21)
('工作', 17)
('成长', 15)
('小组', 14)
('院士', 13)
('进行', 13)
('专业', 13)
('技术', 12)
('研制', 12)
('我国', 12)
('工程', 11)
('访谈', 10)
('第一代', 8)
('介绍', 8)
('主要', 8)
('科学', 8)
('思想', 7)
('人生', 7)
('及其', 7)
('历史', 7)
('传记', 7)
...
('节点', 3)
('叙述', 3)
('计划', 3)
('直接', 3)
```

Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...

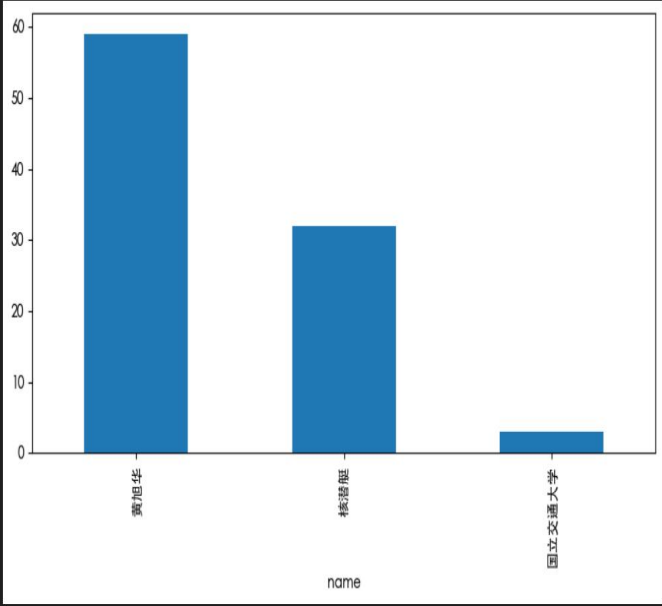
tf1\_full\_text\_sanguo.ipynb • tf2\_the\_three\_kingdoms.ipynb • tf1\_full\_text\_huangxuhua.ipynb D:\88888\... • **tf2\_the\_huangxuhua.ipynb** X • tf1\_full\_text\_sanguo-checkpoint.ipynb • 001-word\_cut\_基本分词...

D: > 88888 > 用户数据 > 用户数据 > 第2讲 感知世界: 词频统计与分析 > tf2\_the\_huangxuhua.ipynb > M4 你没有发现下面的代码很神奇吗? data\_txt.count一下, 就可以统计词! 实际上是字符串匹配的过程! > draw\_dict(terms\_dict)

生成 + 代码 + Markdown | 全部运行 ... 选择内核

```
draw_dict(terms_dict)
```

[18] Python



name	frequency
黄旭华	58
核潜艇	32
国立交通大学	3

# 4

这篇论文用谷歌图书、谷歌学术的数字化资源研究杰出物理学家的科学声誉，发现牛顿、爱因斯坦等虽离世，影响力却持续数世纪。全球范围看，20 世纪中期后爱因斯坦在学术界声誉超牛顿，但存在“群体内偏好”，比如英式英语书籍中牛顿更受关注，爱因斯坦则在美式英语、德语书籍中提及度更高。共现分析显示，牛顿的声誉与万有引力定律、微积分等相关，爱因斯坦则与相对论、量子理论相关。研究还列出 21 世纪知名度前 20 的物理学家，爱因斯坦、普朗克、牛顿居前三。论文认为谷歌相关工具可作替代计量工具衡量科学家的社会影响力，也指出语料库语言偏见、姓名消歧误差等局限性，建议未来拓展研究领域

# 第三讲

## 1

即梦词云



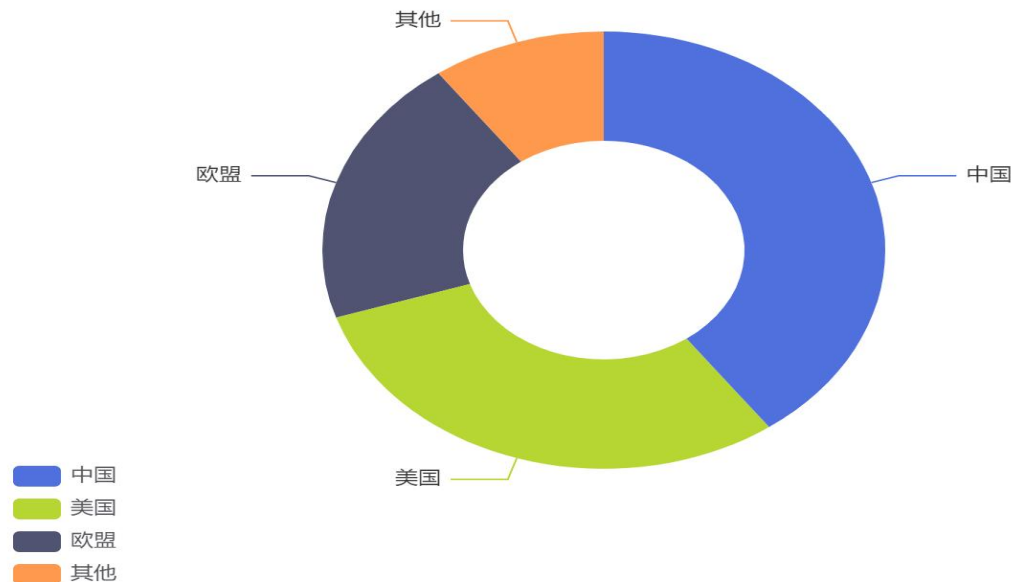
该词云图围绕“人工智能发展”主题，提取了当前AI领域最具代表性的关键词。图中“人工智能”作为核心词汇居中放大，周围环绕着“深度学习”“大模型”“AIGC”等高频热词，体现了技术演进的焦点。整体配色采用蓝绿渐变，象征科技与未来感，突出AI技术的前沿性与广泛影响力。

2

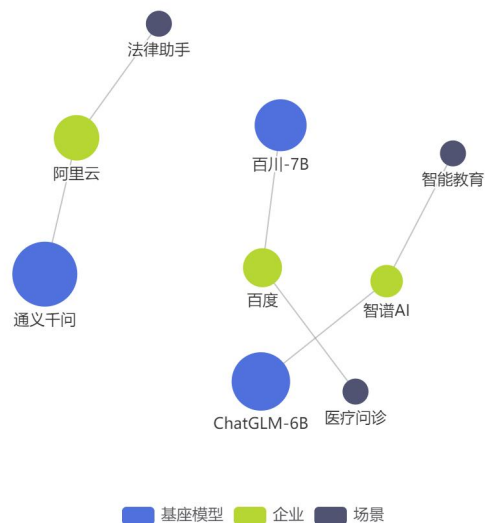
条形图：AI 岗位年薪 45–55 万，模型越新越值钱。  
饼图：中国 AI 论文量全球 4 成，领先美欧。  
关系图：基座—企业—场景三层链，大模型生态一目了然。

2023 全球 AI 论文国家占比

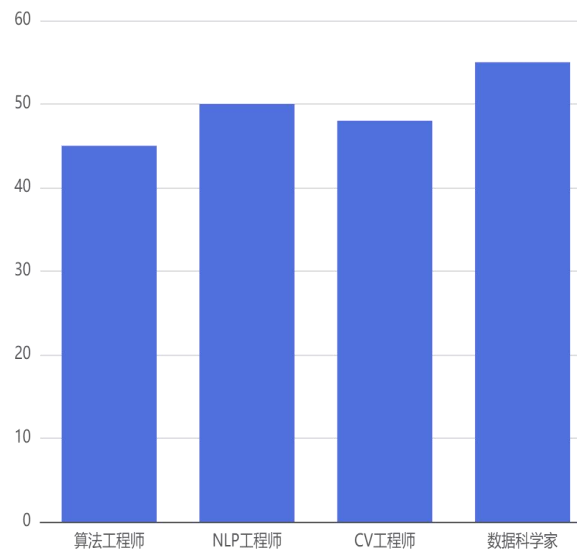
Fake Data



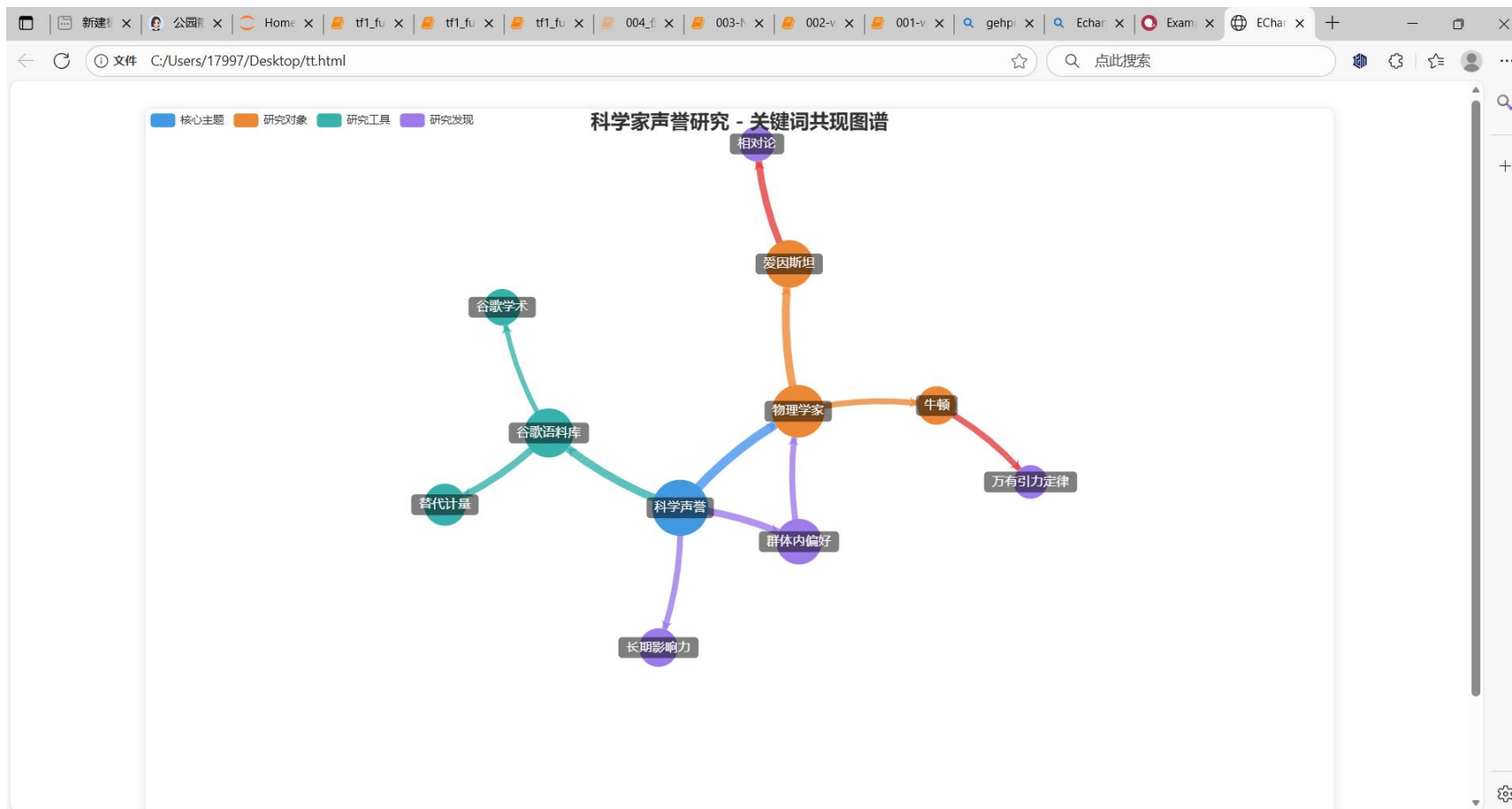
国内大模型生态关系图



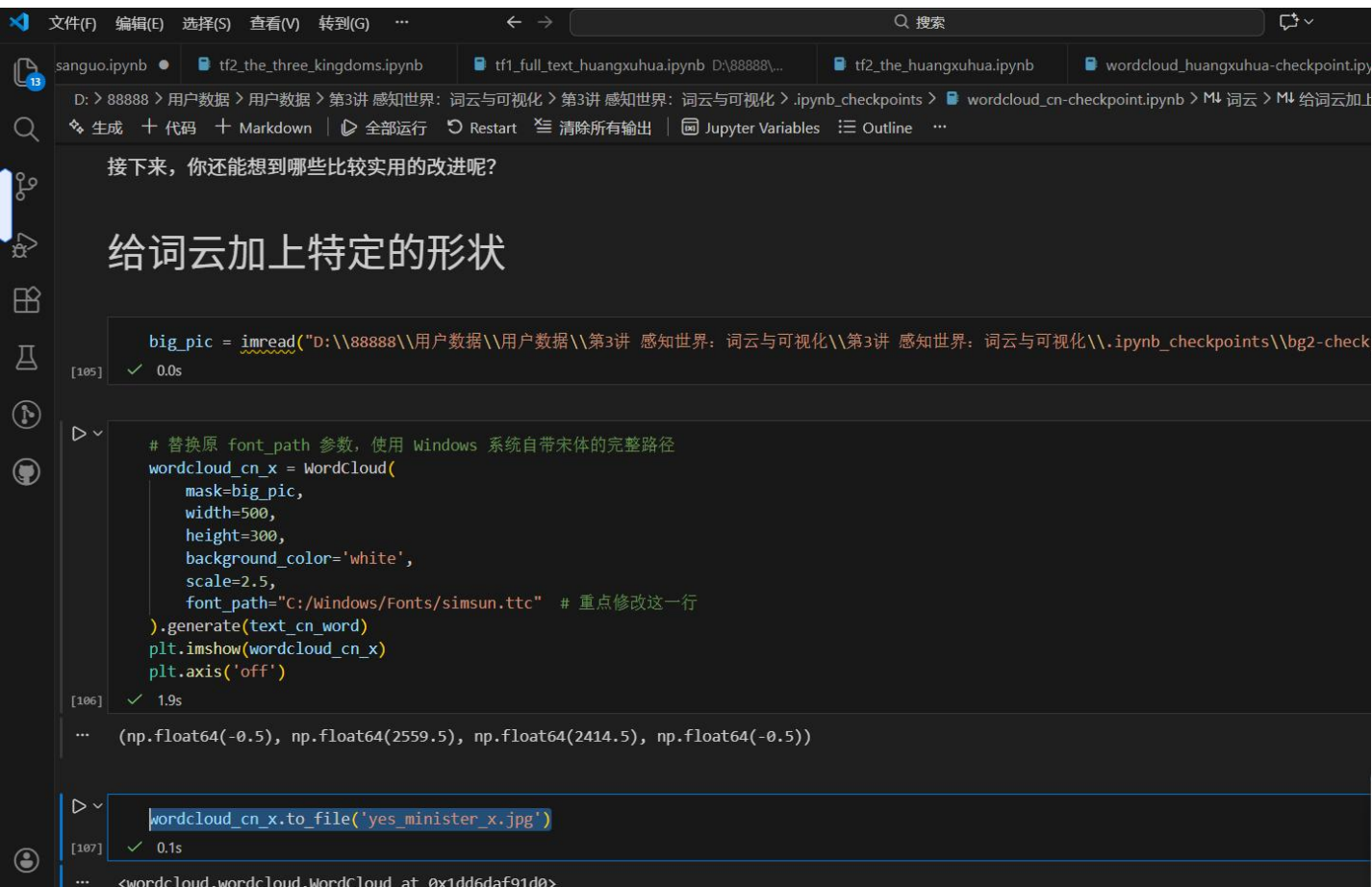
万元



3



4





# 第四讲

## • 1

← ↻ <https://hanlp.hankcs.com/demos/sentiment.html?text=这款+Python+数据分析库的+API+设计得太人性化了，几行代码就能完成数据清洗和可视化，帮我高效搞定了课程大作业的数据分析部分，强烈...> ☆ 1 ☆

在线演示

HanLP 安装 演示 文档 书籍 引用 论坛 博客 API

🔍 搜索

中文分词

词性标注

命名实体识别

依存句法分析

成分句法分析

语义依存分析

语义角色标注

抽象意义表示

指代消解

语义文本相似度

文本风格转换

关键词短语提取

抽取式自动摘要

请输入一段中文文本：

这款 Python 数据分析库的 API 设计得太人性化了，几行代码就能完成数据清洗和可视化，帮我高效搞定了课程大作业的数据分析部分，强烈推荐给入门学习者！

78/1000

情感分析

此页内容

简介

调用方法

创建客户端

情感分析

本地调用

多语种支持

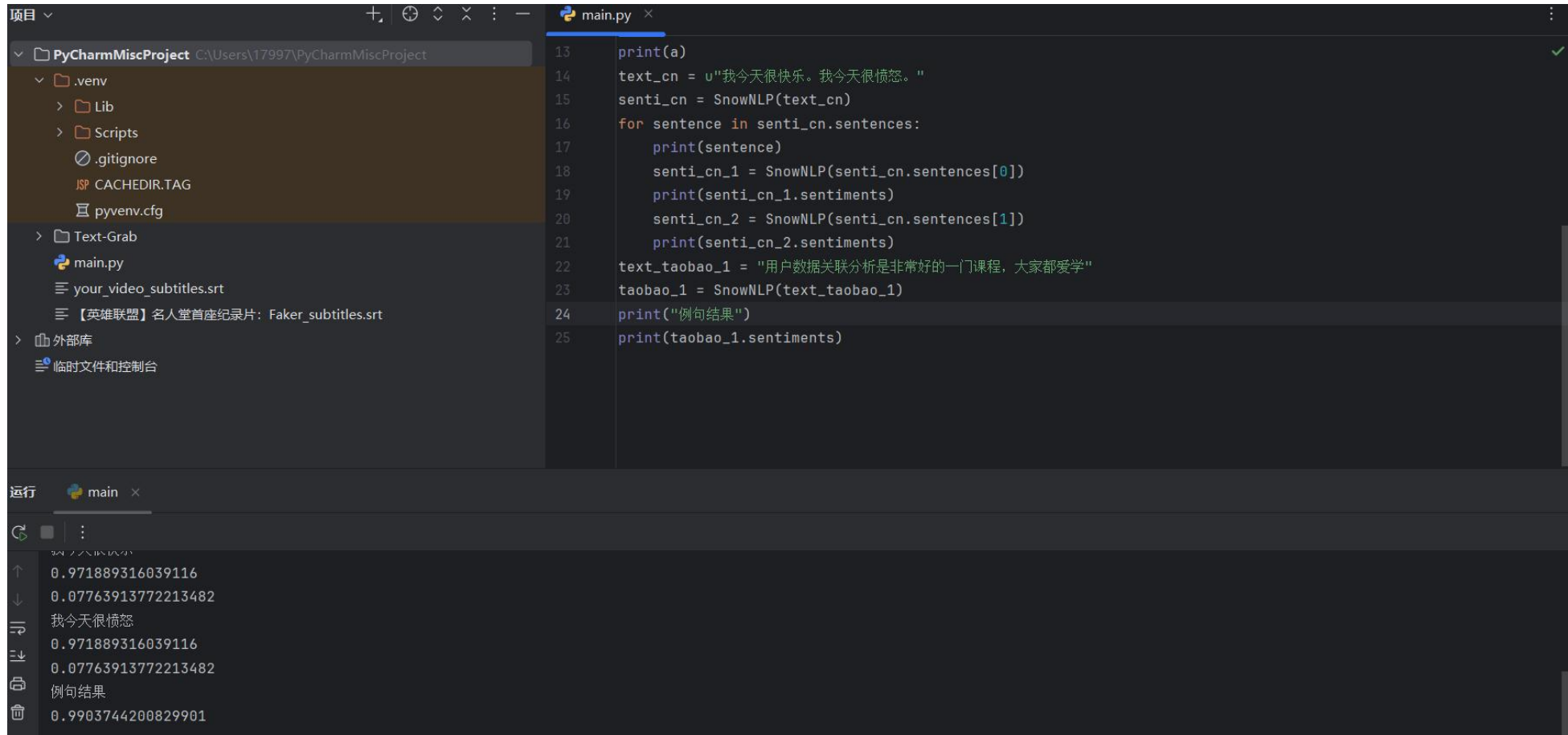
情感极性



情感极性



# 2 sentiment\_analysis\_1



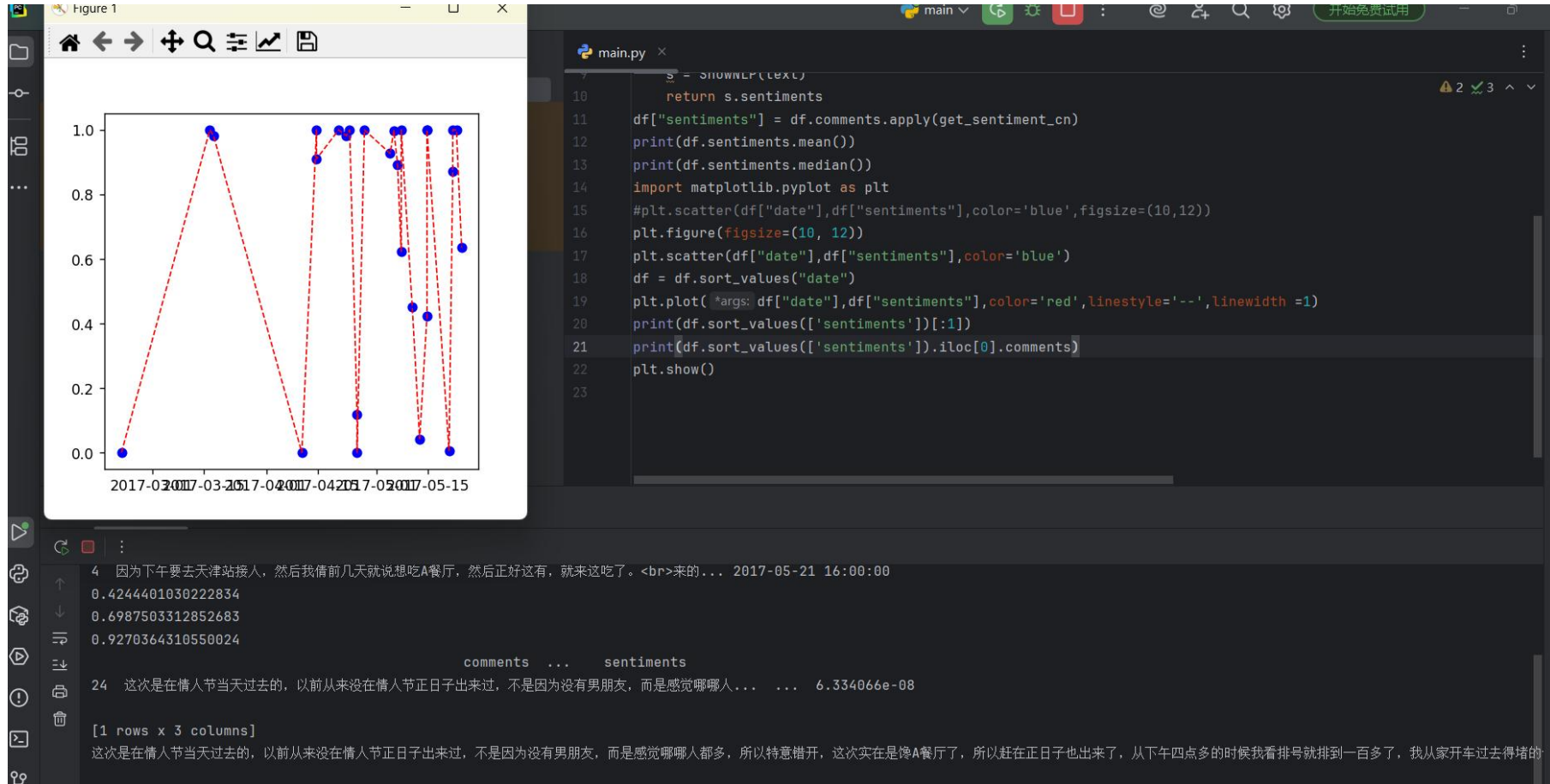
The screenshot displays the PyCharm IDE interface. The left sidebar shows the project structure for 'PyCharmMiscProject', including a virtual environment (.venv) and a file named 'main.py'. The main editor window shows the code in 'main.py', which uses the SnowNLP library for sentiment analysis. The code processes two sentences: '我今天很快乐。我今天很愤怒。' and '用户数据关联分析是非常好的一门课程，大家都爱学'.

```
13 print(a)
14 text_cn = u"我今天很快乐。我今天很愤怒。"
15 senti_cn = SnowNLP(text_cn)
16 for sentence in senti_cn.sentences:
17 print(sentence)
18 senti_cn_1 = SnowNLP(senti_cn.sentences[0])
19 print(senti_cn_1.sentiments)
20 senti_cn_2 = SnowNLP(senti_cn.sentences[1])
21 print(senti_cn_2.sentiments)
22 text_taobao_1 = "用户数据关联分析是非常好的一门课程，大家都爱学"
23 taobao_1 = SnowNLP(text_taobao_1)
24 print("例句结果")
25 print(taobao_1.sentiments)
```

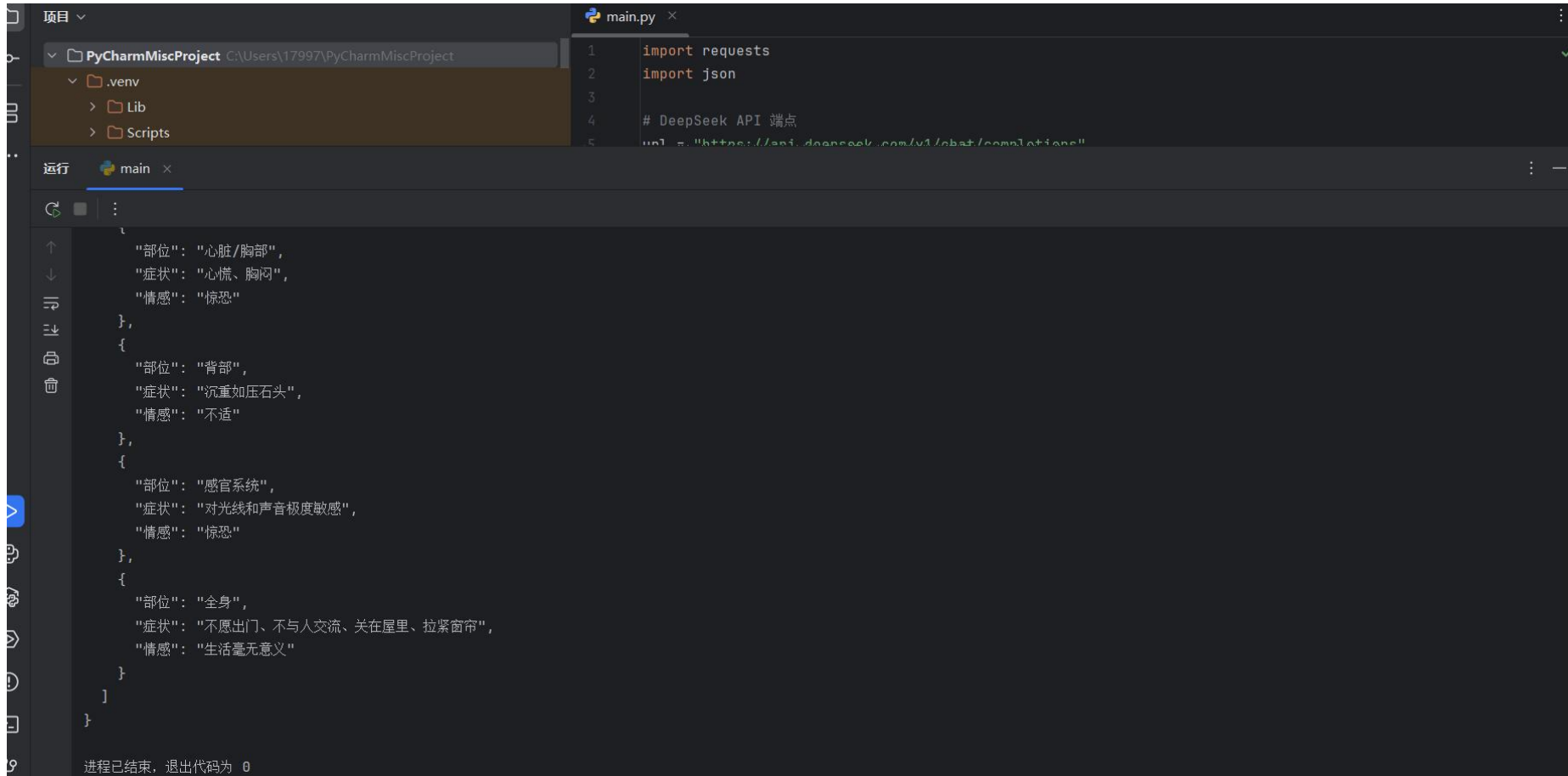
The bottom panel shows the output of the program, displaying sentiment scores and the original sentences.

```
0.971889316039116
0.07763913772213482
我今天很愤怒
0.971889316039116
0.07763913772213482
例句结果
0.9903744200829901
```

# sentiment\_analysis\_2



# sentiment\_analysis\_3

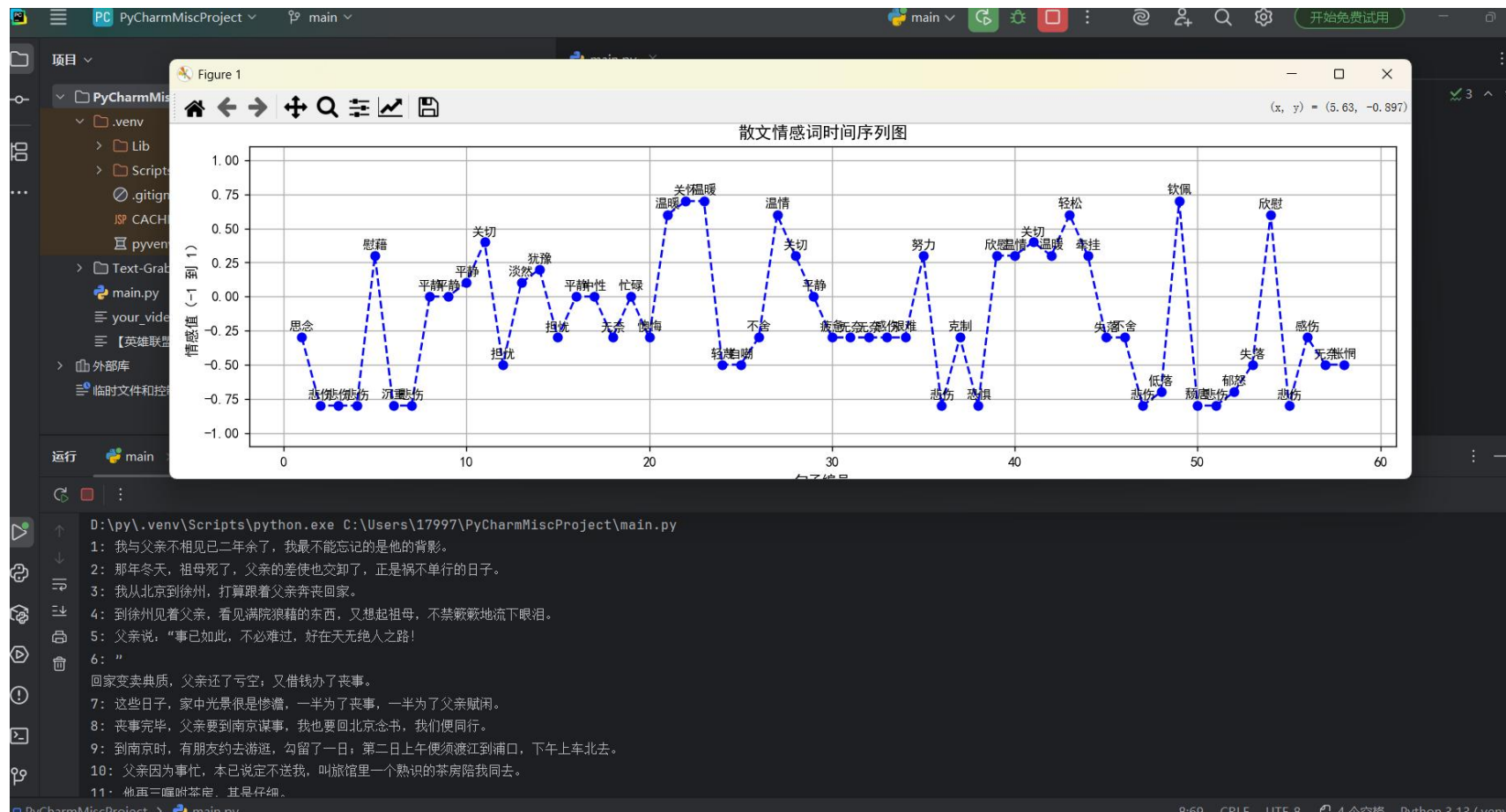


```
main.py
1 import requests
2 import json
3
4 # DeepSeek API 端点
5 url = "https://api.deepseek.com/v1/chat/completions"
```

```
{
 "部位": "心脏/胸部",
 "症状": "心慌、胸闷",
 "情感": "惊恐"
},
{
 "部位": "背部",
 "症状": "沉重如压石头",
 "情感": "不适"
},
{
 "部位": "感官系统",
 "症状": "对光线和声音极度敏感",
 "情感": "惊恐"
},
{
 "部位": "全身",
 "症状": "不愿出门、不与人交流、关在屋里、拉紧窗帘",
 "情感": "生活毫无意义"
}
]
```

进程已结束，退出代码为 0

# sentiment\_analysis\_4



# 总结分析

- 四段代码的功能都是对文本进行情感分析，在运行时发现打印图表时电脑加载比较慢。在将代码移动到pycharm中运行时常常遇到缩进错误等问题，需要手动调整代码缩进。总的来看使用这几段python代码进行情感分析的过程中虽有问题，但都能够很快解决并得出不错的结果，可以看出python在情感分析功能方面具有很好的作用

# 第六讲（1）美团大脑知识图谱生态构建分析

## 核心定位与生态规模定位

美团大脑是本地生活服务领域的超级知识中枢，以“全域数据整合 + AI 技术赋能”为核心，打通人、店、商品、场景的关联，支撑美团从“外卖平台”向“综合性生活服务超级 App”升级。

核心生态规模（截至 2025 年最新数据）

知识覆盖：33 类核心概念（含餐饮、酒店、零售、出行等）

实体量级：30 亿 + 业务实体（商户、商品、用户、景点等）

关联强度：1000 亿 + 三元组（如“用户 - 偏好 - 川菜”“商户 - 提供 - 到店服务”）

数据底座：40 亿 + 累计用户评价 + 即时零售 / 无人配送场景动态数据

# 最新生态构建进展

## 1. 技术架构升级：从“静态图谱”到“全域协同大脑”

融合美团 BERT 自研模型与多模态感知技术，实现“文本理解 + 场景识别”双驱动

打通低空无人机、地面无人车、楼宇配送机器人的数据链路，构建“空地楼”全域决策网络

升级 Listwise 排序模型，在 MRR@10 指标上保持行业领先，支撑复杂搜索需求解析

## 2. 场景生态拓展：从“餐饮外卖”到“全场景即时服务”

关联 13 万 + 可配送商品，实现“关键词搜索 → 需求识别 → 履约匹配”全链路支撑（如搜索“口罩”自动匹配附近药店 + 无人配送）

支持长句搜索，精准拆解场景、人群、品类需求

为校园、园区、机场等提供定制化图谱服务，适配无人配送路径规划

## 3. 数据生态激活：从“被动存储”到“主动赋能”

通过 NLP 技术解析 40 亿 + 用户评价，提取菜品、服务、环境等维度偏好，反哺商户经营优化

动态更新实体关联：实时同步商户营业时间、商品库存、配送范围等信息，保证图谱时效性

# 典型应用案例

## 案例背景

用户搜索“疫情期间适合遛娃的户外景点”，需同时满足“场景安全”“亲子适配”“即时可达”三大需求。

## 图谱赋能流程

需求拆解：通过知识图谱识别“疫情安全”“遛娃”“户外景点”“即时可达”四大核心实体

关联匹配：调用 景点 - 属性 - 安全评级，景点 - 适配人群 - 儿童，景点 - 配送范围 - 3 公里内”等三元组

履约衔接：匹配就近无人车配送站点，提供“景点门票购买 + 亲子套餐预订 + 即时物资配送”一体化服务

## 案例效果

搜索准确率提升 40%，长句需求满足率达 85%

无人配送履约效率提升 25%，核心场景配送时效压缩至 30 分钟内



# 生态评价与展望

## 优势评价

AI 团队与业务团队一体化架构，避免技术浪费，如无人配送数据快速反哺图谱优化  
从被动响应需求 到 主动预判，显著降低服务决策成本

## 现存挑战

跨场景数据治理复杂度高，部分领域（如医疗健康）实体关联精度待提升  
无人配送场景的动态数据更新频率，需匹配图谱实时计算能力  
个性化推荐具有负面效果风险，需平衡用户偏好与探索性需求

## 未来展望

深化 具身智能 + 知识图谱 融合，提升无人配送场景的动态决策能力  
开放部分图谱能力给第三方商户，提供 用户画像 - 商品优化 - 履约效率”定制化解决方案  
拓展 本地生活 + 城市服务 关联，如将交通、政务等实体纳入图谱，打造城市级服务中枢

# 总结与评价

美团大脑的最新生态构建，核心是以知识图谱为纽带，打通 数据 - 技术 - 场景 - 履约 的通路。既支撑了 C 端用户的复杂需求，也能够支持 B 端商户的数字化经营和无人配送的规模化落地。未来若能持续提升跨领域实体关联精度与开放生态建设，有望成为生活服务领域的基础设施中枢。

优点：贴合外卖、到店消费等日常场景，搜店、点餐更精准，还能帮商家算备货量、提经营建议，功能多样且易于使用。

不足：跨领域能力有待完善，还面临用户对于数据安全和隐私保护可靠性的质疑。

# 第六讲 (续)

1

openkg - 搜索 × 开放图谱 - 开放知识图谱 × 开放图谱 - 开放知识图谱 × 孤独症谱系障碍知识库 × 沉默蛇/ASDKB × Zhishi.me - 数据集 - 开 × 图书问答知识图谱 - 数 ×

https://github.com/SilenceSnake/ASDKB

README.md 更新 README.md 3 years ago

阅读 CC-BY-SA-4.0 许可证

## AsdKB: 中文自闭症谱系障碍早期筛查与诊断知识库

(Text) DSM-5

(Structured Knowledge) SNOMED CT The global language of healthcare

(Text) ICD-10 The ICD-10 Classification of Mental and Behavioural Disorders

(Tables) 好大夫在线 https://m.haodf.com

(Tables) 国家统计局 National Bureau of Statistics

Physicians and Hospitals

家庭医生在线 familydoctor.com.cn

Intervention Information

(Tables) The National Clearinghouse on Autism Evidence & Practice NCAEP • BRIDGING SCIENCE AND PRACTICE

ALS LIFE 中国自闭症评估干预平台

CDC Centers for Disease Control and Prevention

OCAALI (Tables)

AUTISM CANADA AUTISME

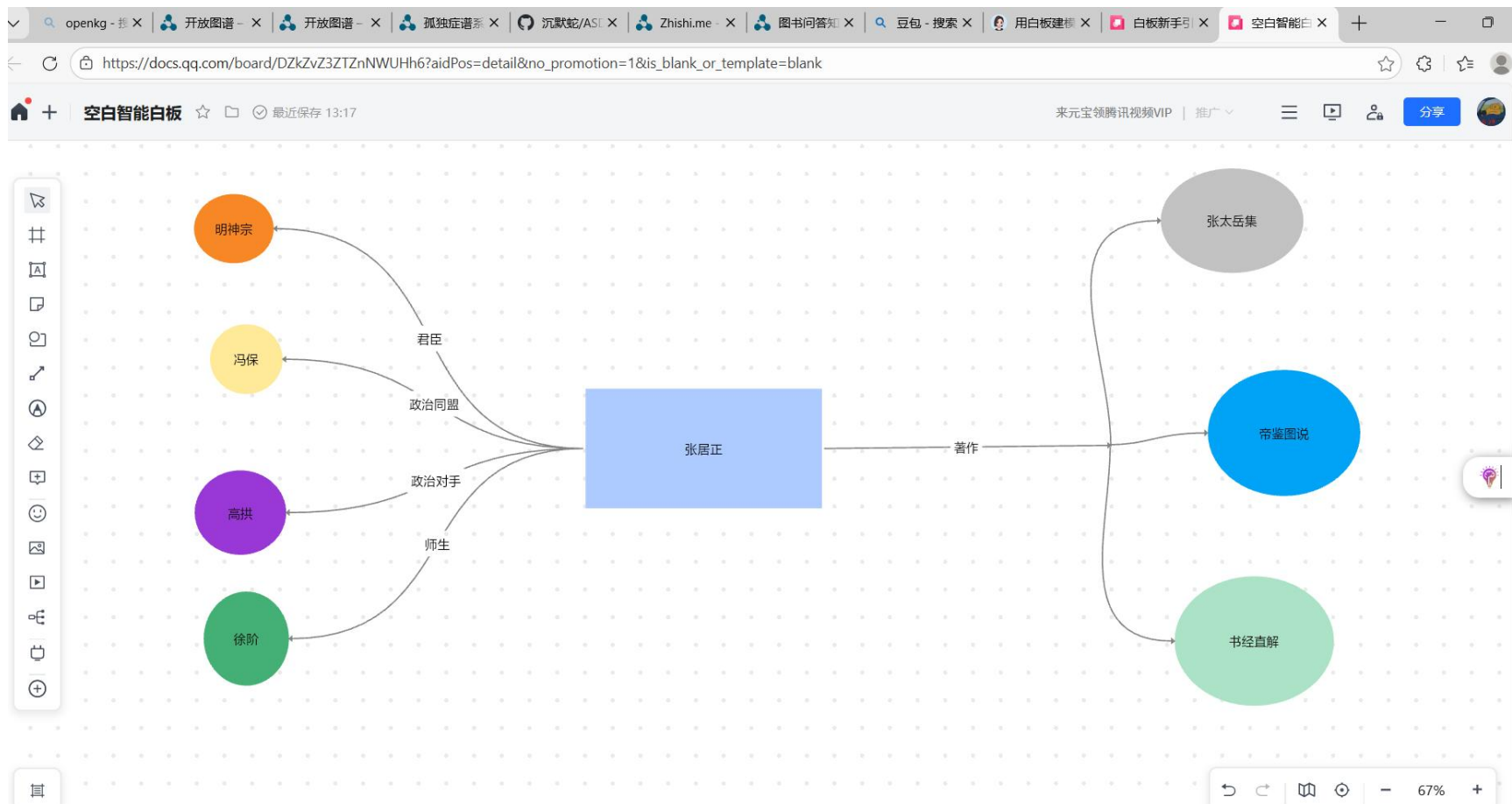
发行作品

未发布任何版本

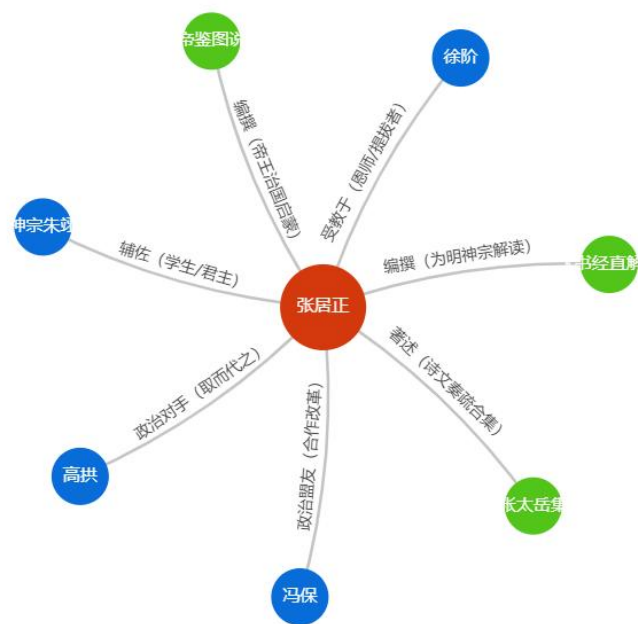
套餐

没有发布任何软件包

2



张居正知识图谱



4(neo4j)

