数据科学导论 -HW 5 报告

孙育泉 10234900421 2025.10.26

总览

Ι	实验	量求・						 	 	 	 	. 1
Π	具体	文现.						 	 	 	 	. 1
	II.1	准备コ	[作					 	 	 	 	. 1
		II.1.1	定义基	本的参	数、豆	复用图	函数	 	 	 	 	. 1
		II.1.2	数据准	备				 	 	 	 	3
	II.2	数据分	分析与建					 	 	 	 	. 4
		II.2.1	分析.					 	 	 	 	. 4
				i模型								
	II.3	学科画	画像					 	 	 	 	7
		II.3.1	分析.					 	 	 	 	. 7
		II.3.2	思路与	i模型 ·				 	 	 	 	8
	II.4	学科技	非名预测	∥模型・				 	 	 	 	12
		II.4.1	分析.					 	 	 	 	12
		II.4.2	思路与	i模型 ·				 	 	 	 	12

I实验要求

- (1) 请结合这些学科排名数据,分析全球高校可以大致分为哪几类?并且分析出与华师大类似 的高校?
- (2) 请通过探索性分析的方式,对华东师范大学做一个学科画像?用尽可能多的角度去做.
- (3) 请利用数据建模的方式,对各学科做一个排名模型,能够较好的预测出排名位置.(可以用 各学科前60%的数据作为训练集,后20%的数据作为测试集)

II 具体实现

II.1 准备工作

II.1.1 定义基本的参数、可复用函数

首先, 我们需要导入必要的库, 并加载数据集.

我们使用 pandas 来处理数据, numpy 来进行数值计算, matplotlib 和 seaborn 来进行数据可视 化,使用 scikit-learn 来进行机器学习建模.

In [8]:

```
python
2 # 单元格 1: 导入所有库
4 import sqlite3
5 import pandas as pd
6 import numpy as np
7 import matplotlib.pyplot as plt
8 import seaborn as sns
```



```
9 import warnings
10
11 # Sklearn 库
12 from sklearn.preprocessing import StandardScaler, OneHotEncoder
13 from sklearn.cluster import KMeans
14 from sklearn.metrics.pairwise import euclidean distances
15 from sklearn.model_selection import train_test_split
16 from sklearn.ensemble import RandomForestRegressor
17 from sklearn.metrics import mean_squared_error, r2_score
18 from sklearn.compose import ColumnTransformer
19 from sklearn.pipeline import Pipeline
21 # 设置 Matplotlib 和 Seaborn 风格
22 sns.set_theme(style="whitegrid")
23 plt.rcParams['font.sans-serif'] = ['Source Han Sans SC'] # 用来正常显示中文标签
24 plt.rcParams['axes.unicode_minus'] = False # 用来正常显示负号
25 warnings.filterwarnings("ignore", category=FutureWarning)
26 warnings.filterwarnings("ignore", category=UserWarning)
27
28 print("所有库导入成功.")
                                                                                   txt
 1 所有库导入成功.
```

接下来,我们定义一下全局的配置,包括数据库文件所在的路径、院校的名称等.

然后,是一些数据预处理的工作,包括从数据库加载数据、预处理数据为:按照每个机构的名字进行聚合数据.



```
9
          :param db file: 数据库文件名
10
          :param table name: 表名
11
          :return: 包含原始数据的 pandas DataFrame
12
13
       print(f"正在从 {db_file} 加载数据...")
14
       try:
          conn = sqlite3.connect(db_file)
15
          df = pd.read_sql_query(f"SELECT * FROM {table_name}", conn)
16
17
          conn.close()
18
          print(f"数据加载成功,共 {len(df)} 条记录.")
19
          return df
       except Exception as e:
20
21
          print(f"数据加载失败: {e}")
22
          return pd.DataFrame() # 返回空 DataFrame
23
24 def get_aggregated_data(df):
25
26
          将原始 DataFrame 按机构聚合,构建高校画像.
27
28
          :param df: 原始 DataFrame (来自 load_data)
          :return: (df_agg, features)
29
                           df_agg: 聚合后的 DataFrame
30
31
                           features: 用于聚合的特征列表
32
33
       print("正在按机构聚合数据以构建高校画像...")
34
       if df.empty:
35
          print("输入数据为空,无法聚合.")
          return pd.DataFrame(), []
36
37
38
       features = [
          'total_subjects', 'avg_rank', 'total_documents',
39
40
           'total_cites', 'total_top_papers', 'avg_cites_per_paper'
41
42
43
       df agg = df.groupby('institution').agg(
44
          total_subjects=pd.NamedAgg(column='research_field', aggfunc='count'),
45
          avg_rank=pd.NamedAgg(column='rank', aggfunc='mean'),
          total_documents=pd.NamedAgg(column='documents', aggfunc='sum'),
46
47
          total_cites=pd.NamedAgg(column='cites', aggfunc='sum'),
          total_top_papers=pd.NamedAgg(column='top_papers', aggfunc='sum'),
48
          avg_cites_per_paper=pd.NamedAgg(column='cites_per_paper', aggfunc='mean')
49
       ).reset_index()
50
51
52
       print(f"数据聚合完成,共 {len(df_agg)} 所机构.")
53
       return df_agg, features
54
55 print("共用函数 (load_data, get_aggregated_data) 定义完成.")
                                                                                  txt
 1 共用函数 (load_data, get_aggregated_data) 定义完成.
 2
```

II.1.2 数据准备

接下来, 我们导入数据库中的数据, 并进行简单的展示.

```
In [11]: python
```



```
1 # -----
2 # 单元格 4: 加载原始数据
4 # 调用单元格 3 中定义的函数
5 df_raw = load_data(DB_FILE, TABLE_NAME)
6
7 # 显示数据样例
8 if not df_raw.empty:
     print("\n数据预览 (前5行):")
10
      print(df_raw.head())
                                                                          txt
1 正在从 esi_rankings.db 加载数据...
 2 数据加载成功,共 34121 条记录.
4 数据预览 (前5行):
5
    id
              research field rank \
6 0
     1 AGRICULTURAL SCIENCES
     2 AGRICULTURAL SCIENCES
     3 AGRICULTURAL SCIENCES
8 2
                              3
     4 AGRICULTURAL SCIENCES
                              4
10 4 5 AGRICULTURAL SCIENCES
11
                      institution country_region documents \
12
                     CHINESE ACADEMY OF SCIENCES CHINA MAINLAND
13 0
                                                               15661
14 1
          CHINESE ACADEMY OF AGRICULTURAL SCIENCES CHINA MAINLAND
                                                               12222
15 2 UNITED STATES DEPARTMENT OF AGRICULTURE (USDA)
                                                               12564
                   CHINA AGRICULTURAL UNIVERSITY CHINA MAINLAND
16 3
                                                               10052
17 4
                                        INRAE
                                                     FRANCE
                                                                9314
18
19
     cites cites_per_paper top_papers
20 0 332254
             21.22
21 1 223855
                    18.32
                                198
                   17.56
22 2 220644
                                105
23 3 207779
                    20.67
                                166
24 4 187838
                    20.17
                                118
25
```

II.2 数据分析与建模

II.2.1 分析

这个任务包含两个子任务:

- 全球高校分类: 我们想知道全球高校可以"大致分为哪几类". 当数据中没有现成的"类别"标签,而我们想根据数据的相似性自动找出"群体"时,这是一个聚类问题.
- 寻找相似高校: 我们想找到与"华师大"特征相似的其他高校. 这是一个相似性匹配或距离计算问题.

II.2.2 思路与模型

我的思路是首先将数据从"学校-学科"粒度聚合到"学校"粒度,为每所高校构建一个"画像",然后基于这个画像同时解决上述两个问题.

• 数据聚合:

▶ 原始数据是 [机构, 学科, 排名, ...].



- ▶ 我需要将其转换为 [机构,总入围学科数,平均排名,总论文数,总引用数,总高被引论文数,平均篇均引用].
- 模型选型 (聚类):
 - ▶ 算法: 我选用 K-Means 聚类算法.
 - ▶ 原因: K-Means 是一种高效且直观的聚类算法,适合处理当前这种"特征画像"数据.
 - ▶ 前置处理: K-Means 是基于距离的,而"总引用数"(数百万)和"总学科数"(数十)的 尺度完全不同. 因此,在聚类前,我必须使用 StandardScaler 对所有特征进行标准化,使 其均值为 0,方差为 1,消除尺度影响.
- 模型选型 (相似性):
 - ▶ 算法: 我选用基于欧氏距离的相似度计算.
 - ► 原因:在聚类时,我们已经对数据进行了标准化.在同一个标准化空间中,欧氏距离是衡量两个点(两所学校)"绝对相似度"的最直观方法.距离越近,代表两所学校的画像越相似.
 - ▶ 步骤:
 - (1) 提取"华师大"的标准化特征向量.
 - (2) 计算它与所有其他学校向量的欧氏距离.
 - (3) 排序, 找出距离最小的(即最相似的).

In [12]:

```
python
2 # 单元格 5: 任务1 - 高校分类与相似高校分析
3 # -----
4 print("--- 开始执行任务1: 高校分类与相似高校分析 ---")
6 # 1. 分析 (Analysis)
7 print("\n[分析]: 这是一个聚类 (Clustering) 和 相似性匹配 (Similarity Matching) 问题.")
9 # 2. 思路 (Approach)
10 print("[思路]:")
11 print(" - 调用 get_aggregated_data() 获取机构画像.")
12 print(" - 使用 StandardScaler 标准化数据 (消除量纲影响). ")
13 print(" - 使用 K-Means (K=5) 进行聚类分析.")
14 print(" - 使用标准化后的欧氏距离 (Euclidean Distance) 寻找相似高校. ")
16 # 3. 代码实现 (Code)
17
18 # 3.1. 获取聚合数据
19 df_agg, agg_features = get_aggregated_data(df_raw)
20
21 if not df_agg.empty:
22
     # 3.2. 准备特征并标准化
23
      df_agg_features = df_agg.set_index('institution')[agg_features].fillna(0)
24
25
      scaler = StandardScaler()
26
      features_scaled = scaler.fit_transform(df_agg_features)
      df scaled = pd.DataFrame(features scaled, index=df agg features.index,
27
      columns=agg features)
```



```
print("\n步骤 1/4: 特征标准化完成.")
28
29
30
      # 3.3. K-Means 聚类
      K = 5
31
32
      kmeans = KMeans(n_clusters=K, random_state=42, n_init=10)
      df_agg['cluster'] = kmeans.fit_predict(features_scaled)
33
34
      print(f"步骤 2/4: K-Means 聚类完成 (K={K}). ")
35
36
      # 3.4. 分析聚类结果
37
      cluster_analysis = df_agg.groupby('cluster')[agg_features].mean()
38
      cluster_analysis['count'] = df_agg['cluster'].value_counts()
39
      print("\n--- 任务1 结果:全球高校分类(聚类中心均值)---")
40
      print(cluster_analysis)
41
42
      # 3.5. 寻找相似高校(欧氏距离)
43
      print(f"\n步骤 3/4: 正在寻找与 '{TARGET_INSTITUTION}' 相似的高校...")
44
      try:
          target vector scaled = df scaled.loc[TARGET INSTITUTION].values.reshape(1, -1)
45
46
47
          distances = euclidean_distances(target_vector_scaled, df_scaled.values)
48
          df_scaled['distance_to_target'] = distances[0]
49
50
          df_similar_euclidean = df_scaled.sort_values(by='distance_to_target',
          ascending=True)
51
52
          # 提取前11名(包含自己)
          top_similar_euclidean = df_similar_euclidean.iloc[0:11]
53
54
55
          # 从 df agg features (原始聚合数据) 中提取信息以提高可读性
          original_data_similar = df_agg_features.loc[top_similar_euclidean.index]
56
          original_data_similar['distance'] =
57
          top_similar_euclidean['distance_to_target']
58
59
          print(f"\n--- 任务1 结果:与 {TARGET INSTITUTION} 最相似的高校(欧氏距离)---")
          print(original_data_similar[['total_subjects', 'avg_rank',
60
          'avg_cites_per_paper', 'distance']])
61
          print("步骤 4/4: 相似高校分析完成.")
62
63
      except KeyError:
64
          print(f"错误:在数据中未找到机构 '{TARGET INSTITUTION}'.")
65
      except Exception as e:
66
          print(f"计算相似高校时出错: {e}")
67 else:
68
      print("聚合数据为空,跳过任务1.")
                                                                              txt
1 --- 开始执行任务1: 高校分类与相似高校分析 ---
2
 3 [分析]: 这是一个聚类 (Clustering) 和 相似性匹配 (Similarity Matching) 问题.
4 [思路]:
 5
    - 调用 get_aggregated_data() 获取机构画像.
6
    - 使用 StandardScaler 标准化数据 (消除量纲影响).
7
    - 使用 K-Means (K=5) 进行聚类分析.
   - 使用标准化后的欧氏距离 (Euclidean Distance) 寻找相似高校.
9 正在按机构聚合数据以构建高校画像...
10 数据聚合完成, 共 9990 所机构.
```



```
12 步骤 1/4: 特征标准化完成.
13 步骤 2/4: K-Means 聚类完成 (K=5).
15 --- 任务1 结果:全球高校分类 (聚类中心均值) ---
total_subjects avg_rank total_documents total_cites \
17 cluster
              2.907097 1519.894761
18 0
                                      2490.407903 6.155025e+04
             20.763889 163.341380 139027.541667 4.048206e+06
19 1
              1.033791 4887.824182
              1.000000 4919.506849 1.493151 7.1000000 7.974747e+03 32236.870041 8.163620e+05
                                      216.190426 7.974747e+03
20 2
             15.757866 587.575101
21 3
22 4
23
   total_top_papers avg_cites_per_paper count
25 cluster
                                 63.168193 6200
26 0
               52.783710
              3679.958333
27 1
                                  29.581056
                                            72
28 2
               7.057022
                                332.615831 2841
29 3
              688.727770
                                 25.208109 731
                               5790.577123 146
30 4
               1.321918
32 步骤 3/4: 正在寻找与 'EAST CHINA NORMAL UNIVERSITY' 相似的高校...
34 --- 任务1 结果:与 EAST CHINA NORMAL UNIVERSITY 最相似的高校(欧氏距离)---
35 total_subjects avg_rank \
36 institution
37 EAST CHINA NORMAL UNIVERSITY
                                            17 601.764706
                                           17 569.411765
38 UNIVERSITY OF GENOA
39 GOETHE UNIVERSITY FRANKFURT
                                           16 498.562500
40 UNIVERSITY OF TENNESSEE KNOXVILLE
                                          17 447.176471
41 HEBREW UNIVERSITY OF JERUSALEM
                                           18 499.500000
42 LAVAL UNIVERSITY
                                           18 421.277778
43 UNIVERSITY OF DUISBURG ESSEN
                                           18 555,277778
                                           18 603.277778
44 NANCHANG UNIVERSITY
45 KYUNG HEE UNIVERSITY
                                           17 520.000000
                                           18 424.333333
46 GRIFFITH UNIVERSITY
47 LEIPZIG UNIVERSITY
                                           16 546.937500
48
               avg_cites_per_paper distance
50 institution
51 EAST CHINA NORMAL UNIVERSITY
                                          20.548824 0.000000
52 UNIVERSITY OF GENOA
                                          21.984118 0.191479
53 GOETHE UNIVERSITY FRANKFURT
                                        24.761250 0.322667
                                        22.390000 0.325080
54 UNIVERSITY OF TENNESSEE KNOXVILLE
55 HEBREW UNIVERSITY OF JERUSALEM
                                        22.439444 0.344547
56 LAVAL UNIVERSITY
                                        23.129444 0.351795
57 UNIVERSITY OF DUISBURG ESSEN
                                         23.712778 0.378091
58 NANCHANG UNIVERSITY
                                         18.113889 0.381699
                                         20.160588 0.383135
59 KYUNG HEE UNIVERSITY
60 GRIFFITH UNIVERSITY
                                         28.681111 0.384137
61 LEIPZIG UNIVERSITY
                                          24.103125 0.387179
62 步骤 4/4: 相似高校分析完成.
```

II.3 学科画像

II.3.1 分析

这个任务要求对华师大做"学科画像",并"用尽可能多的角度".



II.3.2 思路与模型

我的思路是,从原始的"学校-学科"数据中,筛选出所有"华师大"的记录,然后从多个维度进行描述和可视化。

- 数据筛选: 从原始 df 中过滤 institution == 'EAST CHINA NORMAL UNIVERSITY'.
- 分析角度:
 - ▶ 学科广度: 总共有多少个学科入围? (对 research field 计数).
 - ► 学科高度: 哪些学科排名最靠前? (按 rank 升序排序).
 - ▶ 学科影响力:哪些学科的高被引论文 (top papers) 最多? (按 top papers 降序排序).
 - ▶ 学科学术质量:哪些学科的篇均引用 (cites_per_paper) 最高? (按 cites_per_paper 降序排序).

• 可视化:

- ► 算法/工具: 使用 matplotlib 和 seaborn 库.
- ▶ 图表:对上述角度 2, 3, 4,分别使用条形图进行可视化.条形图最适合展示不同类别(学科)在某个数值指标(排名、高被引、篇均引)上的对比.

In [13]:

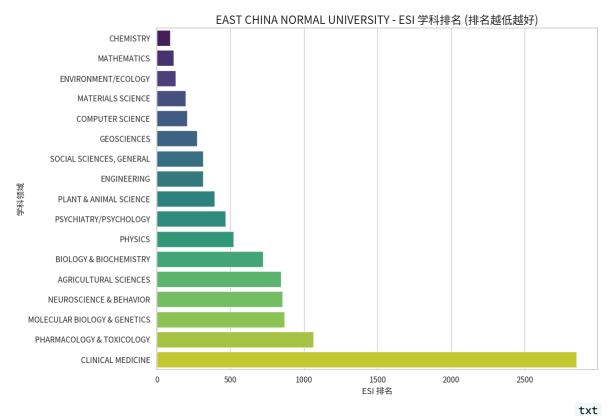
```
python
2 # 单元格 6: 任务2 - 华东师范大学学科画像 (EDA)
3 # -----
4 print("--- 开始执行任务2: 华东师范大学学科画像 ---")
6 # 1. 分析 (Analysis)
7 print("\n[分析]: 这是一个探索性数据分析 (EDA) 问题.")
9 # 2. 思路 (Approach)
10 print("[思路]:")
11 print(" - 从 df raw 筛选出华师大所有记录.")
12 print(" - 从 广度 (学科数)、高度 (排名)、影响力 (Top Papers)、质量 (篇均引) 四个角度分析.
13 print(" - 使用 Seaborn 条形图 (barplot) 进行可视化.")
15 # 3. 代码实现 (Code)
16
17 # 3.1. 筛洗数据
18 df_ecnu = df_raw[df_raw['institution'] == TARGET_INSTITUTION].copy()
19 print(f"\n步骤 1/4: 筛选 '{TARGET INSTITUTION}' 数据...")
20
21 if df_ecnu.empty:
22 print("未找到华师大数据,无法生成画像.")
23 else:
24
      # 3.2. 角度1: 学科广度
25
      total_subjects_ecnu = df_ecnu['research_field'].nunique()
26
      print(f"步骤 2/4: 学科广度分析完成. 共 {total subjects ecnu} 个学科进入 ESI 前 1%. ")
27
28
     # 3.3. 角度2/3/4:排序分析
29
      df_ecnu_sorted_rank = df_ecnu.sort_values(by='rank', ascending=True)
      df ecnu sorted top = df ecnu.sort values(by='top papers', ascending=False)
30
      df_ecnu_sorted_cpp = df_ecnu.sort_values(by='cites_per_paper', ascending=False)
31
32
      print("步骤 3/4: 学科排名、高被引、篇均引数据排序完成.")
```



```
33
34
      print("\n--- 任务2 结果:学科画像详情 ---")
35
      print(f"优势学科 (排名前3): \n{df_ecnu_sorted_rank[['research_field',
       'rank']].head(3)}")
36
      print(f"\n影响力学科 (Top Papers 前3): \n{df_ecnu_sorted_top[['research_field',
       'top_papers']].head(3)}")
37
      print(f"\n高质量学科 (篇均引用 前3): \n{df_ecnu_sorted_cpp[['research_field',
       'cites_per_paper']].head(3)}")
38
      # 3.4. 可视化
39
40
      print("\n步骤 4/4: 正在生成可视化图表...")
41
42
      # 图1: 各学科 ESI 排名
43
      plt.figure(figsize=(12, 8))
44
      sns.barplot(data=df_ecnu_sorted_rank, x='rank', y='research_field',
      palette='viridis')
45
      plt.title(f'{TARGET INSTITUTION} - ESI 学科排名 (排名越低越好)', fontsize=16)
46
      plt.xlabel('ESI 排名', fontsize=12)
47
      plt.ylabel('学科领域', fontsize=12)
48
      plt.tight_layout()
49
      plt.savefig('ecnu_subject_rankings.png')
50
      plt.show() # 在 Notebook 中显示
51
      print(" - 已生成学科排名图: ecnu subject rankings.png")
52
53
      # 图2: 各学科高被引论文数
54
      plt.figure(figsize=(12, 8))
55
      sns.barplot(data=df_ecnu_sorted_top, x='top_papers', y='research_field',
      palette='plasma')
56
      plt.title(f'{TARGET_INSTITUTION} - ESI 学科高被引论文数', fontsize=16)
57
      plt.xlabel('高被引论文数 (Top Papers)', fontsize=12)
58
      plt.ylabel('学科领域', fontsize=12)
59
      plt.tight layout()
      plt.savefig('ecnu_top_papers.png')
60
61
      plt.show() # 在 Notebook 中显示
62
      print(" - 已生成高被引论文图: ecnu top papers.png")
63
64
      # 图3:各学科篇均引用
65
      plt.figure(figsize=(12, 8))
      sns.barplot(data=df_ecnu_sorted_cpp, x='cites_per_paper', y='research_field',
66
      palette='coolwarm')
67
      plt.title(f'{TARGET_INSTITUTION} - ESI 学科篇均引用', fontsize=16)
68
      plt.xlabel('篇均引用 (Cites per Paper)', fontsize=12)
69
      plt.ylabel('学科领域', fontsize=12)
70
      plt.tight_layout()
71
      plt.savefig('ecnu_cites_per_paper.png')
72
      plt.show() # 在 Notebook 中显示
73
      print(" - 已生成篇均引用图: ecnu_cites_per_paper.png")
74
75
      print("\n--- 任务2 完成 ---")
                                                                                txt
 1 --- 开始执行任务2:华东师范大学学科画像 ---
 3 [分析]: 这是一个探索性数据分析 (EDA) 问题.
 4 [思路]:
 5
    - 从 df_raw 筛选出华师大所有记录.
   - 从 广度 (学科数)、高度 (排名)、影响力 (Top Papers)、质量 (篇均引) 四个角度分析.
```



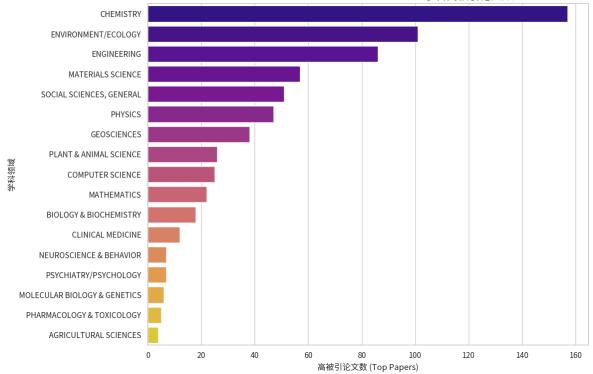
```
7 - 使用 Seaborn 条形图 (barplot) 进行可视化.
8
9 步骤 1/4: 筛选 'EAST CHINA NORMAL UNIVERSITY' 数据...
10 步骤 2/4: 学科广度分析完成. 共 17 个学科进入 ESI 前 1%.
11 步骤 3/4: 学科排名、高被引、篇均引数据排序完成.
12
13 --- 任务2 结果:学科画像详情 ---
14 优势学科 (排名前3):
15
             research_field rank
                CHEMISTRY
                           90
16 3119
17 22230
               MATHEMATICS
                            115
18 16247 ENVIRONMENT/ECOLOGY
                            130
19
20 影响力学科 (Top Papers 前3):
             research_field top_papers
21
22 3119
                 CHEMISTRY
23 16247 ENVIRONMENT/ECOLOGY
                                 101
24 13647
               ENGINEERING
                                  86
25
26 高质量学科 (篇均引用 前3):
                     research_field cites_per_paper
28 24180 MOLECULAR BIOLOGY & GENETICS
                  MATERIALS SCIENCE
                                            34.55
29 20731
30 16247
                ENVIRONMENT/ECOLOGY
                                            31.31
31
32 步骤 4/4: 正在生成可视化图表...
```



1 - 已生成学科排名图: ecnu_subject_rankings.png





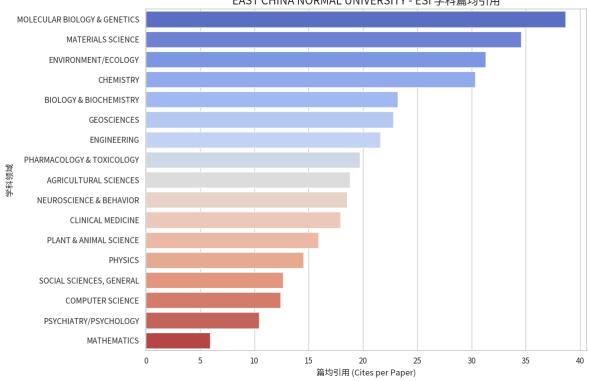


1 - 已生成高被引论文图: ecnu_top_papers.png

EAST CHINA NORMAL UNIVERSITY - ESI 学科篇均引用

txt

txt



- 已生成篇均引用图: ecnu_cites_per_paper.png

1 2



3 --- 任务2 完成 ---

II.4 学科排名预测模型

II.4.1 分析

这是一个典型的监督学习问题. 由于 rank 是一个连续的数值(或至少是高基数的有序数值),这具体来说是一个回归问题.

II.4.2 思路与模型

- 特征 (X) 和 目标 (y)
 - ▶ 目标 (y): rank
 - ▶ 特征 (X): documents, cites, cites_per_paper, top_papers (数值特征),以及 research_field (类别特征).
 - ▶ research_field 很重要因为 ESI 排名是"在特定学科内"排名的. 不同学科的竞争激烈程度 和数据分布(如临床医学有几千个机构上榜,而空间科学可能只有几百个)截然不同. 因此,模型必须知道它正在为哪个学科进行预测.

• 数据拆分

- ▶ 各学科前 60% 的数据作为训练集,后 20% 的数据作为测试集
- ▶ 我的思路是
 - 遍历每一个 research_field.
 - 在同学科内, 按 rank 升序排序.
 - 使用 iloc 按比例切片.
 - 将所有学科的切片重新组合成最终的 X train, y train, X test, y test

• 预处理

- ▶ 我们需要一个 Pipeline 来处理混合数据.
- ▶ 数值特征 (documents, cites 等) 可以直接使用
- ▶ 类别特征必须进行独热编码,将其转换为模型可以理解的 0/1 矩阵.
- ► ColumnTransformer 是实现这一点的最佳工具.

• 模型选型

- ▶ 算法选用随机森林回归.
- ▶ 原因:
 - 鲁棒性: 对数据尺度不敏感, 不需要复杂的特征缩放.
 - 非线性: 它能很好地捕捉 cites 和 rank 之间的非线性关系(例如,引用数加倍,排名可能提升远不止两倍). 交互性: 能自动捕捉特征间的交互作用(例如 cites 对 rank 的影响,在"化学"和"数学"领域是不同的).
 - 可解释性: 可以提供"特征重要性", 让我们知道哪个指标对 ESI 排名的预测贡献最大.



• 评估

- R^2 : 模型解释了多少百分比的方差? (越接近 1 越好).
- ► RMSE: 模型的预测平均偏离了多少个"名次"? (越低越好).

In [14]:

```
python
2 # 单元格 7: 任务3 - 学科排名预测模型 (回归)
3 # -----
4 print("--- 开始执行任务3:学科排名预测模型 ---")
6 # 1. 分析 (Analysis)
7 print("\n[分析]: 这是一个监督学习中的回归 (Regression) 问题.")
9 # 2. 思路 (Approach)
10 print("[思路]:")
11 print(" - 目标 (y): 'rank'. ")
12 print(" - 特征 (X): 'documents', 'cites', 'cites_per_paper', 'top_papers' (数值) +
   'research field' (类别). ")
13 print(" - 拆分: 按学科内排名 60% 训练,后 20% 测试.")
14 print(" - 预处理: 使用 ColumnTransformer 和 Pipeline 对类别特征 OneHotEncode,数值特征
   Passthrough. ")
15 print(" - 模型:使用 RandomForestRegressor,它对非线性和特征交互处理较好.")
16 print(" - 评估: R<sup>2</sup> 和 RMSE.")
18 # 3. 代码实现 (Code)
19
20 # 3.1. 准备数据
21 df_model = df_raw.dropna(subset=[
      'research_field', 'rank', 'documents',
      'cites', 'cites_per_paper', 'top_papers'
23
24 ])
25
26 target = 'rank'
27 numeric_features = ['documents', 'cites', 'cites_per_paper', 'top_papers']
28 categorical_features = ['research_field']
29
30 X = df_model[numeric_features + categorical_features]
31 y = df_model[target]
32 print("步骤 1/7: 准备特征 (X) 和目标 (y) 完成.")
33
34 # 3.2. 自定义数据拆分
35 print("步骤 2/7: 正在按学科内排名拆分 60% 训练集和 20% 测试集...")
36 train_indices = []
37 test_indices = []
38
39 for field in df_model['research_field'].unique():
      df_field = df_model[df_model['research_field'] == field].copy()
40
      df_field_sorted = df_field.sort_values(by='rank', ascending=True)
41
42
43
      n = len(df_field_sorted)
44
      train end = int(n * 0.6)
45
      test_start = int(n * 0.8)
46
47
      train indices.extend(df field sorted.iloc[:train end].index)
48
      test indices.extend(df field sorted.iloc[test start:].index)
49
```



```
50 X_train = X.loc[train_indices]
51 y_train = y.loc[train_indices]
52 X test = X.loc[test indices]
53 y_test = y.loc[test_indices]
54 print(f"数据拆分完成. 训练集大小: {len(X train)}, 测试集大小: {len(X test)}")
55
56 # 3.3. 创建预处理 Pipeline
57 print("步骤 3/7: 正在构建数据预处理 Pipeline...")
58 preprocessor = ColumnTransformer(
59
       transformers=[
           ('num', 'passthrough', numeric_features),
60
61
           ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
62
        ])
63
64 # 3.4. 定义模型
65 model = RandomForestRegressor(
       n estimators=100,
66
67
       random state=42,
68
       n jobs=-1,
69
       max_depth=20,
70
       min_samples_leaf=5
71 )
72 print("步骤 4/7: 定义 RandomForestRegressor 模型.")
74 # 3.5. 组装完整 Pipeline
75 clf = Pipeline(steps=[('preprocessor', preprocessor),
76
                         ('regressor', model)])
77
78 # 3.6. 模型训练
79 print("步骤 5/7: 正在训练模型...")
80 clf.fit(X_train, y_train)
81 print("模型训练完成.")
83 # 3.7. 模型评估
84 print("步骤 6/7: 正在评估模型 (测试集)...")
85 y_pred = clf.predict(X_test)
87 rmse = np.sqrt(mean_squared_error(y_test, y_pred))
88 r2 = r2\_score(y\_test, y\_pred)
89
90 print("\n--- 任务3 结果:模型评估(测试集) ---")
91 print(f"R2 (决定系数): {r2:.4f}")
92 print(f"RMSE (均方根误差): {rmse:.4f} (预测排名平均误差约 {rmse:.0f} 名)")
93
94 # 3.8. 特征重要性分析
95 print("\n步骤 7/7: 正在分析特征重要性...")
96 try:
97
        rf_model = clf.named_steps['regressor']
98
        ohe feature names =
        clf.named_steps['preprocessor'].named_transformers_['cat'].get_feature_names_out(cat')
99
        all_feature_names = numeric_features + list(ohe_feature_names)
100
101
        importances = rf_model.feature_importances_
102
103
        feature_importance_df = pd.DataFrame({
104
            'feature': all_feature_names,
            'importance': importances
105
```



```
106
       }).sort_values(by='importance', ascending=False)
107
108
       print("\n--- 任务3 结果:模型特征重要性 (Top 10) ---")
109
       print(feature importance df.head(10))
110
111 except Exception as e imp:
112
       print(f"分析特征重要性时出错: {e_imp}")
113
114 print("\n--- 任务3 完成 ---")
                                                                            txt
 1 --- 开始执行任务3:学科排名预测模型 ---
 3 [分析]: 这是一个监督学习中的回归 (Regression) 问题.
 4 [思路]:
 5
    - 目标 (y): 'rank'.
 6
     - 特征 (X): 'documents', 'cites', 'cites_per_paper', 'top_papers' (数值) +
    'research_field' (类别).
 7
    - 拆分: 按学科内排名 60% 训练, 后 20% 测试.
 8
     - 预处理: 使用 ColumnTransformer 和 Pipeline 对类别特征 OneHotEncode,数值特征
     Passthrough.
 9
    - 模型: 使用 RandomForestRegressor,它对非线性和特征交互处理较好.
10
   - 评估: R² 和 RMSE.
11 步骤 1/7: 准备特征 (X) 和目标 (y) 完成.
12 步骤 2/7: 正在按学科内排名拆分 60% 训练集和 20% 测试集...
13 数据拆分完成. 训练集大小: 20464, 测试集大小: 6833
14 步骤 3/7: 正在构建数据预处理 Pipeline...
15 步骤 4/7: 定义 RandomForestRegressor 模型.
16 步骤 5/7: 正在训练模型...
17 模型训练完成.
18 步骤 6/7: 正在评估模型 (测试集)...
19
20 --- 任务3 结果:模型评估 (测试集) ---
21 R2 (决定系数): 0.6939
22 RMSE (均方根误差): 1055.9679 (预测排名平均误差约 1056 名)
23
24 步骤 7/7: 正在分析特征重要性...
25
26 --- 任务3 结果:模型特征重要性 (Top 10) ---
27
                                feature importance
28 7
           research field CLINICAL MEDICINE
                                          0.496979
29 1
                                          0.444441
                                  cites
30 10
                research field ENGINEERING
                                          0.017065
31 6
                 research field CHEMISTRY
                                          0.009471
        research_field_ENVIRONMENT/ECOLOGY
32 11
                                          0.006295
         research_field_MULTIDISCIPLINARY
                                          0.004359
33 18
34 3
                              top_papers
                                          0.003504
35 14
           research_field_MATERIALS SCIENCE
                                          0.002655
36 5 research_field_BIOLOGY & BIOCHEMISTRY
                                          0.002603
37 15
                research field MATHEMATICS
                                          0.002272
38
39 --- 任务3 完成 ---
40
```





完整代码见 ./hw5.ipynb.