

数据科学导论 -HW 4 报告

孙育泉 10234900421
2025.10.14

总览

- I 实验要求 1
- II 具体实现 1
 - II.1 建立数据库与导入数据 1
 - II.2 选择合适的 Schema 4
 - II.3 获取 ECNU 在各个学科中的排名 5
 - II.4 获取中国（大陆地区）大学在各个学科中的表现 8
 - II.5 分析全球不同区域在各个学科中的表现 11

I 实验要求

对于 ESI 数据的获取方面，可以查看 hw3 目录下的文件，不再赘述. 这次我们已经获取到了 csv 文件，要求如下：

- (1) 将获取的数据导入到一个关系型数据库系统中（系统可以自选）
- (2) 优化关系型数据，并整理一个合理的 schema
- (3) 通过写 SQL 语句，获取华东师范大学在各个学科中的排名
- (4) 通过写 SQL 语句，获取中国（大陆地区）大学在各个学科中的表现
- (5) 通过写 SQL 语句，分析全球不同区域在各个学科中的表现

II 具体实现

II.1 建立数据库与导入数据

由于本次实验的任务较轻，且便于 github 的展示，这里选用轻量级的 SQLite 作为数据库系统.

首先，引入所需的库文件，检查数据库是否已经存在，如果存在则删除，防止重复插入. 我们先要创建数据库的表结构：

sql

```
1 CREATE TABLE IF NOT EXISTS esi_rankings (  
2     id INTEGER PRIMARY KEY AUTOINCREMENT,  
3     research_field TEXT NOT NULL,  
4     rank INTEGER,  
5     institution TEXT NOT NULL,  
6     country_region TEXT,  
7     documents INTEGER,  
8     cites INTEGER,  
9     cites_per_paper REAL,  
10    top_papers INTEGER  
11 );
```

接着，我们将 csv 文件中的数据条目插入到数据库中. 在这里使用 pandas 读取 csv 文件，然后使用 to_sql 方法将数据插入到数据库中.

这里要注意的是，所给文件的编码是 `latin-1`，并非 `utf-8`，所以在读取时要指定编码格式。以及有些数据可能会有空缺，要进行处理。

In [17]:

python

```
1 import sqlite3
2 import pandas as pd
3 import os
4 import glob
5
6 csv_folder_path = "./data" # 包含所有22个CSV文件的文件夹路径
7 db_path = "./esi_rankings.db" # 希望生成的数据库文件的存放路径和名称
8
9 if os.path.exists(db_path):
10     os.remove(db_path)
11     print(f"已删除旧的数据库文件: {db_path}")
12
13 conn = sqlite3.connect(db_path)
14 cursor = conn.cursor()
15
16 create_table_sql = """
17 CREATE TABLE IF NOT EXISTS esi_rankings (
18     id INTEGER PRIMARY KEY AUTOINCREMENT,
19     research_field TEXT NOT NULL,
20     rank INTEGER,
21     institution TEXT NOT NULL,
22     country_region TEXT,
23     documents INTEGER,
24     cites INTEGER,
25     cites_per_paper REAL,
26     top_papers INTEGER
27 );
28 """
29 cursor.execute(create_table_sql)
30 print(f"数据库 '{db_path}' 和数据表 'esi_rankings' 已成功创建。")
31
32 csv_files = glob.glob(os.path.join(csv_folder_path, "*.csv"))
33
34 if not csv_files:
35     print(f"警告: 在文件夹 '{csv_folder_path}' 中没有找到任何CSV文件。请检查路径。")
36 else:
37     for file_path in csv_files:
38         file_encoding = "latin-1"
39
40         with open(file_path, "r", encoding=file_encoding) as f:
41             first_line = f.readline()
42
43         try:
44             field_value = (
45                 first_line.split("Filter Value(s):")[1].split("Show:")[0].strip()
46             )
47             print(f"正在处理文件: {os.path.basename(file_path)}, 学科: {field_value}")
48
49             df = pd.read_csv(file_path, header=1, encoding=file_encoding)
50
51             # 重命名列以匹配数据库的Schema
52             df.rename(
53                 columns={
54                     "Unnamed: 0": "rank",
```

```

55         "Institutions": "institution",
56         "Countries/Regions": "country_region",
57         "Web of Science Documents": "documents",
58         "Cites": "cites",
59         "Cites/Paper": "cites_per_paper",
60         "Top Papers": "top_papers",
61     },
62     inplace=True,
63 )
64
65 # 关键清洗步骤：移除机构名称为空的无效行（例如CSV末尾的空行）
66 initial_rows = len(df)
67 df.dropna(subset=["institution"], inplace=True)
68
69 # --- d. 将清洗后的数据写入数据库 ---
70 if not df.empty:
71     df["research_field"] = field_value
72     # 筛选并排序 DataFrame 的列，以确保与数据库表结构一致
73     final_df = df[
74         [
75             "research_field",
76             "rank",
77             "institution",
78             "country_region",
79             "documents",
80             "cites",
81             "cites_per_paper",
82             "top_papers",
83         ]
84     ]
85     # 使用 'append' 模式将数据追加到表中
86     final_df.to_sql("esi_rankings", conn, if_exists="append", index=False)
87     print(f"✅ 成功导入 {len(final_df)} 条 '{field_value}' 学科的数据。")
88
89 except Exception as e:
90     print(f"❌ 处理文件 {os.path.basename(file_path)} 时发生错误: {e}")
91
92 # --- 4. 提交更改并关闭数据库连接 ---
93 conn.commit()
94 conn.close()
95 print("\n🎉 所有数据导入完成！数据库 'esi_rankings.db' 已准备就绪。")

```

txt

- 1 已删除旧的数据库文件: ./esi_rankings.db
- 2 数据库 './esi_rankings.db' 和数据表 'esi_rankings' 已成功创建.
- 3 正在处理文件: AGRICULTURAL SCIENCES.csv, 学科: AGRICULTURAL SCIENCES
- 4 ✅ 成功导入 1381 条 'AGRICULTURAL SCIENCES' 学科的数据.
- 5 正在处理文件: BIOLOGY & BIOCHEMISTRY.csv, 学科: BIOLOGY & BIOCHEMISTRY
- 6 ✅ 成功导入 1649 条 'BIOLOGY & BIOCHEMISTRY' 学科的数据.
- 7 正在处理文件: CHEMISTRY.csv, 学科: CHEMISTRY
- 8 ✅ 成功导入 2141 条 'CHEMISTRY' 学科的数据.
- 9 正在处理文件: CLINICAL MEDICINE.csv, 学科: CLINICAL MEDICINE
- 10 ✅ 成功导入 6754 条 'CLINICAL MEDICINE' 学科的数据.
- 11 正在处理文件: COMPUTER SCIENCE.csv, 学科: COMPUTER SCIENCE
- 12 ✅ 成功导入 863 条 'COMPUTER SCIENCE' 学科的数据.
- 13 正在处理文件: ECONOMICS & BUSINESS.csv, 学科: ECONOMICS & BUSINESS
- 14 ✅ 成功导入 543 条 'ECONOMICS & BUSINESS' 学科的数据.

```

15 正在处理文件: ENGINEERING.csv, 学科: ENGINEERING
16 成功导入 2787 条 'ENGINEERING' 学科的数据.
17 正在处理文件: ENVIRONMENT ECOLOGY.csv, 学科: ENVIRONMENT/ECOLOGY
18 成功导入 2066 条 'ENVIRONMENT/ECOLOGY' 学科的数据.
19 正在处理文件: GEOSCIENCES.csv, 学科: GEOSCIENCES
20 成功导入 1175 条 'GEOSCIENCES' 学科的数据.
21 正在处理文件: IMMUNOLOGY.csv, 学科: IMMUNOLOGY
22 成功导入 1177 条 'IMMUNOLOGY' 学科的数据.
23 正在处理文件: MATERIALS SCIENCE.csv, 学科: MATERIALS SCIENCE
24 成功导入 1580 条 'MATERIALS SCIENCE' 学科的数据.
25 正在处理文件: MATHEMATICS.csv, 学科: MATHEMATICS
26 成功导入 395 条 'MATHEMATICS' 学科的数据.
27 正在处理文件: MICROBIOLOGY.csv, 学科: MICROBIOLOGY
28 成功导入 803 条 'MICROBIOLOGY' 学科的数据.
29 正在处理文件: MOLECULAR BIOLOGY & GENETICS.csv, 学科: MOLECULAR BIOLOGY & GENETICS
30 成功导入 1169 条 'MOLECULAR BIOLOGY & GENETICS' 学科的数据.
31 正在处理文件: MULTIDISCIPLINARY.csv, 学科: MULTIDISCIPLINARY
32 成功导入 216 条 'MULTIDISCIPLINARY' 学科的数据.
33 正在处理文件: NEUROSCIENCE & BEHAVIOR.csv, 学科: NEUROSCIENCE & BEHAVIOR
34 成功导入 1298 条 'NEUROSCIENCE & BEHAVIOR' 学科的数据.
35 正在处理文件: PHARMACOLOGY & TOXICOLOGY.csv, 学科: PHARMACOLOGY & TOXICOLOGY
36 成功导入 1389 条 'PHARMACOLOGY & TOXICOLOGY' 学科的数据.
37 正在处理文件: PHYSICS.csv, 学科: PHYSICS
38 成功导入 995 条 'PHYSICS' 学科的数据.
39 正在处理文件: PLANT & ANIMAL SCIENCE.csv, 学科: PLANT & ANIMAL SCIENCE
40 成功导入 1950 条 'PLANT & ANIMAL SCIENCE' 学科的数据.
41 正在处理文件: PSYCHIATRY PSYCHOLOGY.csv, 学科: PSYCHIATRY/PSYCHOLOGY
42 成功导入 1147 条 'PSYCHIATRY/PSYCHOLOGY' 学科的数据.
43 正在处理文件: SOCIAL SCIENCES, GENERAL.csv, 学科: SOCIAL SCIENCES, GENERAL
44 成功导入 2407 条 'SOCIAL SCIENCES, GENERAL' 学科的数据.
45 正在处理文件: SPACE SCIENCE.csv, 学科: SPACE SCIENCE
46 成功导入 236 条 'SPACE SCIENCE' 学科的数据.
47
48 所有数据导入完成! 数据库 'esi_rankings.db' 已准备就绪.
49

```

至此，创建数据库工作完成。

II.2 选择合适的 Schema

根据这个项目要求和数据特点，一个“好”的结构应该是：

- (1) 逻辑清晰：能够直观地反映数据。
- (2) 查询高效：能快速地执行分析查询。
- (3) 避免冗余：在合理的范围内减少重复数据。
- (4) 保证数据完整性：确保存入的数据是有效和一致的。

先来看一下现在的 schema：

列名 (Column Name)	数据类型 (Data Type)	约束/说明 (Constraints/Notes)
id	INTEGER	主键 (<i>PRIMARY KEY</i>), AUTOINCREMENT. 每一行的唯一标识，数据库自动生成。

列名 (Column Name)	数据类型 (Data Type)	约束/说明 (Constraints/Notes)
research_field	TEXT	不能为空 (<i>NOT NULL</i>). 关键的分类字段, 如 “CHEMISTRY”, “ENGINEERING”.
rank	INTEGER	机构在该学科下的排名.
institution	TEXT	不能为空 (<i>NOT NULL</i>). 机构名称.
country_region	TEXT	国家或地区.
documents	INTEGER	Web of Science 文献数.
cites	INTEGER	总引用数.
cites_per_paper	REAL	篇均引用数. 使用 REAL (浮点数) 来精确存储小数.
top_papers	INTEGER	顶尖论文数.

事实上, 这样的安排已经很好了. 虽然, 我们后续要查询“机构”或者“国家”, 但是这些字段并不适合单独拆分成表格, 因为它们的种类太多, 且没有明显的层级关系. 拆分成表格反而会会增加查询的复杂度, 在这张表中, 我们通过简单的 `WHERE` 和 `GROUP BY` 就可以完成. 如果分表, 就需要使用 `JOIN` 操作, 这会让查询变得更复杂.

II.3 获取 ECNU 在各个学科中的排名

为了方便后续的代码编写, 我们先定义如下几个东西:

- `execute_query()`: 执行 SQL 查询的函数
- `get_iso_alpha_3()`: 获取国家的 ISO Alpha-3 代码的函数, 用于最后一步实现全球可视化
- `DATABASE_TITLE`: 数据库的路径

In [18]:

```

1 import sqlite3
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 import plotly.graph_objects as go
7 import pycountry
8
9
10 def execute_query(db_path, query):
11     """
12     功能: 连接指定的SQLite数据库, 执行一条SQL查询语句.
13     参数:
14         db_path (str): 数据库文件的路径.
15         query (str): 需要执行的SQL查询字符串.
16     返回:
17         pandas.DataFrame: 包含查询结果的DataFrame. 如果出错则返回None.
18     """
19     try:
20         with sqlite3.connect(db_path) as conn:
21             df = pd.read_sql_query(query, conn)
22             return df
23     except Exception as e:
24         print(f"执行查询时发生错误: {e}")
25         return None

```

python

```

26
27
28 def get_iso_alpha_3(country_name):
29     """
30     功能：将国家/地区名称转换为标准的ISO 3166-1 alpha-3三字符代码。
31     说明：内置了一个手动映射表来处理数据库中的非标准名称。
32     参数：
33         country_name (str): 数据库中的国家/地区名称。
34     返回：
35         str: 对应的三字符ISO代码，如果找不到则返回None。
36     """
37     # 手动映射表，处理 'pycountry' 无法直接识别的特殊名称
38     manual_map = {
39         "CHINA MAINLAND": "CHN",
40         "USA": "USA",
41         "ENGLAND": "GBR", # 英格兰、苏格兰等均属于联合王国 (GBR)
42         "SCOTLAND": "GBR",
43         "WALES": "GBR",
44         "NORTH IRELAND": "GBR",
45         "GERMANY (FED REP GER)": "DEU",
46         "SOUTH KOREA": "KOR",
47         "RUSSIA": "RUS",
48         "IRAN": "IRN",
49     }
50
51     if country_name in manual_map:
52         return manual_map[country_name]
53
54     try:
55         # 尝试使用pycountry的模糊搜索功能来查找匹配项
56         return pycountry.countries.search_fuzzy(country_name)[0].alpha_3
57     except (LookupError, AttributeError):
58         # 如果找不到任何匹配，则返回None
59         return None
60
61
62 # --- 全局设置 ---
63 # 定义数据库文件路径，后续所有单元格都将使用此变量
64 DATABASE_FILE = "esi_rankings.db"
65 print("✅ 环境准备就绪，核心函数已定义。")

```

txt

```

1 ✅ 环境准备就绪，核心函数已定义。
2

```

这里的处理是，为了方便后续的对比，我们定义了一个所查机构的列表，用于遍历。只需要在查询到 ECNU 的时候进行输出即可。这里的查询语句就是看机构的名字是否与 `institution` 字段匹配。注意，我们使用了 `LIKE` 语句，而不是 `=`，为了提高我们查询的容错。

In [19]:

```

1 universities_to_compare = [
2     "EAST CHINA NORMAL UNIV",
3     "PEKING UNIV",
4     "TSINGHUA UNIV",
5     "FUDAN UNIV",
6     "TONGJI UNIV",

```

python

```

7 | "SHANGHAI JIAO TONG UNIV",
8 | ]
9 |
10 | all_rankings_data = []
11 | print("正在从数据库中查询指定大学的排名数据...")
12 |
13 | for university in universities_to_compare:
14 |     query = f"SELECT research_field, rank FROM esi_rankings WHERE institution LIKE
        |     '{university}%"
15 |     df = execute_query(DATABASE_FILE, query)
16 |     if df is not None and not df.empty:
17 |         df["institution"] = university # 添加大学名称列用于分组
18 |         if university == universities_to_compare[0]:
19 |             print(df)
20 |         all_rankings_data.append(df)

```

txt

```

1 | 正在从数据库中查询指定大学的排名数据...
2 |
3 |      |      |      |      |      |      |      |      |      |      |
4 | 0 |      |      |      |      |      |      |      |      |      |      |
5 | 1 |      |      |      |      |      |      |      |      |      |      |
6 | 2 |      |      |      |      |      |      |      |      |      |      |
7 | 3 |      |      |      |      |      |      |      |      |      |      |
8 | 4 |      |      |      |      |      |      |      |      |      |      |
9 | 5 |      |      |      |      |      |      |      |      |      |      |
10 | 6 |      |      |      |      |      |      |      |      |      |      |
11 | 7 |      |      |      |      |      |      |      |      |      |      |
12 | 8 |      |      |      |      |      |      |      |      |      |      |
13 | 9 |      |      |      |      |      |      |      |      |      |      |
14 | 10 |      |      |      |      |      |      |      |      |      |      |
15 | 11 |      |      |      |      |      |      |      |      |      |      |
16 | 12 |      |      |      |      |      |      |      |      |      |      |
17 | 13 |      |      |      |      |      |      |      |      |      |      |
18 | 14 |      |      |      |      |      |      |      |      |      |      |
19 | 15 |      |      |      |      |      |      |      |      |      |      |
20 | 16 |      |      |      |      |      |      |      |      |      |      |

```

	research_field	rank	institution
0	AGRICULTURAL SCIENCES	845	EAST CHINA NORMAL UNIV
1	BIOLOGY & BIOCHEMISTRY	721	EAST CHINA NORMAL UNIV
2	CHEMISTRY	90	EAST CHINA NORMAL UNIV
3	CLINICAL MEDICINE	2852	EAST CHINA NORMAL UNIV
4	COMPUTER SCIENCE	207	EAST CHINA NORMAL UNIV
5	ENGINEERING	317	EAST CHINA NORMAL UNIV
6	ENVIRONMENT/ECOLOGY	130	EAST CHINA NORMAL UNIV
7	GEOSCIENCES	275	EAST CHINA NORMAL UNIV
8	MATERIALS SCIENCE	196	EAST CHINA NORMAL UNIV
9	MATHEMATICS	115	EAST CHINA NORMAL UNIV
10	MOLECULAR BIOLOGY & GENETICS	867	EAST CHINA NORMAL UNIV
11	NEUROSCIENCE & BEHAVIOR	853	EAST CHINA NORMAL UNIV
12	PHARMACOLOGY & TOXICOLOGY	1064	EAST CHINA NORMAL UNIV
13	PHYSICS	522	EAST CHINA NORMAL UNIV
14	PLANT & ANIMAL SCIENCE	395	EAST CHINA NORMAL UNIV
15	PSYCHIATRY/PSYCHOLOGY	467	EAST CHINA NORMAL UNIV
16	SOCIAL SCIENCES, GENERAL	314	EAST CHINA NORMAL UNIV

为了方便对比，这里罗列了一些高校在所有学科中的排名。（注意：越低越好 hhh）

In [20]:

python

```

1 | if all_rankings_data:
2 |     # 将所有查询结果合并成一个DataFrame
3 |     comparison_df = pd.concat(all_rankings_data, ignore_index=True)
4 |     print("数据获取完毕，正在生成可视化图表...")
5 |
6 |     # 使用Seaborn绘制分组条形图
7 |     plt.figure(figsize=(20, 12))
8 |     sns.set_theme(style="whitegrid", font="Maple Mono Normal NF CN") # 设置支持中文的
        |     字体，如黑体
9 |     plt.rcParams["axes.unicode_minus"] = False # 解决负号显示问题
10 |
11 |     sns.barplot(
12 |         data=comparison_df,
13 |         x="research_field",
14 |         y="rank",
15 |         hue="institution",

```

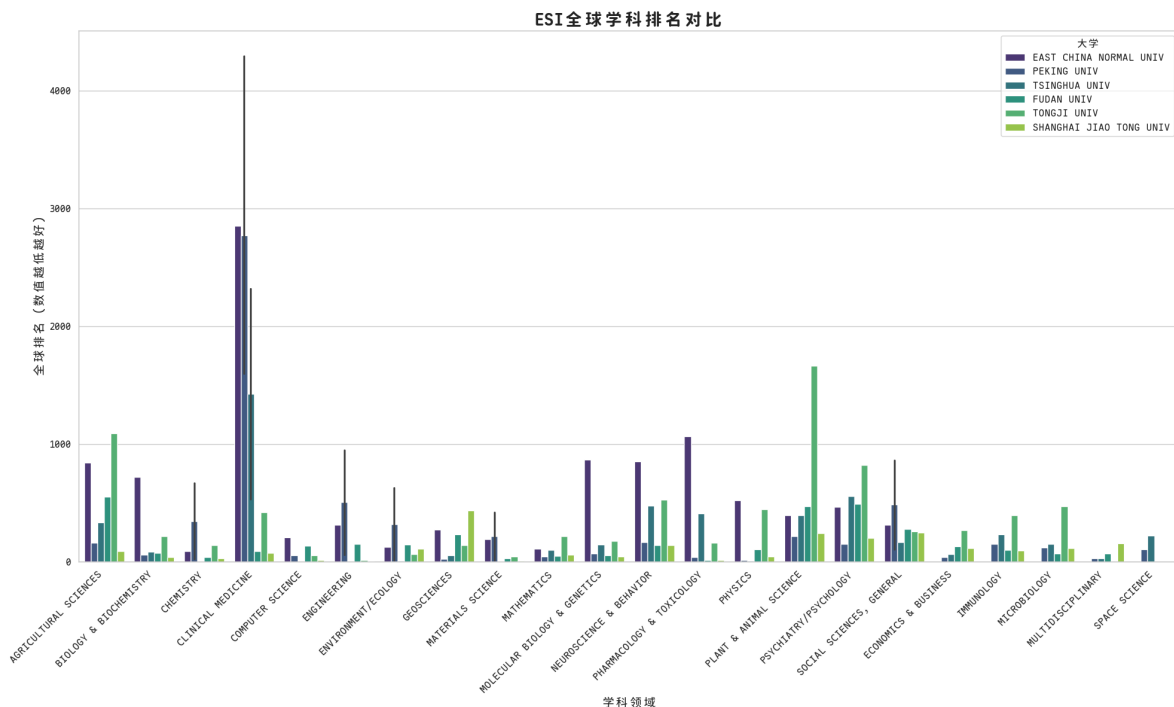
```

16     palette="viridis",
17   )
18
19   # 美化图表
20   plt.title("ESI全球学科排名对比", fontsize=20, weight="bold")
21   plt.xlabel("学科领域", fontsize=14)
22   plt.ylabel("全球排名 (数值越低越好)", fontsize=14)
23   plt.xticks(rotation=45, ha="right", fontsize=12)
24   plt.legend(title="大学", fontsize=12)
25   plt.tight_layout()
26   plt.show()
27 else:
28   print("未能查询到任何指定大学的数据。")
  
```

txt

```

1 数据获取完毕，正在生成可视化图表...
2
  
```



II.4 获取中国（大陆地区）大学在各个学科中的表现

这里，我们使用了 `country_region` 字段来筛选出中国（大陆地区）的大学。然后，我们考虑两种指标：

- (1) 在每个学科中，有多少所中国大陆的大学进入了全球排名前 100 名
- (2) 在每个学科中，排名第一的中国大陆大学

既注重整体，也注重代表。

In [21]:

```

1 print("--- 分析中国（大陆）大学在各学科的表现 ---")
2
3 # --- 分析角度1：集体实力 ---
  
```

python



```
4 # 统计在每个学科中，有多少所中国大陆的大学进入了全球排名前100名。
5 # 这反映了中国在哪些学科领域具有广泛的高水平研究基础。
6 print("\n>>> 分析1: 各学科领域进入全球前100的中国（大陆）大学数量")
7 query_strength = """
8     SELECT
9         research_field AS '学科领域',
10        COUNT(institution) AS '全球前100大学数量'
11     FROM
12         esi_rankings
13     WHERE
14         country_region = 'CHINA MAINLAND' AND rank <= 100
15     GROUP BY
16         research_field
17     ORDER BY
18         '全球前100大学数量' DESC;
19 """
20 df_strength = execute_query(DATABASE_FILE, query_strength)
21 if df_strength is not None:
22     # 使用 display() 函数在Jupyter中可以获得更好的表格格式化效果
23     display(df_strength)
24 else:
25     print("未能查询到相关数据。")
26
27
28 # --- 分析角度2: 顶尖实力 ---
29 # 找出在每个学科中，排名第一的中国大陆大学。
30 # 这反映了中国在各个学科领域的“领头羊”。
31 print("\n>>> 分析2: 各学科领域排名第一的中国（大陆）大学")
32 query_leaders = """
33     SELECT
34         research_field AS '学科领域',
35         institution AS '排名第一的大学',
36         rank AS '全球排名'
37     FROM
38         (
39             SELECT
40                 research_field,
41                 institution,
42                 rank,
43                 ROW_NUMBER() OVER(PARTITION BY research_field ORDER BY rank ASC)
44             as rn
45             FROM
46                 esi_rankings
47             WHERE
48                 country_region = 'CHINA MAINLAND'
49         )
50     WHERE
51         rn = 1
52     ORDER BY
53         research_field ASC;
54 """
55 df_leaders = execute_query(DATABASE_FILE, query_leaders)
56 if df_leaders is not None:
57     display(df_leaders)
58 else:
59     print("未能查询到相关数据。")
```

txt

```
1 --- 分析中国（大陆）大学在各学科的表现 ---
2
3 >>> 分析1: 各学科领域进入全球前100的中国（大陆）大学数量
4
```

txt

1		学科领域	全球前100大学数量
2	0	SPACE SCIENCE	2
3	1	SOCIAL SCIENCES, GENERAL	1
4	2	PLANT & ANIMAL SCIENCE	13
5	3	PHYSICS	10
6	4	PHARMACOLOGY & TOXICOLOGY	29
7	5	NEUROSCIENCE & BEHAVIOR	2
8	6	MULTIDISCIPLINARY	5
9	7	MOLECULAR BIOLOGY & GENETICS	9
10	8	MICROBIOLOGY	11
11	9	MATHEMATICS	24
12	10	MATERIALS SCIENCE	48
13	11	IMMUNOLOGY	4
14	12	GEOSCIENCES	22
15	13	ENVIRONMENT/ECOLOGY	14
16	14	ENGINEERING	46
17	15	ECONOMICS & BUSINESS	8
18	16	COMPUTER SCIENCE	44
19	17	CLINICAL MEDICINE	3
20	18	CHEMISTRY	46
21	19	BIOLOGY & BIOCHEMISTRY	8
22	20	AGRICULTURAL SCIENCES	22

txt

```
1
2 >>> 分析2: 各学科领域排名第一的中国（大陆）大学
3
```

txt

1		学科领域	排名第一的大学	全球排名
2	0	AGRICULTURAL SCIENCES	CHINESE ACADEMY OF SCIENCES	1
3	1	BIOLOGY & BIOCHEMISTRY	CHINESE ACADEMY OF SCIENCES	3
4	2	CHEMISTRY	CHINESE ACADEMY OF SCIENCES	1
5	3	CLINICAL MEDICINE	SHANGHAI JIAO TONG UNIVERSITY	75
6	4	COMPUTER SCIENCE	CHINESE ACADEMY OF SCIENCES	1
7	5	ECONOMICS & BUSINESS	PEKING UNIVERSITY	44
8	6	ENGINEERING	CHINESE ACADEMY OF SCIENCES	1
9	7	ENVIRONMENT/ECOLOGY	CHINESE ACADEMY OF SCIENCES	1
10	8	GEOSCIENCES	CHINESE ACADEMY OF SCIENCES	1
11	9	IMMUNOLOGY	CHINESE ACADEMY OF SCIENCES	36
12	10	MATERIALS SCIENCE	CHINESE ACADEMY OF SCIENCES	1
13	11	MATHEMATICS	CHINESE ACADEMY OF SCIENCES	4
14	12	MICROBIOLOGY	CHINESE ACADEMY OF SCIENCES	3
15	13	MOLECULAR BIOLOGY & GENETICS	CHINESE ACADEMY OF SCIENCES	13
16	14	MULTIDISCIPLINARY	CHINESE ACADEMY OF SCIENCES	3
17	15	NEUROSCIENCE & BEHAVIOR	CHINESE ACADEMY OF SCIENCES	71
18	16	PHARMACOLOGY & TOXICOLOGY	CHINESE ACADEMY OF SCIENCES	4
19	17	PHYSICS	CHINESE ACADEMY OF SCIENCES	1
20	18	PLANT & ANIMAL SCIENCE	CHINESE ACADEMY OF SCIENCES	1
21	19	PSYCHIATRY/PSYCHOLOGY	PEKING UNIVERSITY	154
22	20	SOCIAL SCIENCES, GENERAL	CHINESE ACADEMY OF SCIENCES	47
23	21	SPACE SCIENCE	CHINESE ACADEMY OF SCIENCES	15

II.5 分析全球不同区域在各个学科中的表现

这里，为了以地区为单位，我们在查询的时候，使用了 `country_region` 字段进行分组。我们统计了每个国家在各个学科中的总引用数并进行了可视化。

python

```
1 print("正在获取全球各学科数据用于生成世界地图...")
2 subject_data_query = """
3     SELECT
4         country_region,
5         research_field,
6         SUM(cites) AS total_cites
7     FROM
8         esi_rankings
9     GROUP BY
10        country_region, research_field;
11 """
12 df_map = execute_query(DATABASE_FILE, subject_data_query)
13
14 if df_map is not None:
15     print("数据预处理中：清洗、转换ISO代码...")
16     df_map.dropna(subset=["country_region", "research_field"], inplace=True)
17     df_map["total_cites"] = pd.to_numeric(
18         df_map["total_cites"], errors="coerce"
19     ).fillna(0)
20     df_map["iso_alpha"] = df_map["country_region"].apply(
21         get_iso_alpha_3
22     ) # 调用单元格1的函数
23
24     subjects = sorted(df_map["research_field"].unique())
25
26     print("正在生成交互式地图...")
27     fig = go.Figure()
28
29     for subject in subjects:
30         subject_df = df_map[df_map["research_field"] == subject]
31         log_cites = np.log10(
32             subject_df["total_cites"] + 1
33         ) # 对引用数取对数，以优化颜色显示
34
35         fig.add_trace(
36             go.Choropleth(
37                 locations=subject_df["iso_alpha"],
38                 z=log_cites,
39                 text=subject_df["country_region"],
40                 customdata=subject_df["total_cites"],
41                 hovertemplate="<b>%(text)</b><br>总引用数: %(customdata:,.0f)<extra></extra>",
42                 colorscale="Plasma",
43                 visible=(subject == subjects[0]), # 默认只显示第一个学科
44             )
45         )
46
47     buttons = []
48     for i, subject in enumerate(subjects):
49         visibility_mask = [False] * len(subjects)
50         visibility_mask[i] = True
51         buttons.append(
52             dict(label=subject, method="update", args=[{"visible": visibility_mask}])
53         )
```

```
54
55 fig.update_layout(
56     title_text=f"<b>全球各学科学术影响力地图</b><br><i>请从下拉菜单中选择一个学科</i>",
57     title_x=0.5,
58     title_font_size=22,
59     margin=dict(t=70, l=0, r=0, b=0),
60     updatemenus=[
61         dict(
62             active=0,
63             buttons=buttons,
64             direction="down",
65             pad={"r": 10, "t": 10},
66             showactive=True,
67             x=0.1, # 将菜单的x坐标设置在左侧10%的位置
68             xanchor="left", # x坐标的锚点是菜单的左边缘
69             y=1.1, # 将菜单的y坐标放在新的顶部空间内
70             yanchor="top", # y坐标的锚点是菜单的上边缘
71         )
72     ],
73     geo=dict(
74         showframe=False, showcoastlines=False, projection_type="natural earth"
75     ),
76 )
77
78 fig.show()
79 else:
80     print("未能获取用于生成地图的数据。")
```

这里只截取了 CS 学科的结果截图：

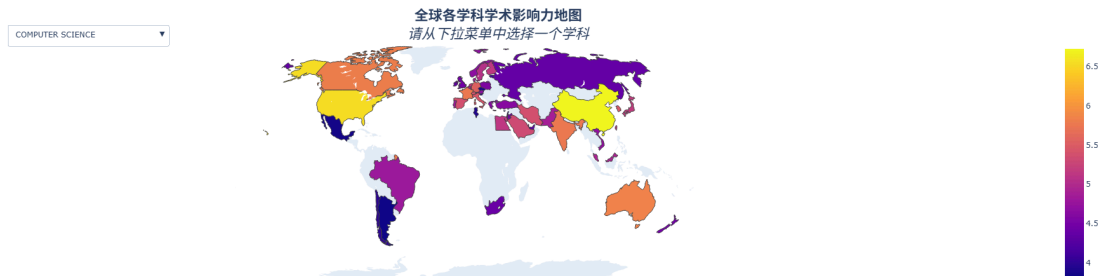


Figure 1: 全球机构 - computer science

Remark

- 完整代码见 `hw4.ipynb`.
- 最后一部分为交互图表，在 pdf 文件中无法展示，可以打开 jupyter notebook 查看.