

数据科学导论 -HW 6 报告

孙育泉 10234900421
2025.10.29

总览

- I 实验要求 1
- II 具体实现 1
 - II.1 准备工作 1
 - II.2 深度学习模型构建与训练 4
 - II.3 聚类分析：寻找与华师大类似的学校 11

I 实验要求

本次作业的要求如下：

- (1) 在上一节课作业的基础上，请利用深度学习方法，对各学科做一个排名模型，能够较好的预测出排名位置，并且利用 MSE, MAPE 等指标来进行评价模型的优劣。
- (2) 对 ESI 的数据进行聚类，发现与华师大类似的学校有哪些，并分析下原因。

II 具体实现

II.1 准备工作

首先，引入所需的库文件，方便之后的数据分析、数据库导入、可视化等工作。

In [48]:

python

```
1 import sqlite3
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import warnings
7
8 # --- 核心 ML/DL 库 ---
9 import tensorflow as tf
10 from tensorflow import keras
11 from tensorflow.keras import layers
12 from tensorflow.keras.models import Sequential
13
14 # --- Sklearn 库 ---
15 from sklearn.preprocessing import StandardScaler, OneHotEncoder, RobustScaler
16 from sklearn.cluster import KMeans
17 from sklearn.metrics.pairwise import euclidean_distances
18 from sklearn.compose import ColumnTransformer
19 from sklearn.pipeline import Pipeline
20 from sklearn.decomposition import PCA
21 from sklearn.metrics import mean_squared_error, mean_absolute_percentage_error # 导入
    MSE 和 MAPE
22
23 # --- 设置 ---
24 sns.set_theme(style="whitegrid")
25 plt.rcParams['font.sans-serif'] = ['Maple Mono Normal NF CN']
```

```
26 warnings.filterwarnings("ignore", category=FutureWarning)
27 warnings.filterwarnings("ignore", category=UserWarning)
28
29 print(f"TensorFlow 版本: {tf.__version__}")
30 print("所有库导入成功。")
```

txt

```
1 TensorFlow 版本: 2.18.1
2 所有库导入成功。
3
```

然后，定义数据库的路径、学校的名称等变量，方便后续使用。

In [49]:

python

```
1 DB_FILE = 'esi_rankings.db'
2 TABLE_NAME = 'esi_rankings'
3 TARGET_INSTITUTION = 'EAST CHINA NORMAL UNIVERSITY'
4
5 print(f"数据库文件: {DB_FILE}")
6 print(f"目标机构: {TARGET_INSTITUTION}")
```

txt

```
1 数据库文件: esi_rankings.db
2 目标机构: EAST CHINA NORMAL UNIVERSITY
3
```

接下来，类似上次作业的代码，定义加载数据库、数据聚合的函数，并且调用。完成数据的加载与预处理工作。

In [50]:

python

```
1 def load_data(db_file, table_name):
2     """
3     从 SQLite 数据库加载数据。
4     """
5     print(f"正在从 {db_file} 加载数据...")
6     try:
7         conn = sqlite3.connect(db_file)
8         df = pd.read_sql_query(f"SELECT * FROM {table_name}", conn)
9         conn.close()
10        print(f"数据加载成功, 共 {len(df)} 条记录。")
11        return df
12    except Exception as e:
13        print(f"数据加载失败: {e}")
14        return pd.DataFrame()
15
16 def get_aggregated_data(df):
17     """
18     将原始 DataFrame 按机构聚合, 构建高校画像。
19     """
20    print("正在按机构聚合数据以构建高校画像...")
21    if df.empty:
22        print("输入数据为空, 无法聚合。")
23        return pd.DataFrame(), []
24
```

```

25 features = [
26     'total_subjects', 'avg_rank', 'total_documents',
27     'total_cites', 'total_top_papers', 'avg_cites_per_paper'
28 ]
29
30 df_agg = df.groupby('institution').agg(
31     total_subjects=pd.NamedAgg(column='research_field', aggfunc='count'),
32     avg_rank=pd.NamedAgg(column='rank', aggfunc='mean'),
33     total_documents=pd.NamedAgg(column='documents', aggfunc='sum'),
34     total_cites=pd.NamedAgg(column='cites', aggfunc='sum'),
35     total_top_papers=pd.NamedAgg(column='top_papers', aggfunc='sum'),
36     avg_cites_per_paper=pd.NamedAgg(column='cites_per_paper', aggfunc='mean')
37 ).reset_index()
38
39 print(f"数据聚合完成, 共 {len(df_agg)} 所机构. ")
40 return df_agg, features
41
42 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

```

In [51]:

```

1 # -----
2 # 单元格 4: 加载原始数据
3 # -----
4 df_raw = load_data(DB_FILE, TABLE_NAME)
5
6 # 预览数据
7 if not df_raw.empty:
8     print("\n数据预览 (前5行):")
9     print(df_raw.head())
10    print("\n数据信息:")
11    df_raw.info()

```

python

txt

```

1 正在从 esi_rankings.db 加载数据...
2 数据加载成功, 共 34121 条记录.
3
4 数据预览 (前5行):
5 |   id   | research_field | rank | \
6 0 | 1 | AGRICULTURAL SCIENCES | 1
7 1 | 2 | AGRICULTURAL SCIENCES | 2
8 2 | 3 | AGRICULTURAL SCIENCES | 3
9 3 | 4 | AGRICULTURAL SCIENCES | 4
10 4 | 5 | AGRICULTURAL SCIENCES | 5
11
12 |   |   |   |   | institution | country_region | documents | \
13 0 |   |   |   |   | CHINESE ACADEMY OF SCIENCES | CHINA MAINLAND | 15661
14 1 |   |   |   |   | CHINESE ACADEMY OF AGRICULTURAL SCIENCES | CHINA MAINLAND | 12222
15 2 |   |   |   |   | UNITED STATES DEPARTMENT OF AGRICULTURE (USDA) | USA | 12564
16 3 |   |   |   |   | CHINA AGRICULTURAL UNIVERSITY | CHINA MAINLAND | 10052
17 4 |   |   |   |   | INRAE | FRANCE | 9314
18
19 |   |   |   |   |   |   |   |
20 0 |   |   |   |   |   |   |   |
21 1 |   |   |   |   |   |   |   |

```

```

22 2 220644          17.56      105
23 3 207779          20.67      166
24 4 187838          20.17      118
25
26 数据信息:
27 <class 'pandas.core.frame.DataFrame'>
28 RangeIndex: 34121 entries, 0 to 34120
29 Data columns (total 9 columns):
30 #   Column          Non-Null Count  Dtype
31 ---  ---
32 0    id              34121 non-null  int64
33 1    research_field  34121 non-null  object
34 2    rank            34121 non-null  int64
35 3    institution     34121 non-null  object
36 4    country_region  30960 non-null  object
37 5    documents       34121 non-null  int64
38 6    cites           34121 non-null  int64
39 7    cites_per_paper 34121 non-null  float64
40 8    top_papers      34121 non-null  int64
41 dtypes: float64(1), int64(5), object(3)
42 memory usage: 2.3+ MB
43

```

II.2 深度学习模型构建与训练

- 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.
- 对数变换: rank 值的分布非常偏态 (从 1 到几千). 神经网络在预测这种大跨度值时表现不佳. 我们通过 $\text{np.log}(y)$ 将其压缩, 模型去预测 $\log(\text{rank})$, 最后再用 $\text{np.exp}()$ 还原回来.
- 数值特征: 必须使用 StandardScaler 进行标准化.
- 类别特征: 必须使用 OneHotEncoder 进行独热编码.

In [52]:

```

1 # -----
2 # 单元格 5: 任务1 - 数据准备 (DL版)
3 # -----
4 print("--- 开始执行任务1: 深度学习排名预测模型 ---")
5
6 # 1. 准备数据
7 df_model = df_raw.dropna(subset=[
8     'research_field', 'rank', 'documents',
9     'cites', 'cites_per_paper', 'top_papers'
10 ])
11
12 target = 'rank'
13 numeric_features = ['documents', 'cites', 'cites_per_paper', 'top_papers']

```

python

```
14 categorical_features = ['research_field']
15
16 X = df_model[numeric_features + categorical_features]
17 # 2. 对数变换目标值 (y)
18 y = np.log(df_model[target]) # 关键: 对 y 进行 log 变换
19 print(f"步骤 1/4: 已对目标 'rank' 进行 Log 变换。")
20
21
22 # 3. 自定义数据拆分 (按学科内排名 60% 训练, 20% 测试)
23 print("步骤 2/4: 正在按学科内排名拆分 60% 训练集和 20% 测试集...")
24 train_indices = []
25 test_indices = []
26
27 for field in df_model['research_field'].unique():
28     df_field = df_model[df_model['research_field'] == field].copy()
29     # 必须按 rank 排序来切分
30     df_field_sorted = df_field.sort_values(by='rank', ascending=True)
31
32     n = len(df_field_sorted)
33     train_end = int(n * 0.6)
34     test_start = int(n * 0.8)
35
36     train_indices.extend(df_field_sorted.iloc[:train_end].index)
37     test_indices.extend(df_field_sorted.iloc[test_start:].index)
38
39 # 使用 .loc 保证索引对齐
40 X_train = X.loc[train_indices]
41 y_train_log = y.loc[train_indices]
42 X_test = X.loc[test_indices]
43 y_test_log = y.loc[test_indices]
44 # 保留原始 y_test 以便评估
45 y_test_orig = df_model.loc[test_indices][target]
46
47 print(f"数据拆分完成. 训练集大小: {len(X_train)}, 测试集大小: {len(X_test)}")
48
49 # 4. 创建数据预处理 Pipeline
50 print("步骤 3/4: 正在构建数据预处理 Pipeline (标准化 + 独热编码)...")
51 preprocessor = ColumnTransformer(
52     transformers=[
53         # 神经网络对尺度敏感, 必须标准化
54         ('num', StandardScaler(), numeric_features),
55         ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
56     ])
57
58 # 5. 应用预处理
59 X_train_processed = preprocessor.fit_transform(X_train)
60 X_test_processed = preprocessor.transform(X_test)
61
62 # Keras 需要密集的 numpy 数组
63 if hasattr(X_train_processed, "toarray"):
64     X_train_processed = X_train_processed.toarray()
65     X_test_processed = X_test_processed.toarray()
66
67 print("步骤 4/4: 数据预处理完成。")
68 print(f"处理后训练集特征维度: {X_train_processed.shape}")
```

txt

1 --- 开始执行任务1: 深度学习排名预测模型 ---

```
2 步骤 1/4: 已对目标 'rank' 进行 Log 变换.
3 步骤 2/4: 正在按学科内排名拆分 60% 训练集和 20% 测试集...
4 数据拆分完成. 训练集大小: 20464, 测试集大小: 6833
5 步骤 3/4: 正在构建数据预处理 Pipeline (标准化 + 独热编码)...
6 步骤 4/4: 数据预处理完成.
7 处理后训练集特征维度: (20464, 26)
8
```

我们将构建一个简单的多层感知机.

- 输入层: 维度必须等于处理后的特征数量.
- 隐藏层: 使用 Dense 层和 relu 激活函数, 加入 Dropout 防止过拟合.
- 输出层: Dense(1) (因为是回归问题), 使用 linear 激活函数 (默认) .

In [53]:

python

```
1 # -----
2 # 单元格 6: 任务1 - 构建深度学习模型 (MLP)
3 # -----
4
5 # 获取输入维度
6 input_dim = X_train_processed.shape[1]
7
8 def build_model(input_shape):
9     model = Sequential([
10         # 输入层
11         layers.Input(shape=(input_shape,)),
12
13         # 第一个隐藏层
14         layers.Dense(128, activation='relu'),
15         layers.Dropout(0.2), # Dropout 正则化
16
17         # 第二个隐藏层
18         layers.Dense(64, activation='relu'),
19         layers.Dropout(0.1),
20
21         # 第三个隐藏层
22         layers.Dense(32, activation='relu'),
23
24         # 输出层 (回归问题, 1个神经元, 线性激活)
25         layers.Dense(1)
26     ])
27
28     # 编译模型
29     # 优化器: Adam 是一个好的起点
30     # 损失函数: MSE 是回归问题的标准损失函数
31     model.compile(optimizer='adam',
32                   loss='mean_squared_error', # 监控 MSE
33                   metrics=['mean_absolute_error']) # 也可以监控 MAE
34
35     return model
36
37 model = build_model(input_dim)
38
39 # 打印模型结构
```

```
40 model.summary()
```

txt

```
1 [1mModel: "sequential_5"[0m
2
```

txt

```
1
2 | [1m [0m[1mLayer (type) | [0m[1m [0m[1mOutput Shape
  | [0m[1m [0m[1m Param #[0m[1m [0m[1m
3 |-----|-----|-----|
4 | dense_20 ([38;5;33mDense[0m) | ([38;5;45mNone[0m,
  | [38;5;34m128[0m) | [38;5;34m3,456[0m |
5 |-----|-----|-----|
6 | dropout_10 ([38;5;33mDropout[0m) | ([38;5;45mNone[0m,
  | [38;5;34m128[0m) | [38;5;34m0[0m |
7 |-----|-----|-----|
8 | dense_21 ([38;5;33mDense[0m) | ([38;5;45mNone[0m,
  | [38;5;34m64[0m) | [38;5;34m8,256[0m |
9 |-----|-----|-----|
10 | dropout_11 ([38;5;33mDropout[0m) | ([38;5;45mNone[0m,
   | [38;5;34m64[0m) | [38;5;34m0[0m |
11 |-----|-----|-----|
12 | dense_22 ([38;5;33mDense[0m) | ([38;5;45mNone[0m,
   | [38;5;34m32[0m) | [38;5;34m2,080[0m |
13 |-----|-----|-----|
14 | dense_23 ([38;5;33mDense[0m) | ([38;5;45mNone[0m,
   | [38;5;34m1[0m) | [38;5;34m33[0m |
15 |-----|-----|-----|
16
```

txt

```
1 [1m Total params: [0m[38;5;34m13,825[0m (54.00 KB)
2
```

txt

```
1 [1m Trainable params: [0m[38;5;34m13,825[0m (54.00 KB)
2
```

txt

```
1 [1m Non-trainable params: [0m[38;5;34m0[0m (0.00 B)
2
```

我们使用 `fit` 方法训练模型。 `validation_split=0.2`: Keras 会自动从训练集中分出 20% 作为验证集，用于在每个 epoch 结束时评估模型，帮助我们监控是否过拟合。

In [54]:

python

```
1 # -----
2 # 单元格 7: 任务1 - 训练模型
3 # -----
4 print("开始训练深度学习模型...")
5
6 # 使用 EarlyStopping 可以在模型不再改进时提前停止训练
7 early_stopping = keras.callbacks.EarlyStopping(
8     monitor='val_loss', # 监控验证集损失
9     patience=10,        # 10个 epoch 内没改进就停止
10    restore_best_weights=True # 恢复到最佳权重
11 )
12
```

```

13 # 训练模型
14 history = model.fit(
15     X_train_processed,
16     y_train_log,
17     epochs=100, # 增加 Epoch 数量, 让 EarlyStopping 来决定何时停止
18     batch_size=32,
19     validation_split=0.2, # 从训练集中分出 20% 作为验证集
20     callbacks=[early_stopping],
21     verbose=1 # 打印训练过程
22 )
23
24 print("模型训练完成. ")

```

txt

```

1 开始训练深度学习模型...
2 Epoch 1/100
3 1m512/512[0m 32m-----[0m[37m[0m 1m2s[0m 1ms/step - loss:
  2.8775 - mean_absolute_error: 0.9195 - val_loss: 1.3348 - val_mean_absolute_error:
  1.0605
4 Epoch 2/100
5 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.3377 - mean_absolute_error: 0.4426 - val_loss: 1.0521 - val_mean_absolute_error:
  0.9314
6 Epoch 3/100
7 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.2554 - mean_absolute_error: 0.3895 - val_loss: 0.7495 - val_mean_absolute_error:
  0.7529
8 Epoch 4/100
9 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.2051 - mean_absolute_error: 0.3485 - val_loss: 0.5036 - val_mean_absolute_error:
  0.5770
10 Epoch 5/100
11 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.1704 - mean_absolute_error: 0.3155 - val_loss: 0.2831 - val_mean_absolute_error:
  0.3675
12 Epoch 6/100
13 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.1413 - mean_absolute_error: 0.2891 - val_loss: 0.2355 - val_mean_absolute_error:
  0.3272
14 Epoch 7/100
15 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.1233 - mean_absolute_error: 0.2657 - val_loss: 0.2009 - val_mean_absolute_error:
  0.3561
16 Epoch 8/100
17 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.1011 - mean_absolute_error: 0.2399 - val_loss: 0.1999 - val_mean_absolute_error:
  0.3663
18 Epoch 9/100
19 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.0898 - mean_absolute_error: 0.2217 - val_loss: 0.2478 - val_mean_absolute_error:
  0.4306
20 Epoch 10/100
21 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.0779 - mean_absolute_error: 0.2032 - val_loss: 0.3057 - val_mean_absolute_error:
  0.4950
22 Epoch 11/100
23 1m512/512[0m 32m-----[0m[37m[0m 1m1s[0m 1ms/step - loss:
  0.0612 - mean_absolute_error: 0.1804 - val_loss: 0.3538 - val_mean_absolute_error:
  0.5329

```



```

24 Epoch 12/100
25 [1m512/512][0m 32m-----[0m][37m][0m 1m1s][0m 1ms/step - loss:
    0.0478 - mean_absolute_error: 0.1581 - val_loss: 0.3303 - val_mean_absolute_error:
    0.5044
26 Epoch 13/100
27 [1m512/512][0m 32m-----[0m][37m][0m 1m1s][0m 1ms/step - loss:
    0.0435 - mean_absolute_error: 0.1480 - val_loss: 0.3333 - val_mean_absolute_error:
    0.5173
28 Epoch 14/100
29 [1m512/512][0m 32m-----[0m][37m][0m 1m1s][0m 1ms/step - loss:
    0.0362 - mean_absolute_error: 0.1352 - val_loss: 0.3093 - val_mean_absolute_error:
    0.4747
30 Epoch 15/100
31 [1m512/512][0m 32m-----[0m][37m][0m 1m1s][0m 1ms/step - loss:
    0.0324 - mean_absolute_error: 0.1247 - val_loss: 0.2836 - val_mean_absolute_error:
    0.4507
32 Epoch 16/100
33 [1m512/512][0m 32m-----[0m][37m][0m 1m1s][0m 1ms/step - loss:
    0.0295 - mean_absolute_error: 0.1172 - val_loss: 0.3364 - val_mean_absolute_error:
    0.5007
34 Epoch 17/100
35 [1m512/512][0m 32m-----[0m][37m][0m 1m1s][0m 1ms/step - loss:
    0.0290 - mean_absolute_error: 0.1134 - val_loss: 0.4499 - val_mean_absolute_error:
    0.5997
36 Epoch 18/100
37 [1m512/512][0m 32m-----[0m][37m][0m 1m1s][0m 1ms/step - loss:
    0.0277 - mean_absolute_error: 0.1115 - val_loss: 0.3668 - val_mean_absolute_error:
    0.5362
38 模型训练完成.
39
  
```

模型预测的是 $\log(\text{rank})$ ，我们必须使用 `np.exp()` 将其还原为 `rank`，才能与 `y_test_orig` (原始排名) 进行比较并计算 MSE 和 MAPE.

In [55]:

python

```

1 # -----
2 # 单元格 8: 任务1 - 评估模型 (MSE, MAPE)
3 # -----
4
5 # 1. 绘制训练过程中的损失
6 print("正在绘制训练和验证损失...")
7 plt.figure(figsize=(10, 6))
8 plt.plot(history.history['loss'], label='训练集损失 (Train Loss)')
9 plt.plot(history.history['val_loss'], label='验证集损失 (Validation Loss)')
10 plt.title('模型训练过程中的损失 (MSE)')
11 plt.xlabel('Epoch')
12 plt.ylabel('Loss (Log Scale)')
13 plt.yscale('log') # 损失值可能下降很快, 用 log 尺度看得更清
14 plt.legend()
15 plt.savefig('dl_loss_curve.png')
16 plt.show()
17
18 # 2. 在测试集上评估
19 print("\n正在测试集上评估模型...")
20 # 预测 log(rank)
21 y_pred_log = model.predict(X_test_processed).flatten() # flatten 转为 1D 数组
  
```

```

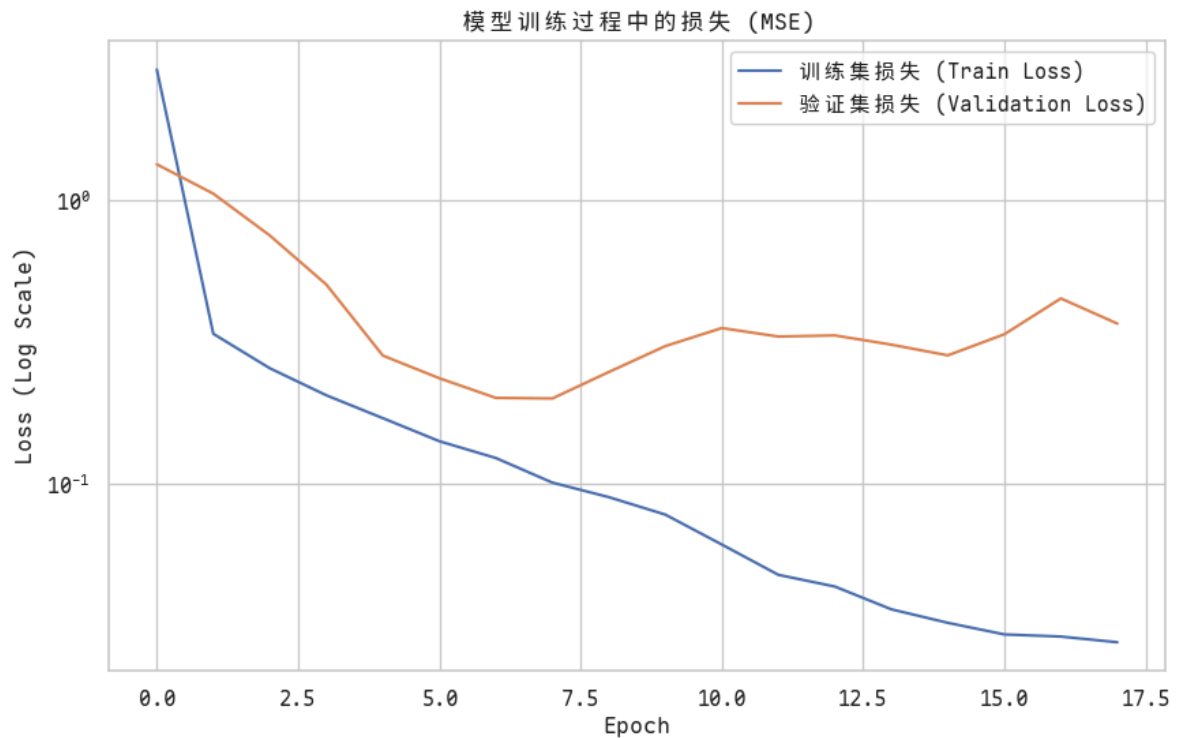
22
23 # 3. 还原预测值
24 # 关键: 将 log(rank) 转换回 rank
25 y_pred_orig = np.exp(y_pred_log)
26
27 # 4. 计算 MSE 和 MAPE
28 # y_test_orig 是我们在单元格 5 中保留的原始排名
29 mse = mean_squared_error(y_test_orig, y_pred_orig)
30 rmse = np.sqrt(mse)
31 mape = mean_absolute_percentage_error(y_test_orig, y_pred_orig)
32
33 print("\n--- 任务1 结果: 深度学习模型评估 (测试集) ---")
34 print(f"R² (决定系数): 无法直接比较, 因为模型是基于 Log 变换的. ")
35 print(f"MSE (均方误差): {mse:.4f}")
36 print(f"RMSE (均方根误差): {rmse:.4f} (预测排名平均误差约 {rmse:.0f} 名)")
37 print(f"MAPE (平均绝对百分比误差): {mape:.4f} (预测平均偏离 {mape*100:.2f}%)")
38
39 # 5. 抽样对比真实值和预测值
40 print("\n--- 预测结果抽样对比 (前20条) ---")
41 df_compare = pd.DataFrame({
42     '真实排名 (Actual Rank)': y_test_orig.values,
43     '预测排名 (Predicted Rank)': y_pred_orig
44 })
45 print(df_compare.head(20))

```

txt

1 正在绘制训练和验证损失...

2



txt

1

2 正在测试集上评估模型...

3 [1m214/214][0m 32m] [0m 37m][0m 1m0s][0m 585us/step

```

4
5 --- 任务1 结果: 深度学习模型评估 (测试集) ---
6 R2 (决定系数): 无法直接比较, 因为模型是基于 Log 变换的.
7 MSE (均方误差): 2062449.0000
8 RMSE (均方根误差): 1436.1229 (预测排名平均误差约 1436 名)
9 MAPE (平均绝对百分比误差): 0.3841 (预测平均偏离 38.41%)
10
11 --- 预测结果抽样对比 (前20条) ---
12 | 真实排名 (Actual Rank) | 预测排名 (Predicted Rank)
13 | 0 | 1105 | 702.721558
14 | 1 | 1106 | 693.597351
15 | 2 | 1106 | 692.361816
16 | 3 | 1106 | 711.410950
17 | 4 | 1109 | 690.580994
18 | 5 | 1110 | 686.135437
19 | 6 | 1111 | 699.151978
20 | 7 | 1112 | 707.306641
21 | 8 | 1113 | 713.290405
22 | 9 | 1114 | 698.660034
23 | 10 | 1115 | 709.850525
24 | 11 | 1116 | 614.558228
25 | 12 | 1117 | 665.700745
26 | 13 | 1118 | 662.359314
27 | 14 | 1119 | 715.645874
28 | 15 | 1120 | 679.368042
29 | 16 | 1121 | 711.627075
30 | 17 | 1122 | 687.464417
31 | 18 | 1123 | 712.421570
32 | 19 | 1124 | 697.000916
33

```

II.3 聚类分析: 寻找与华师大类似的学校

- 方法 1 (聚类): 使用 K-Means 找到 ECNU 所在的“群体”.
- 方法 2 (距离): 使用欧氏距离找到 ECNU 的“近邻”或“统计意义上的双胞胎”.

原因是通过比较 ECNU 所在聚类的中心特征与 ECNU 自身的特征, 我们可以发现它们在“学科广度”、“平均排名”、“科研产出”等画像上具有相似性, 这就是“为什么”它们被分到一类.

我们使用了 `StandardScaler` 进行标准化. `StandardScaler` 通过 (均值) 和 (标准差) 来缩放数据, 这使它对极端异常值非常敏感.

在高校数据中, 像 `total_cites` (总引用) 或 `total_documents` (总论文) 这样的特征是“幂律分布”的——大多数学校产出中等, 但少数机构 (如中科院、哈佛) 的产出是其他学校的几百倍.

这些超级机构的极端数值“拉偏”了 `StandardScaler`, 导致它把 99% 的“普通”大学 (包括华师大) 全都压缩到了一个很小的数值范围内, 使它们看起来都“一样” (都在 Cluster 0), 从而聚类失败.

要解决这个问题, 我们需要“稳健”的预处理方法.

- 策略一: 使用 `RobustScaler` (稳健标准化) `RobustScaler` 不使用均值和标准差, 而是使用四分位数 (IQR) 来缩放数据. 这意味着它在计算缩放比例时会忽略掉那些极端的异常值. 这是解决当前问题的首选方案.

- 策略二：使用 `np.log1p` (对数变换) 对数变换是处理“幂律分布”数据的经典方法. 它可以“压缩”数据的长尾, 将 `[100, 1000, 10000000]` 这样的数据点, 变换到更接近的尺度 `[4.6, 6.9, 16.1]`, 使数据分布更接近正态分布, K-Means 会更有效.

我们将同时使用这两种策略来构建一个优化的聚类流程.

In [56]:

python

```
1 # -----
2 # 单元格 9 (优化版): 任务2 - 稳健聚类与相似性分析
3 # -----
4 print("\n--- 开始执行任务2 (优化版): 稳健聚类与相似高校分析 ---")
5
6 # 1. 获取聚合数据
7 df_agg, agg_features = get_aggregated_data(df_raw)
8
9 if df_agg.empty:
10     print("聚合数据为空, 跳过任务2. ")
11 else:
12     # 复制一份用于变换
13     df_agg_processed = df_agg.set_index('institution')[agg_features].fillna(0).copy()
14
15     # 2. [优化]: 对数变换 (Log Transform)
16     # 我们对具有长尾分布的特征应用 log1p (加1是为了避免 log(0))
17     log_features = ['total_documents', 'total_cites', 'total_top_papers']
18     print(f"步骤 1/5: 正在对 {log_features} 应用 Log(1+p) 变换...")
19
20     for col in log_features:
21         # 确保数据非负
22         df_agg_processed[col] = df_agg_processed[col].apply(lambda x: np.log1p(x if x
23 > 0 else 0))
24
25     # 3. [优化]: 使用 RobustScaler
26     print("步骤 2/5: 正在使用 RobustScaler (稳健标准化)...")
27     scaler = RobustScaler() # <--- 关键优化!
28
29     # 对所有特征 (包括已 log 变换的) 进行稳健标准化
30     features_scaled = scaler.fit_transform(df_agg_processed)
31
32     # 创建标准化的 DataFrame
33     df_scaled = pd.DataFrame(features_scaled, index=df_agg_processed.index,
34 columns=agg_features)
35
36     # 4. K-Means 聚类 (K=5)
37     K = 5
38     kmeans = KMeans(n_clusters=K, random_state=42, n_init=10)
39     df_agg['cluster'] = kmeans.fit_predict(features_scaled)
40     df_scaled['cluster'] = df_agg['cluster'].values
41     print(f"步骤 3/5: K-Means 聚类完成 (K={K}). ")
42
43     # 5. 分析华师大所在的聚类
44     try:
45         ecnu_cluster = df_scaled.loc[TARGET_INSTITUTION]['cluster']
46         print(f"\n--- 任务2 结果 (聚类) ---")
47         print(f"'{TARGET_INSTITUTION}' 位于聚类: {int(ecnu_cluster)}")
48
49     # 6. 分析原因: 查看该聚类的中心特征
50     print("\n[分析原因]: 华师大所在聚类的中心特征 (标准化后) 如下:")
```

```

49 cluster_center = df_scaled[df_scaled['cluster'] == ecnu_cluster]
   [agg_features].mean()
50 print(cluster_center)
51 print(f"\n华师大自身的特征 (标准化后):")
52 print(df_scaled.loc[TARGET_INSTITUTION][agg_features])
53
54 # 7. 寻找相似高校 (欧氏距离) - 同样在优化后的空间中进行
55 print(f"\n步骤 4/5: 正在寻找与 '{TARGET_INSTITUTION}' 最相似的近邻...")
56
57 target_vector_scaled = df_scaled.loc[TARGET_INSTITUTION]
   [agg_features].values.reshape(1, -1)
58 distances = euclidean_distances(target_vector_scaled,
   df_scaled[agg_features].values)
59 df_scaled['distance_to_target'] = distances[0]
60
61 df_similar_euclidean = df_scaled.sort_values(by='distance_to_target',
   ascending=True)
62 top_similar_euclidean = df_similar_euclidean.iloc[0:11]
63
64 # 提取原始聚合数据 (未经 log 变换) 以提高可读性
65 original_data_similar =
   df_agg.set_index('institution').loc[top_similar_euclidean.index]
66 original_data_similar['distance'] =
   top_similar_euclidean['distance_to_target']
67
68 print(f"\n--- T任务2 结果 (近邻) ---")
69 print(f"与 {TARGET_INSTITUTION} 画像最相似的10所高校 (基于优化后的空间):")
70 print(original_data_similar[['total_subjects', 'avg_rank',
   'avg_cites_per_paper', 'distance']])
71 print("步骤 5/5: 相似高校分析完成。")
72
73 except KeyError:
74     print(f"错误: 在数据中未找到机构 '{TARGET_INSTITUTION}'. ")
75 except Exception as e:
76     print(f"执行任务2时出错: {e}")

```

txt

```

1
2 --- 开始执行任务2 (优化版): 稳健聚类与相似高校分析 ---
3 正在按机构聚合数据以构建高校画像...
4 数据聚合完成, 共 9990 所机构.
5 步骤 1/5: 正在对 ['total_documents', 'total_cites', 'total_top_papers'] 应用 Log(1+p)
   变换...
6 步骤 2/5: 正在使用 RobustScaler (稳健标准化)...
7 步骤 3/5: K-Means 聚类完成 (K=5).
8
9 --- 任务2 结果 (聚类) ---
10 'EAST CHINA NORMAL UNIVERSITY' 位于聚类: 0
11
12 [分析原因]: 华师大所在聚类的中心特征 (标准化后) 如下:
13 total_subjects      0.887123
14 avg_rank            0.190960
15 total_documents    0.154672
16 total_cites        0.216767
17 total_top_papers   0.203945
18 avg_cites_per_paper 0.621544
19 dtype: float64
20

```

[illegible]

对高维数据（我们聚类时用了 6 个特征）进行聚类后，如果能将其“拍扁”到二维平面上展示出来，就能非常直观地看到：

- 不同类别是否在空间上被分开了.
- 华师大位于哪个群体的什么位置.

要实现这一点，最经典的方法是使用 PCA 进行降维。

In [57]:

```
1 # -----
2 # 单元格 10: 任务2 - 聚类可视化 (PCA 降维)
3 # (依赖 单元格 9 的运行结果)
4 # -----
```

python

```

5
6 print("\n--- 开始执行 任务2 (可视化部分) ---")
7
8 # 1. 检查所需数据是否存在 (来自 Cell 9)
9 if 'df_scaled' in locals() and 'df_agg' in locals() and 'agg_features' in locals():
10     print("步骤 1/4: 正在使用 PCA 进行降维 (6D -> 2D)...")
11
12     # 2. PCA 降维
13     # 我们对用于聚类的标准化数据 (df_scaled) 进行降维
14     features_to_reduce = df_scaled[agg_features]
15
16     pca = PCA(n_components=2) # 我们希望降到 2 维 (PC1, PC2)
17     pca_result = pca.fit_transform(features_to_reduce)
18
19     # 3. 创建用于可视化的 DataFrame
20     df_viz = pd.DataFrame(data=pca_result, columns=['PC1', 'PC2'])
21     # 附加聚类标签 (来自 df_agg)
22     df_viz['cluster'] = df_agg['cluster'].values
23     # 附加机构名称 (来自 df_agg)
24     df_viz['institution'] = df_agg['institution'].values
25
26     print(f"步骤 2/4: PCA 降维完成. ")
27     print(f" - PC1 解释了 {pca.explained_variance_ratio_[0]:.2%} 的方差")
28     print(f" - PC2 解释了 {pca.explained_variance_ratio_[1]:.2%} 的方差")
29     print(f" - 总共解释了 {np.sum(pca.explained_variance_ratio_):.2%} 的方差")
30
31     # 4. 绘制散点图
32     print("步骤 3/4: 正在绘制聚类散点图...")
33     plt.figure(figsize=(14, 9))
34
35     # 使用 seaborn 绘制散点图, 并按 'cluster' 着色
36     sns.scatterplot(
37         data=df_viz,
38         x='PC1',
39         y='PC2',
40         hue='cluster',          # 按聚类结果上色
41         palette='viridis',      # 使用 'viridis' 色板 (区分度高)
42         s=50,                   # 点的大小
43         alpha=0.6,              # 透明度, 防止点重叠
44         legend='full'
45     )
46
47     # 5. 在图上标记华师大 (ECNU)
48     print(f"步骤 4/4: 正在图上标记 {TARGET_INSTITUTION}...")
49     try:
50         # 找到 ECNU 降维后的坐标
51         ecnu_viz = df_viz[df_viz['institution'] == TARGET_INSTITUTION].iloc[0]
52
53         # 使用一个醒目的标记 (红色 'x') 来突出显示
54         plt.scatter(
55             ecnu_viz['PC1'],
56             ecnu_viz['PC2'],
57             marker='x',          # 标记样式
58             color='red',         # 标记颜色
59             s=200,               # 标记大小
60             label=TARGET_INSTITUTION, # 标签

```

```

61         zorder=5                                # 确保在顶层显示
62     )
63     # 添加文本注释
64     plt.text(
65         ecnu_viz['PC1'] + 0.05, # 在 X 坐标旁偏移一点
66         ecnu_viz['PC2'],
67         f" <-- {TARGET_INSTITUTION}", # 注释内容
68         horizontalalignment='left',
69         color='red',
70         weight='bold',
71         fontsize=12
72     )
73
74     except IndexError:
75         print(f"警告: 未能在可视化数据中找到 {TARGET_INSTITUTION} 进行标记。")
76     except Exception as e:
77         print(f"标记 {TARGET_INSTITUTION} 时出错: {e}")
78
79     plt.title('全球高校 K-Means 聚类可视化 (PCA 降维)', fontsize=18)
80     plt.xlabel(f'第一主成分 (PC1) - 解释 {pca.explained_variance_ratio_[0]:.2%} 方差',
81               fontsize=12)
82     plt.ylabel(f'第二主成分 (PC2) - 解释 {pca.explained_variance_ratio_[1]:.2%} 方差',
83               fontsize=12)
84     plt.legend(title='聚类 (Cluster)')
85     plt.grid(True, linestyle='--', alpha=0.5)
86     plt.savefig('task2_cluster_visualization_pca.png')
87     plt.show()
88 else:
89     print("错误: 未找到 Cell 9 中生成的 'df_scaled' 或 'df_agg'. 请先运行 Cell 9. ")

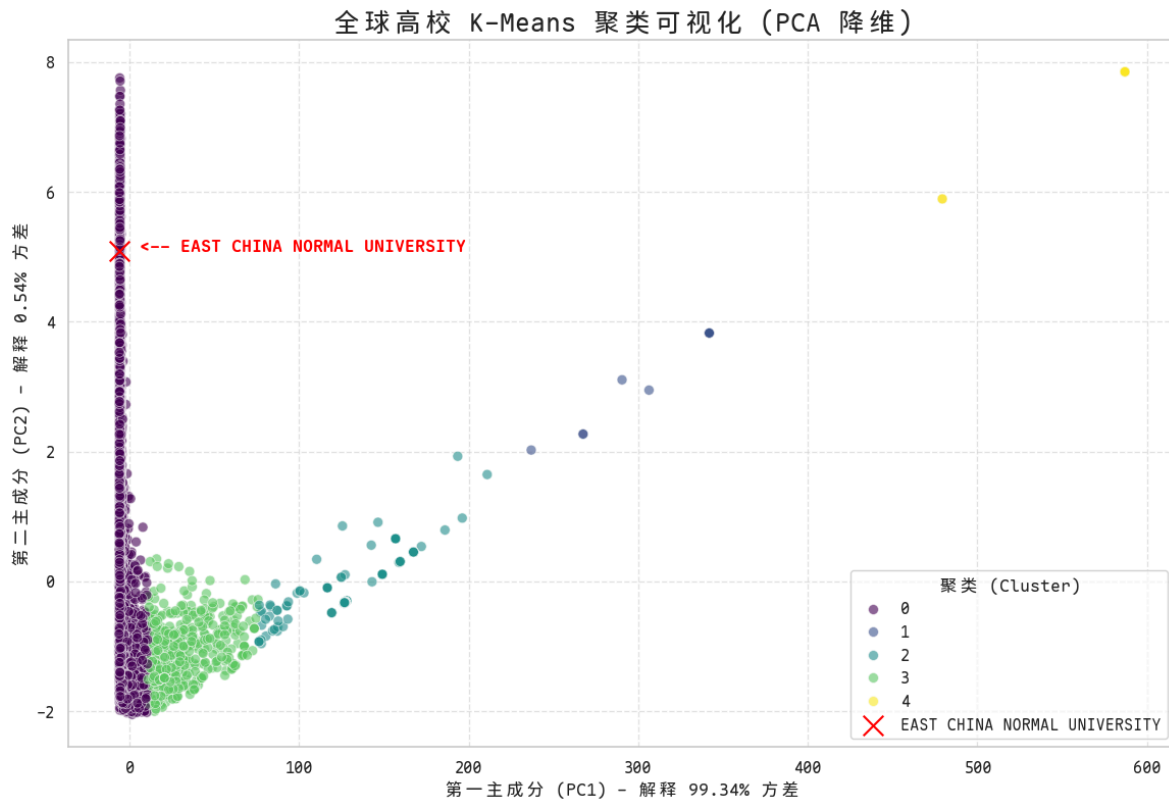
```

txt

```

1
2 --- 开始执行 任务2 (可视化部分) ---
3 步骤 1/4: 正在使用 PCA 进行降维 (6D -> 2D)...
4 步骤 2/4: PCA 降维完成.
5   - PC1 解释了 99.34% 的方差
6   - PC2 解释了 0.54% 的方差
7   - 总共解释了 99.89% 的方差
8 步骤 3/4: 正在绘制聚类散点图...
9 步骤 4/4: 正在图上标记 EAST CHINA NORMAL UNIVERSITY...
10

```

这张图的结果非常理想，它告诉我们：

- 我们成功地将全球高校分成了几个有意义的类别。
- 华师大属于 **Cluster 0** (紫色). 这个群体是最大的群体，代表了全球“高水平、综合性、规模可观”的主流大学。

其他群体：

- 最右侧的黄色群体 (**Cluster 3**)：这很可能是那些在各项指标上都极端领先的“超级机构”（如中科院、哈佛等）
- 左侧的蓝绿色群体 (**Cluster 1**)：这可能代表了另一类特征的机构，比如“小而精”或者在特定领域（如低 `total_cites` 但高 `avg_rank`）有特点的机构。

Remark

完整代码见 `./hw6.ipynb`.