

数据科学导论 -HW 7 报告

孙育泉 10234900421
2025.11.08

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

I 实验要求

本次作业主要是改进上次的作业，这次主要改进的任务是：

(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。

II 具体实现

本次使用表格深度学习的方法。

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 [1]:

```
python
1 # -----
2 # 单元格 1: 导入所有库
3 # -----
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 # --- 核心 DL 库 (用于任务1) ---
12 import tensorflow as tf
13 from tensorflow import keras
14 from tensorflow.keras import layers
15 from tensorflow.keras.models import Sequential
16
17 # --- Sklearn 库 (用于任务1预处理 和 任务2) ---
18 from sklearn.preprocessing import StandardScaler, OneHotEncoder, RobustScaler # 导入两
    种 Scaler
19 from sklearn.cluster import KMeans
20 from sklearn.metrics.pairwise import euclidean_distances
21 from sklearn.compose import ColumnTransformer
22 from sklearn.pipeline import Pipeline
23 from sklearn.metrics import mean_squared_error, mean_absolute_percentage_error # 导入
    MSE 和 MAPE
24 from sklearn.decomposition import PCA # 用于可视化
25
```

```

26 # --- 设置 ---
27 sns.set_theme(style="whitegrid")
28 plt.rcParams['font.sans-serif'] = ['Heiti TC'] # 如果您用中文标签，请取消这行注释
29 # plt.rcParams['axes.unicode_minus'] = False
30 warnings.filterwarnings("ignore", category=FutureWarning)
31 warnings.filterwarnings("ignore", category=UserWarning)
32
33 print(f"TensorFlow 版本: {tf.__version__}")
34 print("所有库导入成功。")

```

txt

```

1 TensorFlow 版本: 2.16.2
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 [2]:

```

1 # -----
2 # 单元格 2: 全局配置
3 # -----
4 DB_FILE = 'esi_rankings.db'
5 TABLE_NAME = 'esi_rankings'
6 TARGET_INSTITUTION = 'EAST CHINA NORMAL UNIVERSITY'
7
8 print(f"数据库文件: {DB_FILE}")
9 print(f"目标机构: {TARGET_INSTITUTION}")

```

python

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 [3]:

```

1 # -----
2 # 单元格 3: 共用函数
3 # -----
4
5 def load_data(db_file, table_name):
6     """
7     从 SQLite 数据库加载数据。
8     """
9     print(f"正在从 {db_file} 加载数据...")
10    try:
11        conn = sqlite3.connect(db_file)
12        df = pd.read_sql_query(f"SELECT * FROM {table_name}", conn)
13        conn.close()
14        print(f"数据加载成功，共 {len(df)} 条记录。")
15        return df
16    except Exception as e:
17        print(f"数据加载失败: {e}")

```

python

```

18         return pd.DataFrame()
19
20 def get_aggregated_data(df):
21     """
22     将原始 DataFrame 按机构聚合，构建高校画像（用于任务2）。
23     """
24     print("正在按机构聚合数据以构建高校画像...")
25     if df.empty:
26         print("输入数据为空，无法聚合。")
27         return pd.DataFrame(), []
28
29     features = [
30         'total_subjects', 'avg_rank', 'total_documents',
31         'total_cites', 'total_top_papers', 'avg_cites_per_paper'
32     ]
33
34     df_agg = df.groupby('institution').agg(
35         total_subjects=pd.NamedAgg(column='research_field', aggfunc='count'),
36         avg_rank=pd.NamedAgg(column='rank', aggfunc='mean'),
37         total_documents=pd.NamedAgg(column='documents', aggfunc='sum'),
38         total_cites=pd.NamedAgg(column='cites', aggfunc='sum'),
39         total_top_papers=pd.NamedAgg(column='top_papers', aggfunc='sum'),
40         avg_cites_per_paper=pd.NamedAgg(column='cites_per_paper', aggfunc='mean')
41     ).reset_index()
42
43     print(f"数据聚合完成，共 {len(df_agg)} 所机构。")
44     return df_agg, features
45
46 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 [4]:

```

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

```

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 | cites cites_per_paper top_papers
20 0 332254 21.22 251
21 1 223855 18.32 198
22 2 220644 17.56 105
23 3 207779 20.67 166
24 4 187838 20.17 118
25

```

II.2 深度学习模型构建与训练

- (1) 对目标 y (rank) 进行 `np.log()` 变换, 使分布更平滑, 便于神经网络学习.
- (2) 对数值特征使用 `StandardScaler` (因为是对原始数据, 不是聚合数据, `StandardScaler` 在这里是合适的).
- (3) 对类别特征使用 `OneHotEncoder`.
- (4) 对数据集进行打乱后, 按照 80% / 20% 划分训练集和测试集, 然后进行分层抽样, 尽可能让训练集与测试集同分布.

In [11]:

```

1 # -----
2 # 单元格 5 (改进版): 任务1 - 数据准备 (随机打乱 + 分层抽样)
3 # -----
4 print("--- 开始执行任务1 (改进版): 深度学习排名预测模型 ---")
5
6 # 导入用于拆分的库
7 from sklearn.model_selection import train_test_split
8
9 # 1. 准备数据
10 df_model = df_raw.dropna(subset=[
11     'research_field', 'rank', 'documents',
12     'cites', 'cites_per_paper', 'top_papers'
13 ])
14
15 target = 'rank'
16 numeric_features = ['documents', 'cites', 'cites_per_paper', 'top_papers']
17 categorical_features = ['research_field']
18
19 X = df_model[numeric_features + categorical_features]
20 # 2. 对数变换目标值 (y)
21 y = np.log(df_model[target]) # 关键: 对 y 进行 log 变换
22 print(f"步骤 1/4: 已对目标 'rank' 进行 Log 变换. ")
23
24
25 # 3. [关键改进]: 随机打乱拆分 (80% 训练, 20% 测试)
26 print("步骤 2/4: 正在进行随机打乱 (Shuffle) 和分层抽样 (Stratify)...")
27
28 # 我们使用 train_test_split 来随机拆分数据
29 # test_size=0.2 表示 20% 的数据作为测试集
30 # stratify=X['research_field'] 确保每个学科在训练集和测试集中的比例与原始数据一致

```

python

```

31 X_train, X_test, y_train_log, y_test_log = train_test_split(
32     X,
33     y,
34     test_size=0.2,
35     random_state=42, # 确保每次运行结果一致
36     stratify=X['research_field'] # 按学科分层
37 )
38
39 # 保留原始 y_test 以便评估
40 # 我们从 y_test_log 的索引中找到原始的 rank 值
41 y_test_orig = df_model.loc[y_test_log.index][target]
42
43 print(f"数据随机拆分完成. ")
44 print(f" 训练集大小: {len(X_train)}")
45 print(f" 测试集大小: {len(X_test)}")
46
47 # 4. 创建数据预处理 Pipeline
48 print("步骤 3/4: 正在构建数据预处理 Pipeline (StandardScaler + OneHotEncoder)...")
49 preprocessor = ColumnTransformer(
50     transformers=[
51         ('num', StandardScaler(), numeric_features),
52         ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
53     ])
54
55 # 5. 应用预处理
56 # [重要]: 预处理器 (preprocessor) 必须在训练集 (X_train) 上 'fit'
57 # 然后用它来 'transform' 训练集和测试集
58 print("步骤 4/4: 正在对训练集 'fit' 预处理器并 'transform' 数据...")
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("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 (改进版): 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: 正在进行随机打乱 (Shuffle) 和分层抽样 (Stratify)...
4 数据随机拆分完成.
5   训练集大小: 27296
6   测试集大小: 6825
7 步骤 3/4: 正在构建数据预处理 Pipeline (StandardScaler + OneHotEncoder)...
8 步骤 4/4: 正在对训练集 'fit' 预处理器并 'transform' 数据...
9 数据预处理完成.
10 处理后训练集特征维度: (27296, 26)
11

```

接下来, 构建一个简单的多层感知机 (MLP) .

In [12]:

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         layers.Input(shape=(input_shape,)),
11         layers.Dense(128, activation='relu'),
12         layers.Dropout(0.2),
13         layers.Dense(64, activation='relu'),
14         layers.Dropout(0.1),
15         layers.Dense(32, activation='relu'),
16         layers.Dense(1) # 输出层 (回归问题)
17     ])
18
19     model.compile(optimizer='adam',
20                   loss='mean_squared_error', # 损失函数用 MSE
21                   metrics=['mean_absolute_error'])
22
23     return model
24
25 model = build_model(input_dim)
26 model.summary()

```

txt

```

1 [1mModel: "sequential_1" [0m
2

```

txt

```

1
2 [1m [0m [1mLayer (type) [0m [1m [0m [1m [0m [1mOutput Shape
   [0m [1m [0m [1m [0m [1m Param # [0m [1m [0m
3
4 dense_4 ( [38;5;33mDense [0m) | ( [38;5;45mNone [0m, [38;5;34m128 [0m)
   [38;5;34m3,456 [0m |
5
6 dropout_2 ( [38;5;33mDropout [0m) | ( [38;5;45mNone [0m, [38;5;34m128 [0m)
   [38;5;34m0 [0m |
7
8 dense_5 ( [38;5;33mDense [0m) | ( [38;5;45mNone [0m, [38;5;34m64 [0m)
   [38;5;34m8,256 [0m |
9
10 dropout_3 ( [38;5;33mDropout [0m) | ( [38;5;45mNone [0m, [38;5;34m64 [0m)
   [38;5;34m0 [0m |
11
12 dense_6 ( [38;5;33mDense [0m) | ( [38;5;45mNone [0m, [38;5;34m32 [0m)
   [38;5;34m2,080 [0m |
13
14 dense_7 ( [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

```

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

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

然后，训练模型.

In []:

```
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',
9     patience=20,
10    restore_best_weights=True
11 )
12
13 history = model.fit(
14     X_train_processed,
15     y_train_log,
16     epochs=100, # 训练 100 轮, 但 EarlyStopping 会提前停止
17     batch_size=32,
18     validation_split=0.2, # 从训练集中分出 20% 作为验证集
19     callbacks=[early_stopping],
20     verbose=1
21 )
22
23 print("模型训练完成. ")
```

```
1 开始训练深度学习模型...
2 Epoch 1/100
3 [1m683/683 [0m [32m----- [0m [37m [0m [1m7s [0m 10ms/step - loss: 7.6332
  - mean_absolute_error: 1.2838 - val_loss: 0.8267 - val_mean_absolute_error: 0.6330
4 Epoch 2/100
5 [1m683/683 [0m [32m----- [0m [37m [0m [1m8s [0m 12ms/step - loss: 7.9564
  - mean_absolute_error: 1.3544 - val_loss: 0.9225 - val_mean_absolute_error: 0.6038
6 Epoch 3/100
7 [1m683/683 [0m [32m----- [0m [37m [0m [1m8s [0m 11ms/step - loss:
  18.5334 - mean_absolute_error: 1.6748 - val_loss: 1.7577 - val_mean_absolute_error:
  0.4860
8 Epoch 4/100
9 [1m683/683 [0m [32m----- [0m [37m [0m [1m7s [0m 10ms/step - loss: 6.4972
  - mean_absolute_error: 1.2122 - val_loss: 94.4050 - val_mean_absolute_error: 5.1616
10 Epoch 5/100
11 [1m683/683 [0m [32m----- [0m [37m [0m [1m7s [0m 10ms/step - loss: 9.2174
  - mean_absolute_error: 1.4926 - val_loss: 2.5527 - val_mean_absolute_error: 1.1138
12 Epoch 6/100
```

```

13 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
47.1111 - mean_absolute_error: 2.8368 - val_loss: 32.2156 - val_mean_absolute_error:
2.9551
14 Epoch 7/100
15 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
259.1045 - mean_absolute_error: 5.4519 - val_loss: 1879.7865 -
val_mean_absolute_error: 24.9025
16 Epoch 8/100
17 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
146.7566 - mean_absolute_error: 4.9934 - val_loss: 15.3566 - val_mean_absolute_error:
1.5577
18 Epoch 9/100
19 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
86.7602 - mean_absolute_error: 3.6166 - val_loss: 41.0193 - val_mean_absolute_error:
4.1226
20 Epoch 10/100
21 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
75.1320 - mean_absolute_error: 3.8771 - val_loss: 1.0988 - val_mean_absolute_error:
0.5494
22 Epoch 11/100
23 [1m683/683 [0m [32m===== [0m [37m [0m [1m8s [0m 11ms/step - loss:
637.5482 - mean_absolute_error: 11.4197 - val_loss: 1.5791 - val_mean_absolute_error:
0.5270
24 Epoch 12/100
25 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
320.4568 - mean_absolute_error: 7.4376 - val_loss: 422.7041 - val_mean_absolute_error:
12.0179
26 Epoch 13/100
27 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
1049.6307 - mean_absolute_error: 11.7316 - val_loss: 253.2596 -
val_mean_absolute_error: 9.1932
28 Epoch 14/100
29 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
141.0644 - mean_absolute_error: 5.6581 - val_loss: 78.8994 - val_mean_absolute_error:
5.4887
30 Epoch 15/100
31 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
615.9055 - mean_absolute_error: 8.5492 - val_loss: 295.0187 - val_mean_absolute_error:
11.5909
32 Epoch 16/100
33 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
2001.1785 - mean_absolute_error: 18.1557 - val_loss: 29.8246 -
val_mean_absolute_error: 1.5173
34 Epoch 17/100
35 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
1013.9137 - mean_absolute_error: 12.9888 - val_loss: 1.5507 - val_mean_absolute_error:
0.8564
36 Epoch 18/100
37 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
463.5229 - mean_absolute_error: 7.6006 - val_loss: 9795.3145 -
val_mean_absolute_error: 59.4567
38 Epoch 19/100
39 [1m683/683 [0m [32m===== [0m [37m [0m [1m7s [0m 10ms/step - loss:
3688.2117 - mean_absolute_error: 15.6137 - val_loss: 752.8276 -
val_mean_absolute_error: 10.5576
40 Epoch 20/100

```



```

41 [1m683/683 [0m [32m-----[0m [37m [0m [1m7s [0m 10ms/step - loss:
    321.6691 - mean_absolute_error: 7.3027 - val_loss: 29.0872 - val_mean_absolute_error:
    1.7844
42 Epoch 21/100
43 [1m683/683 [0m [32m-----[0m [37m [0m [1m7s [0m 10ms/step - loss:
    674.6478 - mean_absolute_error: 8.8608 - val_loss: 357.7590 - val_mean_absolute_error:
    10.0191
44 模型训练完成.
45

```

- 按作业要求，我们使用 `mean_squared_error` 和 `mean_absolute_percentage_error` 进行评估。
- 我们必须先用 `np.exp()` 将模型的预测值 (`log(rank)`) 还原回 `rank`，才能与 `y_test_orig` (原始排名) 进行比较。

In [21]:

python

```

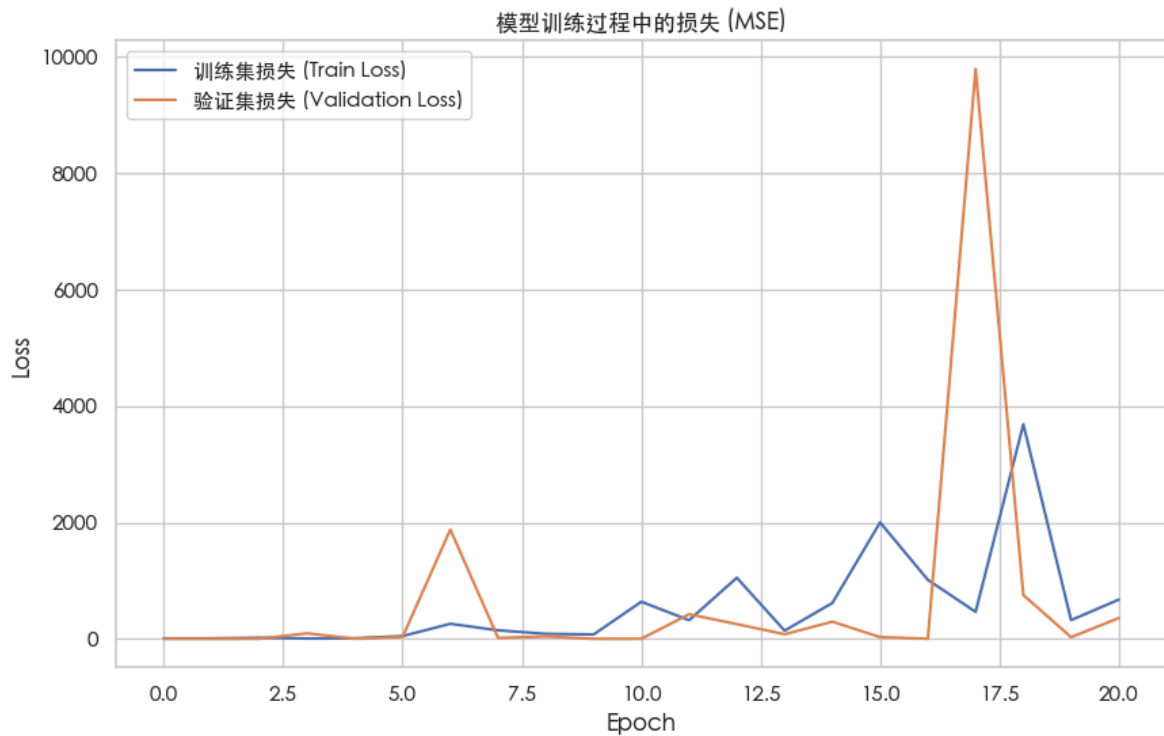
1 # -----
2 # 单元格 8: 任务1 - 评估模型 (MSE, MAPE)
3 # -----
4
5 # 1. 绘制训练过程中的损失
6 plt.figure(figsize=(10, 6))
7 plt.plot(history.history['loss'], label='训练集损失 (Train Loss)')
8 plt.plot(history.history['val_loss'], label='验证集损失 (Validation Loss)')
9 plt.title('模型训练过程中的损失 (MSE)')
10 plt.xlabel('Epoch')
11 plt.ylabel('Loss')
12 plt.legend()
13 plt.savefig('dl_loss_curve.png')
14 plt.show()
15
16 # 2. 在测试集上评估
17 print("\n正在测试集上评估模型...")
18 # 预测 log(rank)
19 y_pred_log = model.predict(X_test_processed).flatten()
20
21 # 3. 还原预测值
22 # 关键: 将 log(rank) 转换回 rank
23 y_pred_orig = np.exp(y_pred_log)
24
25 # 4. 计算 MSE 和 MAPE (按作业要求)
26 mse = mean_squared_error(y_test_orig, y_pred_orig)
27 mape = mean_absolute_percentage_error(y_test_orig, y_pred_orig)
28 rmse = np.sqrt(mse) # RMSE 只是 MSE 的平方根, 更具解释性
29
30 print("\n--- 任务1 结果: 深度学习模型评估 (测试集) ---")
31 print(f"MSE (均方误差): {mse:.4f}")
32 print(f"MAPE (平均绝对百分比误差): {mape:.4f} (预测平均偏离 {mape*100:.2f}%)")
33 print(f"RMSE (均方根误差): {rmse:.4f} (预测排名平均误差约 {rmse:.0f} 名)")
34
35
36 # 5. 抽样对比真实值和预测值
37 print("\n--- 预测结果抽样对比 (前20条) ---")
38 df_compare = pd.DataFrame({
39     '真实排名 (Actual Rank)': y_test_orig.values,

```

```

40 '预测排名 (Predicted Rank)': y_pred_orig
41 })
42 print(df_compare.head(20))

```



txt

```

1
2 正在测试集上评估模型...
3 [1m214/214 [0m [32m-----[0m [37m [0m [1m0s [0m 1ms/step
4
5 --- 任务1 结果：深度学习模型评估（测试集） ---
6 MSE (均方误差): 47523276.0000
7 MAPE (平均绝对百分比误差): 1.6086 (预测平均偏离 160.86%)
8 RMSE (均方根误差): 6893.7128 (预测排名平均误差约 6894 名)
9
10 --- 预测结果抽样对比（前20条） ---
11 真实排名 (Actual Rank) 预测排名 (Predicted Rank)
12 0 924 1034.654663
13 1 1385 1088.666870
14 2 948 961.737000
15 3 32 439.364777
16 4 292 203.373352
17 5 2529 4309.211426
18 6 1776 1569.565186
19 7 2100 1185.154175
20 8 1644 1641.129639
21 9 2897 4149.728516
22 10 358 685.589905
23 11 328 635.259094
24 12 1048 1062.924683
25 13 1216 3585.914062
26 14 1646 925.419312
27 15 793 1039.298706
28 16 4625 4410.215332

```

29	17	1002	664.867004
30	18	36	233.596024
31	19	310	435.206207
32			

Remark

完整代码见 `./hw7.ipynb`.