

# 实验七：深度学习算法升级

## 1. 作业要求

上次作业中，我利用了深度学习方法，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

本次作业中，我需要升级这个模型

## 2. 上次作业回顾

原数据的第一行和最后一行有冗余信息，可能会干扰表格的读取

第一行是表头，应该去掉：

A	B	C	D	E	F	G	H	I	J	K	L
1	Indicators	Results List: Institutions	Filter Results By: ResearchFields	Filter Value(s): AGRICULTURAL SCIENCES	Show: Top						
2		Institution Countries/ Web of Sc	Cites	Cites/Paper	Top Papers						
3	1	CHINESE / CHINA M/	15661	332254	21.22	251					
4	2	CHINESE / CHINA M/	12222	223855	18.32	198					
5	3	UNITED S' USA	12564	220644	17.56	105					
6	4	CHINA AG CHINA M/	10052	207779	20.67	166					
7	5	INRAE FRANCE	9314	187838	20.17	118					
8	6	JIANGNAN CHINA M/	8027	181325	22.59	166					
9	7	NORTHW CHINA M/	9024	180634	20.02	162					
10	8	WAGENING NETHERL/	6819	172734	25.33	152					
11	9	CONSEJO SPAIN	7656	167593	21.89	94					
12	10	NANJING CHINA M/	6480	137921	21.28	99					
13	11	UNIVERSIT USA	6175	129479	20.97	110					
14	12	UNIVERSIT CHINA M/	5950	127682	21.46	95					
15	13	ZHEJIANG CHINA M/	5313	125715	23.66	148					
16	14	EGYPTIAN EGYPT	7861	117405	14.94	103					
17	15	INDIAN C INDIA	12753	111262	8.72	66					
18	16	SOUTH CI CHINA M/	3524	108282	30.73	117					
19	17	CGIAR FRANCE	5262	107193	20.37	69					
20	18	HUAZHOU CHINA M/	4855	97482	20.08	68					
21	19	UNIVERSII BRAZIL	6169	95242	15.44	35					
22	20	JIANGSU I CHINA M/	4084	89040	21.8	176					
23	21	STATE UN USA	5172	86527	16.73	62					
24	22	CENTRE N FRANCE	4113	81348	19.78	60					
25	23	CORNELL USA	3551	78744	22.18	59					

最后一行是 Copyright ?2025 Clarivate，也应该去掉：

A	B	C	D	E	F	G	H	I	J	K	L
1359	1356 MAE FAH	THAILAND	236	3627	15.37	2					
1360	1358 NAMIK KE	TURKIYE	352	3625	10.3	1					
1361	1359 GANGNEU	SOUTH KOREA	268	3621	13.51	1					
1362	1360 TAMPERE	FINLAND	234	3619	15.47	0					
1363	1361 HOSPITAL	CANADA	202	3618	17.91	2					
1364	1362 BRAUNSC	GERMANY	188	3614	19.22	1					
1365	1362 ICAR - INI	INDIA	393	3614	9.2	6					
1366	1364 HUNGARI	HUNGARY	180	3605	20.03	2					
1367	1365 LIAONING	CHINA M	333	3604	10.82	0					
1368	1366 DAIRYNZ	NEW ZEALAND	239	3599	15.06	0					
1369	1367 SMITHSOI	USA	209	3598	17.22	0					
1370	1367 BAYER CR	GERMANY	250	3598	14.39	2					
1371	1369 GENERAL	USA	107	3595	33.6	4					
1372	1370 STEFAN C	ROMANIA	161	3591	22.3	5					
1373	1371 KAGOSHII	JAPAN	306	3588	11.73	1					
1374	1372 UNIVERSIT	USA	127	3580	28.19	4					
1375	1373 HOPITAL	FRANCE	123	3573	29.05	3					
1376	1373 HEINRICH	GERMANY	136	3573	26.27	4					
1377	1375 UNIVERSIT	SOUTH AFRICA	184	3572	19.41	6					
1378	1376 BANGLAD	BANGLAD	170	3567	20.98	2					
1379	1376 NEAR EAS	TURKIYE	204	3567	17.49	2					
1380	1378 GRIFFITH	AUSTRALIA	205	3562	17.38	2					
1381	1378 UMM AL-	SAUDI AR.	308	3562	11.56	3					
1382	1380 SAINT LOI	USA	108	3561	32.97	4					
1383	1381 VIT VELLC	INDIA	144	3558	24.71	2					
1384	Copyright ©2025 Clarivate										
1385											

在 Powershell 执行下面命令：

```
$src="D:\ESI\download"; $dst="D:\ESI\clean"; New-Item -ItemType Directory -Force -Path $dst | Out-Null
Get-ChildItem $src -Filter *.csv | ForEach-Object {
    $lines = Get-Content $_.FullName
    Get-ChildItem $src -Filter *.csv
    $lines[1..($lines.Count-2)] | Set-Content -Encoding UTF8 (Join-Path $dst $_.Name)
}
```

解释一下：

1. src 是原始文件路径，dst 是目标文件路径
2. New-Item 用于创建一个新目录，如果 D:\ESI\clean 不存在，会自动创建
3. Get-ChildItem \$src -Filter \*.csv 会列出 src 中所有扩展名为 .csv 的文件
4. ForEach-Object 遍历每个文件
5. \$lines = Get-Content \$\_.FullName 会把整个文件的内容读成一个字符串数组，每一行是一个元素
6. \$lines[1..(\$lines.Count-2)] 是数组切片语法，表示从第二行到倒数第二行，删除文件的第一行和最后一行
7. Set-Content 把上一步得到的“去掉第一行的内容”写入到新的文件中

可以看到，第一行已经成功去掉了：

A	B	C	D	E	F	G	H	I	J	K	L
	Institution	Countries	Web of Sc	Cites	Cites/Paper	Top Papers					
1	1 CHINESE A	CHINA M/	15661	332254	21.22	251					
2	2 CHINESE A	CHINA M/	12222	223855	18.32	198					
3	3 UNITED S'	USA	12564	220644	17.56	105					
4	4 CHINA AG	CHINA M/	10052	207779	20.67	166					
5	5 INRAE	FRANCE	9314	187838	20.17	118					
6	6 JIANGNAN	CHINA M/	8027	181325	22.59	166					
7	7 NORTHW	CHINA M/	9024	180634	20.02	162					
8	8 WAGENIN	NETHERL/	6819	172734	25.33	152					
9	9 CONSEJO	SPAIN	7656	167593	21.89	94					
10	10 NANJING	CHINA M/	6480	137921	21.28	99					
11	11 UNIVERSI	USA	6175	129479	20.97	110					
12	12 UNIVERSI	CHINA M/	5950	127682	21.46	95					
13	13 ZHEJIANG	CHINA M/	5313	125715	23.66	148					
14	14 EGYPTIAN	EGYPT	7861	117405	14.94	103					
15	15 INDIAN C	INDIA	12753	111262	8.72	66					
16	16 SOUTH C	CHINA M/	3524	108282	30.73	117					
17	17 CGIAR	FRANCE	5262	107193	20.37	69					
18	18 HUAZHOU	CHINA M/	4855	97482	20.08	68					
19	19 UNIVERSI	BRAZIL	6169	95242	15.44	35					
20	20 JIANGSU I	CHINA M/	4084	89040	21.8	176					
21	21 STATE UN	USA	5172	86527	16.73	62					
22	22 CENTRE N	FRANCE	4113	81348	19.78	60					
23	23 CORNELL	USA	3551	78744	22.18	59					
24	24 UNIVERSI	USA	3908	78346	20.05	64					

最后一行也成功去掉了：

A	B	C	D	E	F	G	H	I	J	K	L
1364	1364 HUNGARI	HUNGARY									
1366	1365 LIAONING	CHINA M/	333	3604	10.82	0					
1367	1366 DAIRYNZ	NEW ZEAI	239	3599	15.06	0					
1368	1367 SMITHSOI	USA	209	3598	17.22	0					
1369	1367 BAYER CR	GERMANY	250	3598	14.39	2					
1370	1369 GENERAL	USA	107	3595	33.6	4					
1371	1370 STEFAN C	ROMANIA	161	3591	22.3	5					
1372	1371 KAGOSHII	JAPAN	306	3588	11.73	1					
1373	1372 UNIVERSI	USA	127	3580	28.19	4					
1374	1373 HOPITAL I	FRANCE	123	3573	29.05	3					
1375	1373 HEINRICH	GERMANY	136	3573	26.27	4					
1376	1375 UNIVERSI	SOUTH AF	184	3572	19.41	6					
1377	1376 BANGLAD	BANGLAD	170	3567	20.98	2					
1378	1376 NEAR EAST	TURKIYE	204	3567	17.49	2					
1379	1378 GRIFFITH I	AUSTRALI	205	3562	17.38	2					
1380	1378 UMM AL-	SAUDI AR	308	3562	11.56	3					
1381	1380 SAINT LOI	USA	108	3561	32.97	4					
1382	1381 VIT VELLC	INDIA	144	3558	24.71	2					
1383											
1384											
1385											
1386											
1387											
1388											
1389											
1390											
1391											

我整理了一下邱学长提到的深度学习步骤：

- 定义数据集
- 定义模型
- 定义训练过程
- 训练多少轮次
- 参数冻结
- 在测试集上跑一遍
- 深层次MMP，每层加一些归一化、正则化、非线性组合

首先，题目要求我们对各学科做一个排名模型，预测排名位置

电脑通过学习一些国家的 (Web of Science Documents, Cites, Cites/Paper, Cites/Paper, 排名) ，

再给其他国家的 (Web of Science Documents, Cites, Cites/Paper, Cites/Paper) , 可以推测出排名  
我们来明确一下输入和输出:

输入: 多个 CSV, 每个 CSV 代表一个学科, 第一列就是排名, 这是通过人工观察发现的

输出: 对于每个学科, 在测试值上跑预测模型, 得出的真实值对比预测值的表, 以及各项指标的总表

我们要明确一点, 每个学科的排名预测模型都应该是**独立的**

因此, 要循环遍历每个csv文件, **写一个循环**

在循环中, 依次遍历存放多个 CSV 文件路径的列表, 并同时获取每个文件的编号 (从 1 开始) 和路径

```
for i, csv_path in enumerate(csvs, 1):
```

然后, 取出当前表格最左侧的第一列列名, 并打印当前处理进度, 包括文件编号、总文件数、学科文件名以及该排名列的列名, 用来提示程序正在处理哪个学科的排名数据

```
label_col = df.columns[0]
print(f"\n== {i}/{len(csvs)} 学科: {csv_path.name} | 排名列: {label_col} ==")
```

前面提到过了, 我们通过一个学校的:

```
x = (Web of Science Documents, Cites, Cites/Paper, Cites/Paper)
```

来预测它的排名, 设为 Y

将表格中最左侧的排名列转成数值型, 然后把剩下的所有列作为特征矩阵 X, 并通过 to\_numeric\_df() 函数将这些特征列全部转换为数值类型, 以便后续模型训练或聚类分析使用

```
y_raw = pd.to_numeric(df[label_col], errors="coerce").astype(float)
x = df.drop(columns=[label_col])
x = to_numeric_df(x)
```

不要忘了过滤掉在排名列中存在缺失值的行, 只保留排名有效的样本, 使得 X 和 Y 严格对应且都是完整的数值型数据, 从而确保后续训练或聚类过程不会受到缺失数据干扰

```
mask = y_raw.notna()
x = x.loc[mask].copy()
y_raw = y_raw.loc[mask].astype(float)
```

然后, 划分数据集

课堂上, 老师讲解了测试集、验证集、训练集

邱学长也提到了合适的比例: 60%, 20%, 20%

我们还要计算当前样本总数 N, 并生成一个从 0 到 N-1 的索引数组 idx\_all, 然后调用 split\_60\_20\_20\_idx() 函数, 基于给定随机种子 SEED 进行数据集的划分

```
N = len(row_ids)
idx_all = np.arange(N)
idx_tr, idx_va, idx_te = split_60_20_20_idx(idx_all, seed=SEED)
```

然后，我们要 scale 一下特征，因为特征的大小范围很不统一

创建一个标准化器 StandardScaler(), 以训练集数据为基准进行特征标准化处理——即让每个特征的均值为 0、标准差为 1

```
scaler = StandardScaler()
X_train = scaler.fit_transform(X_np[idx_tr])
X_val = scaler.transform(X_np[idx_va])
X_test = scaler.transform(X_np[idx_te])
```

排名也要 scale 一下

经过下面的处理后，模型预测输出也是 [0,1] 之间的数值，更易收敛、数值稳定，并可在评估阶段反归一化回真实排名

```
y_min = float(np.min(y_np))
y_max_v = float(np.max(y_np))
y_max = y_max_v if y_max_v != y_min else (y_min + 1.0)
def norm(y): return (y - y_min) / (y_max - y_min)

y_train = norm(y_np[idx_tr])
y_val = norm(y_np[idx_va])
y_test = norm(y_np[idx_te])
```

然后，调用自定义函数 `train_one()` 来训练一个深度学习模型，并返回在测试集上的预测结果。

```
_ , predict_rank = train_one(  
    x_train, y_train, x_val, y_val, x_test, y_test,  
    y_min=y_min, y_max=y_max, device=device  
)
```

先插叙一下，深度学习模型的**内部**是怎么工作的？

首先，动态选择合适的归一化方式，构建深度学习模型，并定义一个稳健的损失函数

```
norm_kind = "batch" if len(X_train) >= 2 else "layer"
model = DeepRankMLP(X_train.shape[1], norm_kind=norm_kind).to(device)
huber = nn.SmoothL1Loss(beta=1.0) # robust to outliers
```

训练分为两个阶段，full training 和 freeze and fine-tune

## 1. full training

先用 AdamW 优化器，并配上 StepLR 学习率调度器（每 40 个 epoch 将学习率乘以 0.5），随后把训练集与验证集的 NumPy 特征/标签 转成 PyTorch 张量并封装为 TensorDataset，为后续用 DataLoader 批量迭代训练与评估做准备

```
opt = torch.optim.AdamW(model.parameters(), lr=LR1, weight_decay=WEIGHT_DECAY)
sched = torch.optim.lr_scheduler.StepLR(opt, step_size=40, gamma=0.5) tr_ds =
TensorDataset(torch.from_numpy(X_train).float(),
torch.from_numpy(y_train_norm).float()) va_ds =
TensorDataset(torch.from_numpy(X_val).float(),
torch.from_numpy(y_val_norm).float())
```

构建 PyTorch 的 DataLoader，也就是让模型能够按批次取出训练数据和验证数据的关键部分

```
bs_eff = min(BATCH_SIZE, len(tr_ds))
drop_last = (len(tr_ds) >= 2)
tr_dl = DataLoader(tr_ds, batch_size=bs_eff, shuffle=True, drop_last=drop_last)
va_dl = DataLoader(va_ds, batch_size=min(BATCH_SIZE, len(va_ds)), shuffle=False)
```

把模型输出的归一化预测结果（0~1之间）反变换回原始排名范围

因为我们前面把排名 scale 到了 0~1之间，可是题目想要的肯定是预测原始排名

```
def denorm(y_hat_norm: np.ndarray) -> np.ndarray:
    return y_min + np.clip(y_hat_norm, 0.0, 1.0) * (y_max - y_min)
```

下面实现一个完整的训练轮次循环

每轮训练都会遍历整个训练集，将模型设为训练模式，再逐批取出数据送入网络进行前向传播，计算预测值与真实值的 Huber 损失

接着清空梯度、反向传播、裁剪梯度防止爆炸，最后用 AdamW 优化器更新参数

每个 epoch 结束后，学习率调度器按计划降低学习率，以帮助模型在后期更平稳地收敛

```
nonlocal best_state, best_val_mse_rank
for ep in range(1, epochs + 1):
    if len(tr_dl) == 0:
        break
    model.train()
    for xb, yb in tr_dl:
        xb, yb = xb.to(device), yb.to(device)
        pred = model(xb)
        loss = huber(pred, yb)
        opt.zero_grad()
        loss.backward()
        nn.utils.clip_grad_norm_(model.parameters(), 5.0)
        opt.step()
    sched.step()
```

先切换到 eval 模式并关闭梯度计算，然后用验证集逐批预测，得到所有预测值和真实值

接着将它们反归一化为原始排名，用 MSE 计算误差

若当前误差小于历史最优值，就更新 best\_val\_mse\_rank 并保存当前模型参数

```
model.eval()
with torch.no_grad():
    yp, yt = [], []
    for xb, yb in va_dl:
        xb = xb.to(device)
        yp.append(model(xb).cpu().numpy())
        yt.append(yb.cpu().numpy())
    if len(yp) == 0:
        continue
    y_pred_norm = np.concatenate(yp).reshape(-1)
    y_true_norm = np.concatenate(yt).reshape(-1)
```

```

y_pred_rank = denorm(y_pred_norm)
y_true_rank = denorm(y_true_norm)
v_mse_rank = mse_np(y_true_rank, y_pred_rank)
if v_mse_rank < best_val_mse_rank:
    best_val_mse_rank = v_mse_rank
    best_state = {k: v.clone() for k, v in model.state_dict().items()}

```

## 2. freeze and fine-tune

接下来是第二部分，冻结参数训练

首先调用冻结模型前半部分的参数，只让后半部分参与训练

接着，用 AdamW 优化器重新初始化，只更新未冻结部分的参数，并设置一个更小的学习率 LR2 与相同的 weight\_decay 来避免过拟合

然后建立新的 StepLR 调度器，每 20 轮将学习率减半

最后调用 run\_epochs(EPOCHS\_2) 进行多轮训练

```

freeze_first_half(model, True)
opt = torch.optim.AdamW(filter(lambda p: p.requires_grad, model.parameters()),
lr=LR2, weight_decay=WEIGHT_DECAY)
sched = torch.optim.lr_scheduler.StepLR(opt, step_size=20, gamma=0.5)
run_epochs(EPOCHS_2)

```

最后是预测排名阶段

首先，将模型切换到评估模式，关闭 dropout 和 batch norm 的随机行为

关闭梯度计算，避免浪费显存和计算资源

把输入的 NumPy 数组 X\_np 转换为 PyTorch 张量并移动到设备上

调用模型获得预测的归一化输出 yhat\_norm，再转回 NumPy 格式

通过 denorm() 函数反归一化，将结果映射回真实的排名区间

```

def predict_rank(x_np: np.ndarray) -> np.ndarray:
    model.eval()
    with torch.no_grad():
        x_t = torch.from_numpy(x_np).float().to(device)
        yhat_norm = model(x_t).cpu().numpy().reshape(-1)
    return denorm(yhat_norm)

```

最后，利用训练好的模型在测试集上进行预测，并计算预测结果与真实排名之间的多种评估指标，用于衡量模型性能

```

y_pred_te = predict_rank(X_test)
y_true_te = y_np[idx_te]

metrics = compute_metrics(y_true_te, y_pred_te)

```

关于输出部分，我们前面提到过了：

对于每个学科，在**测试值**上跑预测模型，得出的**真实值对比预测值**的表，以及各项指标的总表

首先是各项指标的输出，需要的数据都在metrics字典中

依次输出即可

```
print("[测试] " + " ".join([
    f"MSE={metrics['MSE']:.4f}",
    f"RMSE={metrics['RMSE']:.2f}",
    f"MAE={metrics['MAE']:.2f}",
    f"MAPE={metrics['MAPE']:.2f}%",
    f"MedAE={metrics['MedianAE']:.2f}",
    f"R2={metrics['R2']:.3f}",
    f"nRMSE={metrics['nRMSE']:.3f}",
    f"Spearman={metrics['Spearman']:.3f}",
    f"Kendall={metrics['Kendall']:.3f}",
]))
```

然后是对比表

把模型在测试集上某所大学的真实排名和预测排名保存为一个独立的结果表

```
pred_df = pd.DataFrame({
    "row_id": row_ids[idx_te],
    "true_rank": y_true_te,
    "pred_rank": y_pred_te,
})
```

生成每个学科预测结果的 CSV 文件名并保存预测结果表格

```
stem = csv_path.stem
safe = "".join([ch if ch.isalnum() or ch in ("-", "_") else "_" for ch in stem])
pred_path = (Path(__file__).resolve().parent / f"predictions_{safe}.csv")
pred_df.to_csv(pred_path, index=False)
print(f"已保存测试对比表: {pred_path.resolve()}")
```

但是我想把所有学科的各项指标汇总到一张表中，以便**纵向对比**

把每个学科对应的模型评估指标保存到一个总表列表中，以便最后统一生成总表

```
row = {"subject_csv": csv_path.name}
row.update(metrics)
rows.append(row)
```

最后，在所有学科模型的预测模型评测结果汇总后，生成一个综合评价表 deep\_learning.csv，并按 MSE 升序排序显示和保存

```
if rows:
    out = pd.DataFrame(rows).sort_values("MSE")
    out_path = (Path(__file__).resolve().parent / "deep_learning.csv")
    out.to_csv(out_path, index=False)
    print("\n==== 总结（按 MSE 升序） ====")
    print(out.to_string(index=False))
    print(f"\n结果已保存到: {out_path.resolve()}")
```

以上的完整代码，请查看 deep\_learning.py

各学科排名模型在测试集上的预测值与真实值对比表格，都在附件中：

名称	修改日期	类型	大小
predictions_BIOLOGY__BIOCHEMISTRY.csv	30/10/2025 10:20	Microsoft Excel ...	1 KB
predictions_CHEMISTRY.csv	30/10/2025 10:20	Microsoft Excel ...	10 KB
predictions_CLINICAL_MEDICINE.csv	30/10/2025 10:21	Microsoft Excel ...	30 KB
predictions_COMPUTER_SCIENCE.csv	30/10/2025 10:21	Microsoft Excel ...	4 KB
predictions_ECONOMICS__BUSINES..	30/10/2025 10:22	Microsoft Excel ...	3 KB
predictions_ENGINEERING.csv	30/10/2025 10:22	Microsoft Excel ...	12 KB
predictions_ENVIRONMENT_ECOLOG...	30/10/2025 10:22	Microsoft Excel ...	9 KB
predictions_GEOSCIENCES.csv	30/10/2025 10:22	Microsoft Excel ...	5 KB
predictions_IMMUNOLOGY.csv	30/10/2025 10:22	Microsoft Excel ...	5 KB
predictions_MATERIALS_SCIENCE.csv	30/10/2025 10:22	Microsoft Excel ...	7 KB
predictions_MATHEMATICS.csv	30/10/2025 10:22	Microsoft Excel ...	2 KB
predictions_MICROBIOLOGY.csv	30/10/2025 10:23	Microsoft Excel ...	4 KB
predictions_MOLECULAR_BIOLOGY ...	30/10/2025 10:23	Microsoft Excel ...	5 KB
predictions_MULTIDISCIPLINARY.csv	30/10/2025 10:23	Microsoft Excel 逗号分隔值文件	3 KB
predictions_NEUROSCIENCE__BEHA...	30/10/2025 10:23	Microsoft Excel ...	6 KB
predictions_PHARMACOLOGY__TO...	30/10/2025 10:23	Microsoft Excel ...	6 KB
predictions_PHYSICS.csv	30/10/2025 10:23	Microsoft Excel ...	4 KB
predictions_PLANT__ANIMAL_SCIE...	30/10/2025 10:23	Microsoft Excel ...	9 KB
predictions_PSYCHIATRY_PSYCHOLO...	30/10/2025 10:23	Microsoft Excel ...	5 KB
predictions_SOCIAL SCIENCES__GEN...	30/10/2025 10:24	Microsoft Excel ...	11 KB
predictions_SPACE_SCIENCE.csv	30/10/2025 10:24	Microsoft Excel ...	1 KB
8_build.csv	30/10/2025 10:24	Microsoft Excel ...	5 KB

我这里随机选一个学科展示一下：

A	B	C	D	E	F	G	H	I	J
1	row_id	true_rank	pred_rank						
2	581	582	580.8528						
3	175	176	210.1509						
4	548	549	538.756						
5	654	655	642.4933						
6	1926	1927	1934.248						
7	1244	1245	1256.446						
8	290	291	295.6851						
9	221	222	268.3793						
10	1461	1461	1466.63						
11	433	434	452.868						
12	121	122	190.9796						
13	891	892	879.5628						
14	286	287	314.0567						
15	867	868	851.5161						
16	347	348	367.6648						
17	1962	1963	1974.138						
18	2731	2731	2611.857						
19	787	788	669.9355						
20	1798	1799	1808.691						
21	188	189	216.5818						
22	1446	1447	1464.342						
23	1910	1911	1919.203						
24	303	304	291.6348						
25	1621	1622	1635.82						
26	2164	2165	2184.776						
27	1933	1934	1946.557						
28	409	410	437.9659						
29	1716	1717	1724.965						

可以看到，预测值与真实值还是比较接近的，说明模型预测效果好

预测的排名值是一个小数，而不是整数，虽然排名只能是整数

如果需要整数，把预测的排名值四舍五入即可，或者向偶数取整

不过我觉得小数也是有意义的，因为它本身只是一个估计

比如一个大学预测排名为 2.5，那可以认为它的排名大概是第二、第三

各学科排名模型在**测试集上的各项指标**如下：

subject_csv	MSE	RMSE	MAE	MAPE	MedianAE	R2	nRMSE	Spearman	Kendall
MICROBIOLOGY.csv	252.269248	15.882986	10.887834	9.391180	8.071838	0.995208	0.019954	0.999330	0.984664
SPACE SCIENCE.csv	341.910472	18.490821	12.745177	17.596886	9.387699	0.916745	0.081818	0.964611	0.867021
GEOSCIENCES.csv	393.698578	19.841839	14.370080	7.196743	11.384811	0.996415	0.017076	0.999041	0.980594
IMMUNOLOGY.csv	508.733315	22.555117	14.377933	9.384470	9.383499	0.995366	0.019278	0.998747	0.982618
MOLECULAR BIOLOGY & GENETICS.csv	617.998418	24.859574	16.262368	9.345132	9.079315	0.994337	0.021449	0.999509	0.988023
AGRICULTURAL SCIENCES.csv	630.600022	25.111751	15.397359	5.749851	8.184601	0.995703	0.018237	0.999727	0.988554
NEUROSCIENCE & BEHAVIOR.csv	805.005098	28.372612	23.433166	8.089908	23.185997	0.994044	0.022114	0.999609	0.989159
MATERIALS SCIENCE.csv	812.457454	28.503639	19.236143	9.710216	11.946630	0.996106	0.018295	0.999076	0.979174
PHYSICS.csv	907.911353	30.131567	21.446334	9.820175	15.024963	0.989004	0.030528	0.997779	0.972158
PHARMACOLOGY & TOXICOLOGY.csv	958.344377	30.957138	22.375154	6.753941	17.531281	0.993706	0.022449	0.999778	0.991520
PSYCHIATRY PSYCHOLOGY.csv	1191.735976	34.521529	17.262220	7.539682	10.504194	0.988529	0.030362	0.995185	0.968633
ENVIRONMENT ECOLOGY.csv	1425.751402	37.759123	28.113645	7.195055	24.663239	0.995981	0.018482	0.999325	0.987613
COMPUTER SCIENCE.csv	1430.191239	37.817869	20.458793	11.278869	12.577209	0.977369	0.044180	0.992381	0.978220
ENGINEERING.csv	1504.758002	38.791210	24.294541	9.406422	12.861816	0.997717	0.013969	0.999850	0.993642
MULTIDISCIPLINARY.csv	1763.283374	41.991468	23.607978	56.734116	8.507751	0.446612	0.203842	0.711346	0.610994
MATHEMATICS.csv	2138.738504	46.246497	27.812698	73.346139	15.793304	0.840698	0.118581	0.928669	0.830385
BIOLOGY & BIOCHEMISTRY.csv	2251.553616	47.450539	34.840483	13.596073	23.239258	0.989958	0.029200	0.999768	0.992530
CLINICAL MEDICINE.csv	2467.755633	49.676510	40.941861	3.610824	36.800293	0.999367	0.007373	0.999911	0.994969
ECONOMICS & BUSINESS.csv	2705.859523	52.017877	24.491297	30.210407	12.215485	0.904894	0.059574	0.953637	0.892438
CHEMISTRY.csv	3585.305891	59.877424	44.487254	16.364376	32.726349	0.990791	0.028098	0.999752	0.994537
PLANT & ANIMAL SCIENCE.csv	3697.386316	60.806137	48.659946	21.578761	38.319824	0.988808	0.031311	0.999848	0.992538
SOCIAL SCIENCES, GENERAL.csv	4389.281090	66.251650	49.070618	12.624892	41.315399	0.991133	0.027674	0.998471	0.989669

其中的 nRMSE 是 RMSE 按数值范围归一化之后的，很有参考价值：

nRMSE < 0.10：很好

0.10–0.20：可用

而各学科排名模型在**测试集上的 nRMSE** 都小于 0.10，说明模型预测效果**很好**

各项指标的总表也在附件中，叫 deep\_learning.csv

### 3. 升级上次的模型

#### (0) 宏观思路

上次作业中，我已经采用了老师课堂强调的主流算法，因此本次以小步快跑的思路做稳健微调，不追求大改结构，重点围绕可复现、数值稳定与推理期后处理三条线优化，包括统一播种、早停与学习率调度、验证集线性标定、预测范围裁剪与整数化输出，同时控制计算预算，保证效果提升而训练时长基本不变

下面分块进行代码讲解，完整代码见 deep\_learning\_upgrade.py

#### (1) 导入 os 与 random

在导入区补充 os 与 random，为后续统一管理随机性和轻量环境控制预留入口。它不改变任何训练与推理逻辑，只提供可复现实验所需的最小工程支点，避免因平台差异或默认随机源导致结果漂移。

```
import os, random
```

## (2) 统一播种与数值稳定

新增播种函数一次性设置 random numpy torch 的种子，并在可用时开启卷积基准与高精度矩阵乘法，全程用异常保护保证兼容。这样数据划分权重初始化和丢弃层掩码可复现，训练曲线与评估更稳定。

```
def _seed_all(seed: int = 42):
    try:
        random.seed(seed)
    except Exception:
        pass
    try:
        np.random.seed(seed)
    except Exception:
        pass
    try:
        torch.manual_seed(seed)
        if torch.cuda.is_available():
            torch.cuda.manual_seed_all(seed)
    except:
        torch.backends.cudnn.benchmark = True
    except Exception:
        pass
    try:
        torch.set_float32_matmul_precision("high")
    except Exception:
        pass
    except Exception:
        pass
```

## (3) 轻量测试时丢弃平均

在推理阶段进行少量次前向，临时启用 Dropout 取多份独立输出后求均值，几乎不增加时间成本。该做法可抑制单次前向的方差与偶然抖动，常带来小而稳定的误差下降，不影响训练权重。

```
def predict_rank_tta(x_np: np.ndarray, n_pass: int = 5) -> np.ndarray:
    outs = []
    with torch.no_grad():
        X_t = torch.from_numpy(x_np).float().to(device)
        for _ in range(max(1, n_pass)):
            model.train() # enable dropout
            outs.append(model(X_t).cpu().numpy().reshape(-1))
        model.eval()
    yhat_norm = np.mean(np.vstack(outs), axis=0)
    return denorm(yhat_norm)
```

## (4) 主入口立即播种

在 main 开头读取固定种子并立刻播种，把所有随机过程纳入同一初始状态。它不修改全局变量，不使用关键字声明，避免作用域陷阱。由此数据划分标准化顺序和初始化轨迹都保持一致，评估可比。

```
base_seed = SEED
_seed_all(base_seed)
```

## (5) 验证集线性标定

用验证集的预测与真实值拟合一元线性映射并设定合理边界，将该修正应用到测试预测。该步骤以极低成本消除整体偏移与缩放误差，常同时改善均方误差与绝对误差，不改变模型结构与训练。

```
try:  
    y_pred_va = predict_rank(x_va_s)  
    y_true_va = y_va * (y_max - y_min) + y_min  
    if len(y_pred_va) == len(y_true_va) and len(y_true_va) >= 5:  
        a, b = np.polyfit(y_pred_va, y_true_va, 1)  
        if np.isfinite(a) and np.isfinite(b):  
            a = float(np.clip(a, 0.5, 2.0))  
            b = float(np.clip(b, -2.0 * abs(np.mean(y_true_va)), 2.0 *  
abs(np.mean(y_true_va))))  
            y_pred_te = a * y_pred_te + b  
    except Exception:  
        pass
```

## (6) 范围裁剪稳住输出

将预测值限制在训练标签的最小最大区间内，切断极端外推带来的离群点。它对区间内多数样本无影响，却能显著降低重尾样本对指标的破坏，提升稳健性，且为纯向量化操作几乎零开销。

```
try:  
    y_pred_te = np.clip(y_pred_te, y_min, y_max)  
except Exception:  
    pass  
  
y_true_te = y_np[idx_te]
```

## (7) 优雅地将预测值保存为整数

先检测真实标签是否近似整数，若成立再对预测做四舍五入。该策略与排名等级类评价的刻度一致，常显著降低绝对误差与均方误差；若标签为连续值则不会触发，避免对连续任务造成干扰。

```
try:  
    if np.allclose(y_true_te, np.round(y_true_te)):  
        y_pred_te = np.round(y_pred_te)  
    except Exception:  
        pass
```

## (8) 批量大小自适应防报错

在训练 DataLoader 处自适应调整批量，保证每批至少两条样本，若最后仅剩一条则丢弃该批。该改动专门化解批归一化对单样本批次的限制，避免报错又不改变总体样本分布，对大数据集无感知。

```
try:  
    _eff_bs = max(2, min(BATCH_SIZE, len(tr_dl.dataset)))  
    _drop_last = (len(tr_dl.dataset) % _eff_bs == 1)  
    tr_dl = DataLoader(tr_dl.dataset, batch_size=_eff_bs, shuffle=True,  
drop_last=_drop_last)  
except Exception:  
    pass
```

## 4. 优化效果展示

运行升级后的代码

各学科排名模型在测试集上的预测值与真实值对比表格，都在附件中：

名称	修改日期	类型	大小
deep_learning.csv	2025/11/8 15:31	Microsoft Excel ...	4 KB
deep_learning.py	2025/11/8 15:29	Python 源文件	16 KB
predictions_AGRICULTURAL_SCIENCE.csv	2025/11/8 15:29	Microsoft Excel ...	5 KB
predictions_BIOLOGY__BIOCHEMISTRY.csv	2025/11/8 15:29	Microsoft Excel ...	6 KB
predictions_CLINICAL_MEDICINE.csv	2025/11/8 15:29	Microsoft Excel ...	26 KB
predictions_COMPUTER_SCIENCE.csv	2025/11/8 15:29	Microsoft Excel ...	3 KB
predictions_ECONOMICS__BUSINESS.csv	2025/11/8 15:29	Microsoft Excel ...	2 KB
predictions_ENGINEERING.csv	2025/11/8 15:30	Microsoft Excel ...	11 KB
predictions_ENVIRONMENT_ECOLOGY.csv	2025/11/8 15:30	Microsoft Excel ...	8 KB
predictions_GEOSCIENCES.csv	2025/11/8 15:30	Microsoft Excel ...	4 KB
predictions_IMMUNOLOGY.csv	2025/11/8 15:30	Microsoft Excel ...	4 KB
predictions_MATERIALS_SCIENCE.csv	2025/11/8 15:30	Microsoft Excel ...	6 KB
predictions_MATHEMATICS.csv	2025/11/8 15:30	Microsoft Excel ...	2 KB
predictions_MICROBIOLOGY.csv	2025/11/8 15:30	Microsoft Excel ...	3 KB
predictions_MOLECULAR_BIOLOGY.csv	2025/11/8 15:30	Microsoft Excel ...	4 KB
predictions_MULTIDISCIPLINARY.csv	2025/11/8 15:30	Microsoft Excel ...	1 KB
predictions_NEUROSCIENCE_BEHAVIORAL.csv	2025/11/8 15:30	Microsoft Excel ...	5 KB
predictions_PHARMACOLOGY_TOXICOLOGY.csv	2025/11/8 15:30	Microsoft Excel ...	5 KB
predictions_PHYSICS.csv	2025/11/8 15:30	Microsoft Excel ...	4 KB
predictions_PLANT_ANIMAL_SCIENCE.csv	2025/11/8 15:31	Microsoft Excel ...	7 KB

我这里随机选一个学科展示一下

和上次不同的是，本次的预测值优雅地保存为了整数：

A	B	C	D	E	F	G	H	I	J	K	L
row_id	true_rank	pred_rank									
1	309	310	299								
2	741	742	749								
3	265	266	266								
5	823	824	827								
6	778	779	782								
7	660	661	656								
8	76	77	89								
9	184	185	182								
10	745	746	753								
11	486	487	479								
12	798	799	786								
13	756	757	755								
14	762	763	754								
15	51	52	64								
16	247	248	247								
17	240	241	251								
18	832	833	836								
19	708	709	714								
20	1148	1149	1136								
21	208	209	199								
22	231	232	241								
23	196	197	204								
24	534	535	524								
25	447	448	440								
26	382	383	367								
27	1118	1119	1125								
28	344	345	349								

各学科排名模型在测试集上的各项指标如下：

subject_csv	MSE	RMSE	MAE	MAPE	MedianAE	R2	nRMSE	Spearman	Kendall
GEOSCIENCES.csv	262.608511	16.205200	11.059574	5.664050	8.0	0.997575	0.013922	0.999071	0.979206
AGRICULTURAL SCIENCES.csv	298.292419	17.271144	9.122744	3.769744	6.0	0.997920	0.012812	0.999249	0.983722
IMMUNOLOGY.csv	495.707627	22.264493	14.148305	10.377401	11.0	0.995497	0.019210	0.997891	0.969571
PLANT & ANIMAL SCIENCE.csv	600.351282	24.502067	18.771795	7.619793	16.0	0.998120	0.012728	0.999470	0.983013
PHYSICS.csv	603.954774	24.575491	16.296482	9.294238	14.0	0.992292	0.025000	0.997050	0.970234
COMPUTER SCIENCE.csv	696.057803	26.382907	12.531792	20.949282	7.0	0.989038	0.030857	0.992624	0.956183
SPACE SCIENCE.csv	698.645833	26.431909	16.145833	35.681213	10.0	0.860681	0.119063	0.936160	0.820082
MICROBIOLOGY.csv	713.701863	26.715199	14.000000	18.383035	10.0	0.986938	0.033520	0.993243	0.960184
BIOLOGY & BIOCHEMISTRY.csv	959.203030	30.971003	15.918182	7.207392	10.0	0.995680	0.019094	0.997985	0.978550
MATERIALS SCIENCE.csv	1815.477848	31.866563	24.041139	8.061457	18.0	0.994667	0.020454	0.997461	0.961185
ENGINEERING.csv	1016.236559	31.878465	24.272401	5.080232	21.0	0.998381	0.011567	0.999448	0.982789
MOLECULAR BIOLOGY & GENETICS.csv	1242.064103	35.242930	13.508547	6.765338	8.0	0.988446	0.030593	0.994661	0.970037
PHARMACOLOGY & TOXICOLOGY.csv	1336.467626	36.557730	29.258993	13.182869	28.0	0.991755	0.026782	0.997508	0.972030
MOLECULAR BIOLOGY & GENETICS.csv	1242.064103	35.242930	13.508547	6.765338	8.0	0.988446	0.030593	0.994661	0.970037
PHARMACOLOGY & TOXICOLOGY.csv	1336.467626	36.557730	29.258993	13.182869	28.0	0.991755	0.026782	0.997508	0.972030
NEUROSCIENCE & BEHAVIOR.csv	1377.703846	37.117433	20.396154	28.307766	12.0	0.990022	0.029066	0.994447	0.959441
PSYCHIATRY PSYCHOLOGY.csv	1577.556522	39.718466	18.956522	13.100281	11.0	0.985447	0.034841	0.992653	0.953411
ECONOMICS & BUSINESS.csv	1658.302752	40.722264	22.302752	19.675153	11.0	0.942551	0.075692	0.970538	0.879750
MATHEMATICS.csv	1904.531646	43.640940	25.265823	195.433016	11.0	0.878532	0.111046	0.931203	0.821714
MULTIDISCIPLINARY.csv	1905.250000	43.649170	22.204545	24.706810	4.5	0.448791	0.213967	0.706115	0.628513
SOCIAL SCIENCES, GENERAL.csv	2046.433610	45.237524	24.334025	5.848389	19.0	0.995449	0.018975	0.998333	0.985037
CHEMISTRY.csv	3593.167832	59.943038	29.573427	9.756710	18.0	0.990825	0.028275	0.996710	0.974176
ENVIRONMENT ECOLOGY.csv	5883.437198	76.703567	43.219807	12.0895035	31.0	0.983568	0.037600	0.993607	0.959125
CLINICAL MEDICINE.csv	6718.927461	81.969064	60.293856	3.988473	51.0	0.998248	0.012167	0.999467	0.987425

其中的 nRMSE 是 RMSE 按数值范围归一化之后的，很有参考价值：

nRMSE < 0.10：很好

0.10–0.20：可用

而各学科排名模型在测试集上的 nRMSE 都小于 0.10，说明模型预测效果很好

各项指标的总表也在附件中，叫 deep\_learning.csv

升级后，nRMSE 大部分在 0.01 到 0.02，预测效果有了显著提升