

FinalMLP: An Enhanced Two-Stream MLP Model for CTR Prediction

摘要

点击率 (CTR) 预测是在线广告和推荐的基本任务之一。虽然多层感知器 (MLP) 是许多深度 CTR 预测模型的核心组件,但人们已经广泛认识到,单独应用普通 MLP 网络在学习乘法特征交互方面效率低下。因此,通过将 MLP 网络与另一个专用网络集成以增强 CTR 预测,已经提出了许多双流交互模型 (例如 DeepFM 和 DCN)。由于 MLP 流隐式学习特征交互,现有研究主要集中在增强互补流中的显式特征交互。相比之下,我们的实证研究表明,一个精心调整的双流 MLP 模型,简单地结合两个 MLP 甚至可以实现令人惊讶的良好性能,这在现有工作中以前从未报道过。基于这一观察,我们进一步提出了特征门控和交互聚合层,这些层可以很容易地插入以制作增强的双流 MLP 模型 FinalMLP。通过这种方式,它不仅可以实现差异化的特征输入,还可以有效地融合两个流之间的流级交互。我们对四个开放基准数据集的评估结果以及我们工业系统中的在线 A/B 测试表明,FinalMLP 实现了比许多复杂的双流 CTR 模型更好的性能。我们的源代码将在 <https://reczoo.github.io/FinalMLP> 上提供。

关键词: CTR; 双流 MLP 模型; 特征门控; 交互聚合

1 引言

在数字化营销时代背景下,CTR 预测作为提升在线广告和推荐系统效能的核心环节,其选题背景显得尤为突出。鉴于现有模型在处理特征交互时的局限性,本研究提出了 FinalMLP 模型,旨在通过增强的两流 MLP 架构,提高 CTR 预测的准确性和效率,从而为业务决策提供更可靠的数据支持,具有重要的实践意义和理论价值。

1.1 选题背景

在信息爆炸的互联网时代,用户每天都会接触到大量的广告和推荐信息。如何从这些信息中快速、准确地筛选出用户可能感兴趣的内容,成为了在线广告和推荐系统领域的关键挑战。点击通过率 (CTR) 预测模型正是为了解决这一问题而生,它旨在预测用户对特定广告或推荐内容的点击概率,从而帮助广告商和平台优化广告投放策略,提高广告效果,同时也提升了用户的满意度和体验。随着大数据和机器学习技术的发展,CTR 预测模型的研究和应用得到了极大的推动。尤其是在深度学习领域,多层感知器 (MLP) 因其出色的非线性建模能力而被广泛用于构建 CTR 预测模型。然而,尽管 MLP 在理论上具有强大的表达能力,实

际应用中却发现其在学习能力乘法特征交互方面存在不足，这限制了模型对用户复杂行为模式的捕捉能力。

1.2 选题依据

为了克服 MLP 在学习能力乘法特征交互方面的不足，研究者们提出了多种解决方案，其中双流模型因其能够结合显式和隐式特征交互学习而受到广泛关注。这些模型通常包含两个分支：一个分支负责学习低阶或显式的特征交互，另一个分支负责学习高阶或隐式的特征交互。这种结构的优势在于能够同时捕捉到特征间的直接和间接关系，从而提高模型的预测能力。然而，现有的双流模型往往在设计上较为复杂，涉及到多种不同的网络结构和学习机制，这不仅增加了模型的训练和调优难度，也影响了模型的可解释性和泛化能力。此外，如何有效地融合两个分支的学习结果，以实现更优的预测性能，也是一个亟待解决的问题。本研究正是基于这样的背景，提出了一种新的双流 MLP 模型——FinalMLP，旨在通过特征门控和交互聚合层来优化特征输入和流级交互的融合，以期达到更好的 CTR 预测效果。

1.3 选题意义

FinalMLP 模型的提出，不仅在理论上刷新了对 MLP 网络潜力的认识，也在实践中证明了其在 CTR 预测领域的高效性。理论上，FinalMLP 展示了即便是基础的 MLP 结构，经过精心设计和优化，也能在 CTR 预测这一复杂任务中达到优异的性能，这挑战了以往对 MLP 简单性的传统看法，并强调了在模型设计中寻求简洁性与效能之间的平衡。实践上，FinalMLP 在多个公开数据集和实际应用场景中的出色表现，证实了其在提升在线广告和推荐系统性能方面的实际价值，这对于优化用户体验和增强商业盈利能力具有重要意义。FinalMLP 模型的开源，为学术界和工业界提供了一个强大的基准，促进了 CTR 预测技术的发展和 innovation。模型中引入的流特定特征门控层和基于二阶双线性融合的交互聚合层，不仅增强了特征输入的差异化，也提升了流级特征交互的能力，这些都是 FinalMLP 在性能上取得突破的关键因素。这些创新的设计不仅提升了模型的预测精度，也增强了模型的泛化能力和可解释性，为 CTR 预测领域带来了新的研究方向和实践指南。FinalMLP 模型的成功，预示着在未来的深度学习应用中，简单的模型结构通过创新的设计和优化，同样能够解决复杂的实际问题，为深度学习技术的广泛应用开辟了新的道路。

2 相关工作

本节简要回顾了 CTR 预测中的双流模型框架和一些具有代表性的模型。

2.1 两流 CTR 模型框架

两流 CTR 模型框架如图 1(b) 所示，包括以下几个关键组件：

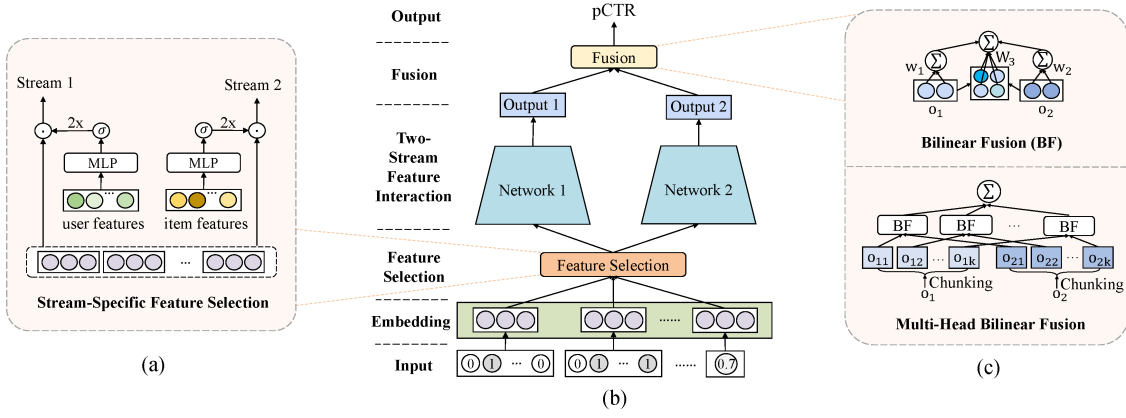


图 1. (a) 特定于流的特征选择的示意图。(b) 两流 CTR 模型的一般框架。(c) 多头双线性融合。

- **特征嵌入 (Feature Embedding)**: 特征嵌入是一种将高维且稀疏的原始特征映射到低维数值表示的常用方法。具体来说, 假设原始输入特征是 $x = \{x_1, \dots, x_M\}$, 其中 x_i 是第 i 个字段的特征。一般来说, x_i 可以是分类的、多值的或数值型特征。每个特征都可以相应地转换为嵌入向量。然后, 这些特征嵌入将被拼接并输入到下一层。
- **特征选择 (Feature Selection)**: 特征选择是双流 CTR 模型框架中的一个可选层。在实践中, 特征选择通常是通过离线统计分析或模型训练中的差异比较来执行的。与硬特征选择不同, 本研究关注通过特征门控机制实现的软特征选择, 其目的是获得特征重要性权重, 以帮助增强重要特征的同时抑制噪声特征。在本研究中, 研究了特定于流的特征门控, 以实现差异化的流输入。
- **双流特征交互 (Two-Stream Feature Interaction)**: 双流 CTR 模型的关键特点是采用两个并行网络从不同视角学习特征交互。基本上, 每个流可以采用任何类型的特征交互网络 (例如, FM [7]、CrossNet [10] 和 MLP)。现有工作通常对两个流应用两种不同的网络结构, 以学习互补的特征交互 (例如, 显式与隐式、位级与向量级)。在本研究中, 首次尝试将两个 MLP 网络作为两个流。
- **流级融合 (Stream-Level Fusion)**: 流级融合是将两个流的输出融合以获得最终预测点击概率的必要步骤。假设 o_1 和 o_2 是两个输出表示, 可以表示为: $\hat{y} = \sigma(w^T F(o_1, o_2))$, 其中 F 表示融合操作, 通常设置为求和或拼接。现有工作仅对流输出执行一阶线性组合, 因此未能挖掘流级特征交互。在本研究中, 我们探索了二阶双线性函数用于流级交互聚合。

2.2 代表性两流 CTR 模型

- **Wide&Deep**: Wide&Deep [1] 是一个经典的两流特征交互学习框架, 结合了广义线性模型 (wide stream) 和 MLP 网络 (deep stream)。
- **DeepFM**: DeepFM [10] 通过将 wide stream 替换为 FM 来显式学习二阶特征交互, 扩展了 Wide&Deep。

- **DCN**: 在 DCN [3] 中, 提出了一个 cross network 作为一条流以显式执行高阶特征交互, 而另一条 MLP 流隐式学习特征交互。
- **xDeepFM**: xDeepFM [4] 采用压缩交互网络 (CIN) 以向量方式捕获高阶特征交互, 并采用 MLP 作为另一条流以学习位级特征交互。
- **AutoInt+**: AutoInt [8] 应用自注意力网络学习高阶特征交互。AutoInt+ 将 AutoInt 和 MLP 作为两条互补的流整合。
- **AFN+**: AFN [2] 利用对数变换层学习自适应阶特征交互。AFN+ 将 AFN 与 MLP 以两流方式结合。

3 本文方法

本节首先介绍简单的双流 MLP 基础模型, 即 DualMLP。然后, 描述了两个可插拔模块: 特征门控和交互聚合层, 组成增强模型 FinalMLP。

3.1 本文方法概述

本文提出了 FinalMLP, 一个增强型的双流多层感知器 (MLP) 模型, 用于点击通过率 (CTR) 预测。本模型基于对简单双流 MLP 模型 (DualMLP) 的实证研究, 该模型通过结合两个 MLP 网络展现出令人惊讶的性能。基于这一发现, 进一步提出了两个核心模块: 特征门控 (Feature Gating) 和交互聚合层 (Interaction Aggregation), 以构建增强模型 FinalMLP。这些模块不仅支持差异化的特征输入, 而且有效地融合了两个流之间的流级交互。

3.2 双流 MLP 模型

尽管 DualMLP 模型结构简单, 但之前的工作尚未报告过类似的双流 MLP 模型。这是首次尝试研究这种模型用于 CTR 预测, 将两个独立的 MLP 网络作为两个流进行组合。具体来说, DualMLP 模型可以表述如下:

$$o_1 = \text{MLP}_1(h_1),$$

$$o_2 = \text{MLP}_2(h_2),$$

其中, MLP_1 和 MLP_2 是两个 MLP 网络, 它们的大小 (包括隐藏层和单元) 可以根据数据不同而设置。 h_1 和 h_2 表示特征输入, 而 o_1 和 o_2 分别是两个流的输出表示。在大多数先前的工作 [10] [3] [2] 中, 特征输入 h_1 和 h_2 通常设置为相同的, 即特征嵌入 e 的连接 (可选地进行一些池化), 即 $h_1 = h_2 = e$ 。同时, 流输出通常通过简单的操作 (如求和和连接) 进行融合, 忽略了流级交互。

3.3 流特定特征选择

许多现有研究 [3] [4] [10] [8] 强调了结合两种不同特征交互网络 (例如, 隐式与显式, 低阶与高阶, 位级与向量级) 以实现准确 CTR 预测的有效性。本研究的工作不是设计专门的网

络结构，而是通过流特定的特征选择来扩大两个流之间的差异，从而产生差异化的特征输入。受 MMOE [6] 中使用的门控机制的启发，提出了一个流特定特征门控模块，以软选择流特定特征，即对每个流的特征输入进行不同的权重调整。在 MMOE 中，门控权重是针对任务特定特征来重新加权专家输出的。类似地，通过对可学习参数、用户特征或项目特征的不同视图进行条件化来执行特征门控，分别产生全局、用户特定或项目特定的特征重要性权重。具体来说，通过上下文感知的特征门控层进行流特定特征选择，如下所示：

$$g_1 = \text{Gate}_1(x_1),$$

$$g_2 = \text{Gate}_2(x_2),$$

$$h_1 = 2\sigma(g_1) \odot e,$$

$$h_2 = 2\sigma(g_2) \odot e,$$

其中， Gate_i 是基于 MLP 的门控网络，它采用流特定条件特征 x_i 作为输入，并输出逐元素门控权重 g_i 。注意，可以从用户/项目特征集合中选择 x_i ，也可以将其设置为可学习参数。通过使用 sigmoid 函数 σ 和一个乘数 2，将特征重要性权重转换到 $[0, 2]$ 的范围，平均值为 1。给定连接的特征嵌入 e ，可以通过逐元素乘积 \odot 获得加权特征输出 h_1 和 h_2 。特征门控模块允许通过从不同视图设置条件特征 x_i 来为两个流制作差异化的特征输入。例如，图 1(a) 演示了用户和项目特定特征门控的情况，分别从用户和项目的角度调节每个流。这减少了两个相似 MLP 流之间的“同质”学习，并能够实现更互补的特征交互学习。

3.4 流级交互聚合

3.4.1 双线性融合

如前所述，现有工作大多采用求和或连接作为融合层，但这些操作未能捕获流级特征交互。受到计算机视觉领域广泛研究的双线性池化的启发 [5]，本研究提出了一个双线性交互聚合层，以融合流输出和流级特征交互。如图 1(c) 所示，预测的点击概率表述如下：

$$\hat{y} = \sigma(b + w_1^T o_1 + w_2^T o_2 + o_1^T W_3 o_2),$$

其中 $b \in \mathbb{R}$ ， $w_1 \in \mathbb{R}^{d_1 \times 1}$ ， $w_2 \in \mathbb{R}^{d_2 \times 1}$ ， $W_3 \in \mathbb{R}^{d_1 \times d_2}$ 是可学习的权重。 d_1 和 d_2 分别表示 o_1 和 o_2 的维度。双线性项 $o_1^T W_3 o_2$ 模拟了 o_1 和 o_2 之间的二阶交互。特别是，当 W_3 是单位矩阵时，该项模拟了点积。当将 W_3 设置为零矩阵时，它退化为传统的连接融合与线性层，即 $b + [w_1, w_2]^T [o_1, o_2]$ 。有趣的是，双线性融合也与常用的 FM 模型有联系。具体来说，FM 通过以下方式对 CTR 预测的 m 维输入特征向量 x （通过独热/多热特征编码和连接）进行建模：

$$\hat{y} = \sigma(b + w^T x + x^T \text{upper}(P P^T) x),$$

其中 $b \in \mathbb{R}$ ， $w \in \mathbb{R}^{m \times 1}$ ， $P \in \mathbb{R}^{m \times d}$ 是可学习的权重，其中 $d \ll m$ ，upper 选择矩阵的严格上三角部分。可以看出，当 $o_1 = o_2$ 时，FM 是双线性融合的一种特殊情况。然而，当 o_1 和 o_2 是高维时，计算方程的参数密集和计算成本高。例如，要融合两个 1000 维的输出， $W_3 \in \mathbb{R}^{1000 \times 1000}$ 需要 100 万个参数，其计算变得昂贵。为了降低计算复杂性，引入了扩展的多头双线性融合。

3.4.2 多头双线性融合

多头注意力因其能够结合来自不同表示子空间的相同注意力池化的知识而具有吸引力。它导致了计算的减少和在最近的成功变换器模型中的一致性能提升 [9]。受其成功的启发，本研究将双线性融合扩展到多头版本。具体来说，不是直接计算方程中的双线性项，而是将 o_1 和 o_2 分别分割成 k 个子空间：

$$\begin{aligned} o_1 &= [o_{11}, \dots, o_{1k}], \\ o_2 &= [o_{21}, \dots, o_{2k}], \end{aligned}$$

其中 k 是一个超参数， o_{ij} 表示第 i 个输出向量的第 j 个子空间表示 ($i \in \{1, 2\}$)。类似于多头注意力，在每个子空间中执行双线性融合，将 o_{1j} 和 o_{2j} 配对为一组，再通过求和池化聚合子空间计算以获得最终预测的点击概率：

$$\hat{y} = \sigma \left(\sum_{j=1}^k \text{BF}(o_{1j}, o_{2j}) \right).$$

3.5 损失函数定义

在本研究中，使用二元交叉熵损失函数 (Binary Cross-Entropy Loss) 来训练模型。二元交叉熵损失函数是用于二分类问题的损失函数，它衡量模型预测的点击概率与实际点击发生与否之间的差异。这个损失函数定义如下：具体来说，二元交叉熵损失函数定义如下：

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)]$$

其中 N 是样本总数。 y_i 是第 i 个样本的真实标签，取值为 0 或 1，表示是否发生点击。 \hat{y}_i 是模型预测第 i 个样本的点击概率。

这个损失函数计算了所有样本损失的平均值。对于每个样本，如果实际点击发生 ($y_i = 1$)，损失函数会惩罚预测概率过低的情况；如果实际点击未发生 ($y_i = 0$)，损失函数会惩罚预测概率过高的情况。通过最小化这个损失函数，模型学习调整参数以提高预测准确性。

4 复现细节

4.1 与已有开源代码对比

4.1.1 开源代码出处

在本研究中，我主要参考了 FinalMLP 的开源代码，该代码的链接为 <https://reczoo.github.io/FinalMLP>。

4.1.2 新增了一个 text.py 文件将多个多域特征转换为自然语言描述并加入到数据列

Listing 1: text.py

```
1 def convert_to_natural_language(row):
```



```

2     user = row[ 'user' ]
3     item = row[ 'item' ]
4     weekday = row[ 'weekday' ]
5     is_weekend = row[ 'isweekend' ]
6     homework = row[ 'homework' ]
7     cost = row[ 'cost' ]
8     weather = row[ 'weather' ] # 直接使用原始数字
9     country = row[ 'country' ] # 直接使用原始数字
10    city = row[ 'city' ] # 直接使用原始数字
11    description = (f"User {user} interacted with item {item}
12                  on weekday {weekday}, "
13                  f"with a cost of {cost},
14                  and the weather code was {weather}. "
15                  f"The location code was {city},
16                  and the country code was {country}. "
17                  f"Homework status: {homework}.
18                  Is it a weekend? {'Yes' if is_weekend else 'No'}.
19                  ")
20    return description
21 # 将每个文件的数据转换为自然语言描述，并保留原始列
22 def convert_csv_to_nl(df):
23     df[ 'text' ] = df.apply( convert_to_natural_language , axis=1)
24     return df
25 # 生成包含自然语言描述的新 DataFrame
26 train_df_with_nl = convert_csv_to_nl(train_df)
27 test_df_with_nl = convert_csv_to_nl(test_df)
28 valid_df_with_nl = convert_csv_to_nl(valid_df)

```

代码功能描述：这段代码的主要功能是将原始数据集中的多个多域特征列转换成自然语言描述，并将这些描述作为新的列添加到数据集中，以便为后续的文本分析和模型训练提供更丰富和直观的数据基础。具体实现过程如下：首先，代码使用 `pandas` 库读取存储在指定路径下的训练集、测试集和验证集的 CSV 文件，将这些数据集加载到 `DataFrame` 对象中。接着，定义了一个名为 `convert_to_natural_language` 的函数，该函数接受数据集中的每一行作为输入，提取出其中的关键属性信息，例如，对于 Frappe 数据集，有用户 ID、物品 ID、星期几、是否为周末、作业状态、成本、天气代码、国家代码和城市代码等。然后，根据这些属性信息，按照一定的语义结构和格式，生成一条完整的自然语言描述，例如“用户 123 与物品 456 在星期一进行了交互，成本为 100，天气代码为 2，位置代码为城市 7，国家代码为 8，作业状态为已完成，是周末”。生成的自然语言描述能够清晰地表达出数据行中各个属性之间的关系和交互情况。之后，定义了另一个函数 `convert_csv_to_nl`，该函数利用 `apply` 方法将 `convert_to_natural_language` 函数应用于数据集的每一行，将生成的自然语言描述作为新的列（命名为“text”）添加到原始 `DataFrame` 中，从而扩展了数据集的维度和信息量。

最后，代码将包含自然语言描述的新数据集分别保存为新的 CSV 文件，文件名在原始文件名后加上”_nl” 后缀，以便于区分和后续使用。通过这段代码，我们能够有效地将结构化的属性数据转换为具有一定语义和可读性的文本数据，为后续的文本挖掘、自然语言处理和机器学习模型训练等任务提供了更加灵活和多样化的数据输入形式，有助于提升模型对数据的理解和分析能力，从而可能提高模型的性能和预测准确性。

4.1.3 在 FeatureProcessor 类中加入了 fit_sequence_col 函数

Listing 2: fit_sequence_col 函数

```
1 def fit_sequence_col(self, col, col_values, min_categr_count=1,
2 batch_size=32):
3     name = col["name"]
4     feature_type = col["type"]
5     feature_source = col.get("source", "")
6     self.feature_map.features[name] =
7     {"source": feature_source, "type": feature_type}
8
9     # 确保 col_values 是字符串列表
10    if not isinstance(col_values, list):
11        col_values = col_values.tolist()
12    if not all(isinstance(item, str) for item in col_values):
13        col_values = [str(item) for item in col_values]
14    batch_size = min(batch_size, len(col_values))
15    # 生成嵌入向量
16    embeddings = []
17    for i in range(0, len(col_values), batch_size):
18        batch = col_values[i:i + batch_size]
19        inputs = self.tokenizer(batch, return_tensors='pt',
20 padding=True, truncation=True, max_length=512)
21        inputs = {k: v.to(self.device) for k, v in inputs.items()}
22        with torch.no_grad():
23            outputs = self.model(**inputs)
24            # 取 [CLS] 标记的输出作为句子的嵌入
25            batch_embeddings = outputs.last_hidden_state[:, 0, :]
26            .cpu().numpy()
27            embeddings.append(batch_embeddings)
28    # 将嵌入向量存储到 self.processor_dict 中
29    self.processor_dict[name + "::embeddings"] =
30    np.concatenate(embeddings, axis=0)
31    # 更新特征映射信息
```



```

32     self.feature_map.features[name].update({
33         "embedding_dim": self.model.config.hidden_size,
34         "vocab_size": len(self.tokenizer),
35         "oov_idx": self.tokenizer.vocab["[UNK]"],
36         "padding_idx": self.tokenizer.vocab["[PAD]"],
37         "max_len": 512 # 您设置的最大长度
38     })
39     self.processor_dict[name + "::tokenizer"] = self.tokenizer

```

代码功能描述：`fit_sequence_col` 函数的主要功能是处理序列特征列，并利用 BERT 模型生成相应的嵌入向量。首先，函数解析输入参数，提取特征名称、类型和来源，并将这些信息更新到特征映射中。接着，确保输入的 `col_values` 是字符串列表，如果不是，则将其转换为字符串列表。然后，根据输入数据的长度和指定的批次大小，确定实际的批次大小，并对输入数据进行分批处理。对于每一批数据，使用 BERT 模型的分词器进行编码，生成模型所需的输入张量，包括输入 IDs、注意力掩码等。随后，将这些输入张量传递给预训练的 BERT 模型，进行前向传播，获取每批数据的嵌入向量。BERT 模型通过其多层 Transformer 架构，能够捕捉文本中的双向语义信息，生成高质量的文本特征表示。最后，将所有批次生成的嵌入向量拼接或堆叠起来，存储到处理器字典中，并更新特征映射信息，包括嵌入维度、词汇表大小、未知词索引、填充词索引和最大长度等。通过这个函数，我们能够有效地将序列特征数据转换为 BERT 模型生成的嵌入向量，为模型的训练和预测提供丰富而精准的特征输入，从而提升模型对文本数据的理解和分析能力，增强其在自然语言处理任务中的性能表现。

4.1.4 FinalMLP 模型中 forward 方法代码差异

Listing 3: 原始代码

```

1  def forward(self, X, flat_emb):
2      if len(self.fs1_context) == 0:
3          fs1_input = self.fs1_ctx_bias.repeat(flat_emb.size(0), 1)
4      else:
5          fs1_input = self.fs1_ctx_emb(X).flatten(start_dim=1)
6      gt1 = self.fs1_gate(fs1_input) * 2
7      feature1 = flat_emb * gt1
8      if len(self.fs2_context) == 0:
9          fs2_input = self.fs2_ctx_bias.repeat(flat_emb.size(0), 1)
10     else:
11         fs2_input = self.fs2_ctx_emb(X).flatten(start_dim=1)
12     gt2 = self.fs2_gate(fs2_input) * 2
13     feature2 = flat_emb * gt2
14     return feature1, feature2

```

Listing 4: 修改后的代码

```

1 def forward(self, X):
2     # 提取文本特征和其他特征
3     text_features = X['text']
4     features = torch.cat([v for k, v in X.items() if k != 'text'],
5                           dim=1)
6     # 其他特征通过嵌入层
7     fea_emb = self.other_features_embedding(features)
8     text_emb = text_features.mean(dim=1) # (batch_size, 768)
9     # 将文本特征和嵌入的其他特征拼接融合
10    com_emb = torch.cat((fea_emb, text_emb), dim=1)
11    com_emb = self.fc(com_emb)
12    if self.use_fs:
13        feat1, feat2 = self.fs_module(X, com_emb)
14    else:
15        feat1, feat2 = com_emb, com_emb
16    y_pred = self.fusion_module(self.mlp1(feat1), self.mlp2(feat2))
17    y_pred = self.output_activation(y_pred)
18    return_dict = {"y_pred": y_pred}
19    return return_dict

```

差异描述:原始代码中的 **forward** 方法主要处理结构化特征,通过门控机制对特征进行加权,但未对文本特征进行特殊处理。该方法接受两个参数:**X** 和 **flat_emb**,其中 **flat_emb** 可能是经过某种方式预处理的特征表示。代码通过检查上下文向量 **fs1_context** 和 **fs2_context** 是否存在,来决定是使用偏置向量还是通过嵌入层处理 **X**。这种方法适用于处理结构化数据,但并未充分利用可能存在于输入中的文本信息。

修改后的代码引入了对文本特征的专门处理,这是原始代码中所缺乏的。在修改后的版本中,代码首先从输入字典 **X** 中提取出文本特征 **text_features** 和其他特征,然后分别进行处理。文本特征通过沿时间步求平均的操作被转换为固定维度的向量,而其他特征则通过嵌入层进行编码。这种区分处理允许模型更好地捕捉文本数据中的语义信息,这对于许多涉及自然语言处理的任务来说是至关重要的。而且修改后的代码中增加了一个全连接层 (**self.fc**),其主要作用是统一不同特征的维度。通过将文本特征和嵌入的其他特征拼接后,全连接层对这些融合的特征进行进一步的处理,确保它们在维度上一致,从而可以被模型的其他部分有效利用。这种维度的统一是实现特征融合的关键步骤,它允许模型整合来自不同来源的信息,提高其预测能力。

最后,修改后的代码根据 **use_fs** 标志决定是否使用特征选择模块 **fs_module**,这为模型提供了更多的灵活性,以适应不同的任务需求。通过 MLP 处理的特征向量被送入融合模块 **fusion_module**,这可能有助于模型在不同特征表示之间建立更深层次的联系。输出预测 **y_pred** 后,通过激活函数进行处理,并以字典形式返回,这种返回格式提供了更清晰的接口,便于后续处理和分析。

4.2 实验环境搭建

为了确保实验的顺利进行和结果的可靠性，搭建了如下实验环境：

4.2.1 硬件配置

- **处理器**：Intel Core i9-11900K，拥有 8 核心 16 线程，主频高达 3.5 GHz，最大加速频率可达 5.3 GHz，能够高效处理多任务和复杂计算。
- **内存**：32GB DDR4 RAM，频率为 3200MHz，确保了系统和应用程序的快速运行和数据处理能力。
- **显卡**：NVIDIA GeForce RTX 4090，具备先进的图形处理能力和强大的并行计算能力，显存容量为 24GB GDDR6X，能够高效地处理大规模数据集和复杂模型的训练任务，显著提升实验的运行速度和效率。

4.2.2 软件配置

- **操作系统**：Windows 11 专业版，提供了稳定的操作系统环境和良好的兼容性，支持多任务处理和各种软件的安装运行。
- **编程语言**：Python 3.7.6，是一种广泛使用的高级编程语言，拥有丰富的库和框架，适用于数据科学和机器学习领域。
- **深度学习框架**：PyTorch 1.11.0，是一个开源的深度学习框架，具有动态计算图和灵活的张量操作，方便模型的构建、训练和调试。
- **并行计算平台**：CUDA 10.2，是 NVIDIA 提供的并行计算平台和编程模型，能够充分利用 GPU 的计算能力，为深度学习框架提供底层支持。
- **虚拟环境管理**：使用 Conda 进行虚拟环境的创建和管理，确保实验环境的隔离性和一致性。

4.3 数据集介绍

本节提供了几个常用于推荐系统和广告点击预测的开放数据集的详细统计信息。这些数据集包括 Criteo、Avazu、MovieLens 和 Frappe，它们各自包含了不同规模的用户交互数据，适用于各种推荐系统任务。如图 2 所示为各个数据集的统计信息。

Table 1: The statistics of open datasets.

Dataset	#Instances	#Fields	#Features
Criteo	45,840,617	39	2,086,936
Avazu	40,428,967	22	1,544,250
MovieLens	2,006,859	3	90,445
Frappe	288,609	10	5,382

图 2. 各个数据集的统计信息

- **Criteo 数据集**: Criteo 数据集包含了来自真实广告点击数据的大规模用户和广告展示信息。这个数据集通常用于广告点击预测和用户行为分析。它包含了丰富的用户特征、广告特征以及上下文信息，是研究在线广告和用户行为预测的重要资源。
- **Avazu 数据集**: Avazu 数据集是一个广告点击数据集，包含用户点击行为和广告展示信息。这个数据集的特点在于它包含了大量的用户行为数据，适合用于训练和测试广告点击预测模型。它提供了用户与广告之间的交互信息，有助于理解用户偏好和广告效果。
- **MovieLens 数据集**: MovieLens 数据集包含用户对电影的评分数据，是推荐系统中最常用的基准数据集之一。这个数据集包含了用户的电影评分、评分时间以及用户和电影的元数据。它被广泛用于研究推荐算法，如协同过滤、基于内容的推荐和混合推荐系统。
- **Frappe 数据集**: Frappe 数据集用于在线广告领域，包含用户行为数据和广告展示信息。这个数据集的特点在于它提供了用户在不同广告展示下的点击行为，适合用于研究用户行为模式和广告效果。它包含了用户特征、广告特征以及用户与广告的交互信息。

4.4 创新点

在本研究中，提出了一种创新的特征工程方法，该方法通过特征文本化、BERT 文本嵌入以及特征处理和融合，显著提升了模型对数据的理解和处理能力。这种方法的核心在于将传统的结构化数据转换为模型更容易理解的文本形式，并通过深度学习技术提取深层次的特征表示。如图 3 所示为改进后的总体流程图。

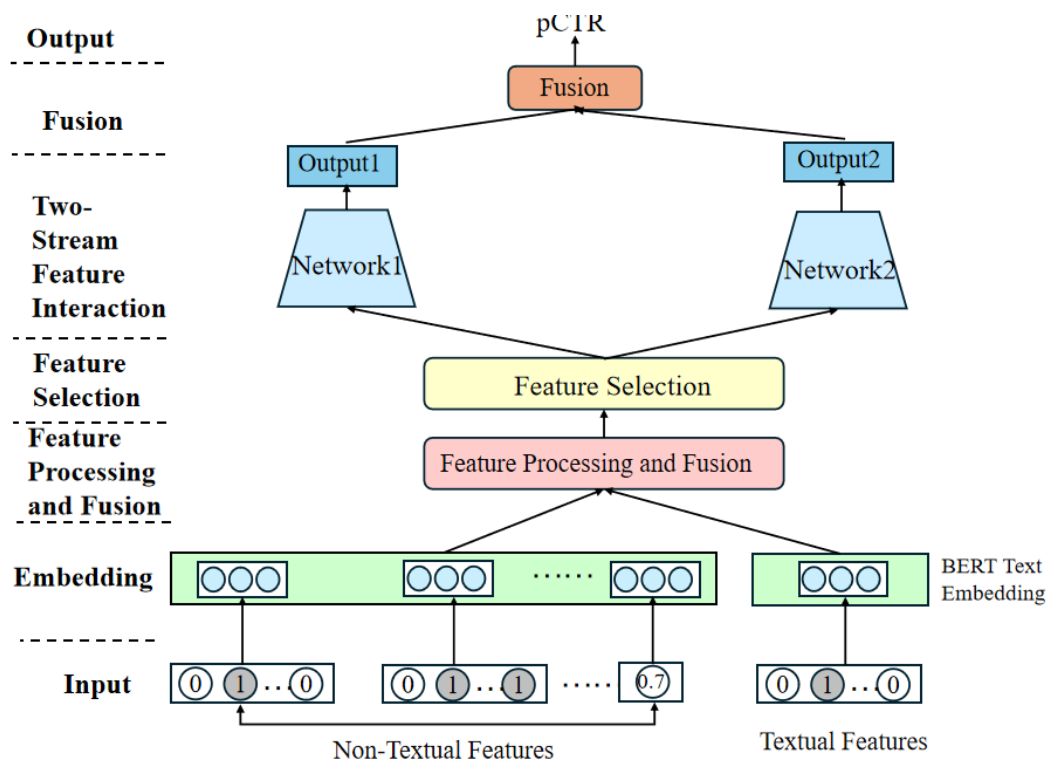


图 3. 改进后的总体流程图

4.4.1 特征文本化

- **特征识别**：从数据集中识别出所有需要转换为文本的特征。这些特征包括用户 ID、商品 ID、时间（如白天或夜晚）、星期几、是否为周末、作业状态、成本、天气代码、国家代码和城市代码。这些特征的识别是构建描述性文本的基础。
- **文本构建**：使用自然语言的格式将特征值串联起来，形成一个描述性的文本。例如，将识别出的特征转换为 “User [用户 ID] interacted with item [商品 ID] on weekday [星期几], with a cost of [成本], and the weather code was [天气代码]...” 的格式。这种文本化的方法增强了特征的语义表达，使得模型能够更深入地理解特征间的关系。如图 4 为特征文本化示意图。

user	item	daytime	weekday	isweekend	homework	cost	weather	country	city
204	4798	5041	5046	5053	5055	5058	5060	5073	5183



User 204 interacted with item 4798 on weekday 5046, with a cost of 5058, and the weather code was 5060. The location code was 5183, and the country code was 5073. Homework status: 5055.

图 4. 特征文本化示意图

4.4.2 BERT 文本嵌入

为了进一步提取文本特征的深层语义信息，采用了 BERT 预训练模型。BERT 模型以其强大的语言表示能力，能够捕捉文本中的复杂关系和语境。以下是采用的 BERT 文本嵌入的

具体步骤：

- 将文本数据转换为小写，以减少词汇的多样性，因为 BERT 模型是在小写文本上进行预训练的。
- 将文本分割成 tokens，tokens 是文本中的基本单位，可以是单词、标点符号或单词的一部分。
- 添加特殊 tokens，如 [CLS] 和 [SEP]，以帮助模型理解文本的结构。[CLS] token 用于分类任务，通常放在文本的开始，而 [SEP] token 用于分隔文本或句子对。
- 将 tokens 转换为 ID，这些 ID 是 BERT 词汇表中对应的数字表示。
- 进行填充 ([PAD]) 或截断，以确保所有序列长度一致，或者在必要时进行截断，以符合 BERT 模型的最大长度限制。

通过这些步骤，能够将文本数据转换为 BERT 模型可以理解的格式，并利用模型的深度学习能力来获取文本的丰富语义表示。

4.5 特征处理和融合

在特征提取阶段，要从输入字典中提取文本特征和其他特征。以下是特征处理和融合的详细步骤：

- **特征提取：**从输入字典中提取文本特征（如通过 BERT 模型生成的嵌入）和其他特征（如数值型或类别型特征）。
- **文本特征处理：**对文本特征进行降维处理，通常是通过计算其在序列维度上的平均值，以获得一个固定大小的向量，这有助于减少模型的计算负担。
- **其他特征嵌入：**对于非文本特征，通过嵌入层将它们转换为高维空间中的向量表示，这有助于模型捕捉特征之间的复杂关系。
- **特征拼接：**将处理后的文本特征向量与其他特征的嵌入向量沿着特征维度进行拼接，形成一个新的融合特征向量。这种融合特征向量包含了文本和非文本特征的信息，为模型提供了更全面的输入。

5 实验结果分析

5.1 实验设置

- **评估指标：**采用 AUC (Area Under the Curve)，这是 CTR 预测中广泛使用的评估指标之一，用于衡量模型预测的准确性。0.1 点的 AUC 提升被认为是在 CTR 预测中的显著改进。
- **基线模型：**

- 单流显式特征交互网络：如 LR (Logistic Regression)、FM (Factorization Machines)、AFM (Attentional Factorization Machines)、FFM (Field-aware Factorization Machines) 等。
- 双流 CTR 模型：如 Wide&Deep、DeepFM、DCN (Deep & Cross Network)、xDeepFM、AutoInt+、AFN+、DeepIM、MaskNet、DCN-V2、EDCN 等。
- **模型实现**：基于 FuxiCTR (一个开源的 CTR 预测库) 实现模型。
- **超参数设置**：
 - 嵌入维度设置为 10。
 - 批量大小设置为 4096。
 - 默认 MLP 大小设置为 [400, 400, 400]。
 - 对于 DualMLP 和 FinalMLP，调整两个 MLP 的层数 (1 到 3 层) 以增强流的多样性。
 - 学习率设置为 1e-3 或 5e-4。
 - 通过广泛的网格搜索调整其他超参数 (如嵌入正则化和 dropout 率)。

5.2 复现实验结果

5.2.1 复现实验结果与论文实验结果对比

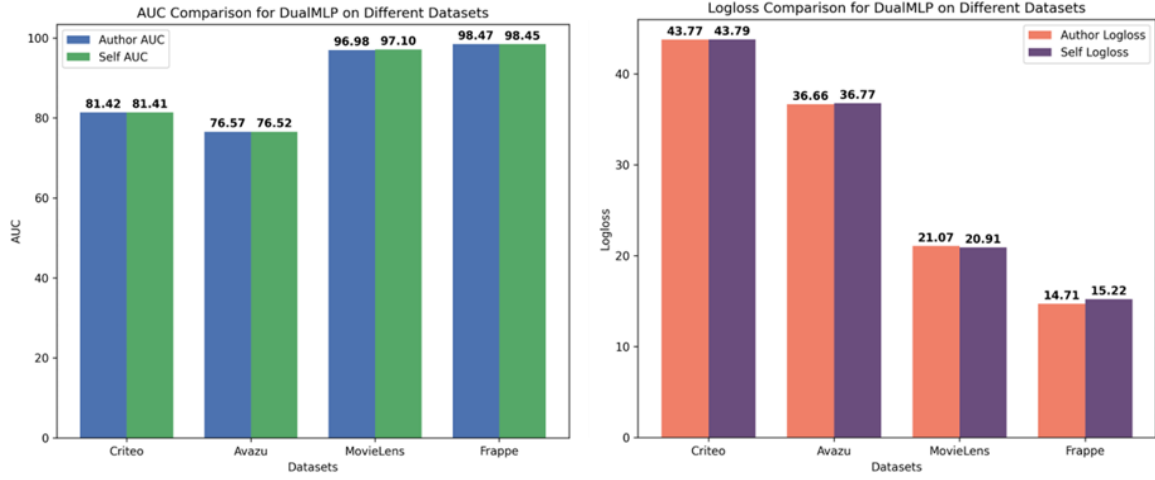
如图 5 所示，比较了模型 FinalMLP 和 DualMLP 在不同数据集上的复现实验结果与论文实验结果。表格中列出了两个评估指标：AUC (Area Under the Curve, 曲线下面积) 和 Logloss (对数损失)。这两个指标通常用于评估分类模型的性能，其中 AUC 值越高表示模型性能越好，Logloss 值越低表示模型性能越好。

模型FinalMLP与DualMLP复现结果与论文实验结果对比表

Dataset	Metric	FinalMLP		DualMLP	
Criteo	AUC(%)	81.49	81.49	81.42	81.41
	Logloss(%)	43.7	43.72	43.77	43.79
Avazu	AUC(%)	76.66	76.61	76.57	76.52
	Logloss(%)	36.58	36.59	36.66	36.71
MovieLens	AUC(%)	97.2	97.23	96.98	97.1
	Logloss(%)	19.66	20.92	21.07	20.91
Frappe	AUC(%)	98.61	98.69	98.47	98.45
	Logloss(%)	14.87	16.12	14.71	15.22

图 5. 模型 FinalMLP 与 DualMLP 复现实验结果与论文实验结果对比表

图 6(a) 展示了不同数据集上 DualMLP 的 AUC 比较，而图 6(b) 则展示了相应的 Logloss 比较。通过这些对比柱状图，我们可以直观地观察到模型在各个数据集上的性能差异。

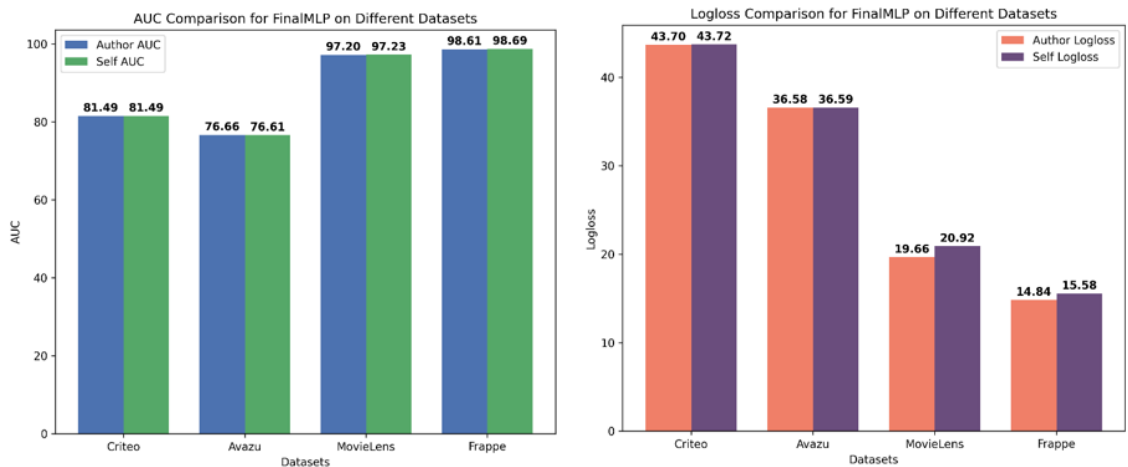


(a) 不同数据集上 DualMLP 的 AUC 比较

(b) 不同数据集上 DualMLP 的 Logloss 比较

图 6. DualMLP 在不同数据集上的性能比较。

图 7(a) 展示了 FinalMLP 在不同数据集上的 AUC 比较，而图 7(b) 展示了 Logloss 的比较。这些对比柱状图清晰地揭示了 FinalMLP 模型在各个数据集上的优劣。



(a) 不同数据集上 FinalMLP 的 AUC 比较

(b) 不同数据集上 FinalMLP 的 Logloss 比较

图 7. FinalMLP 在不同数据集上的性能比较。

5.2.2 复现实验结果分析

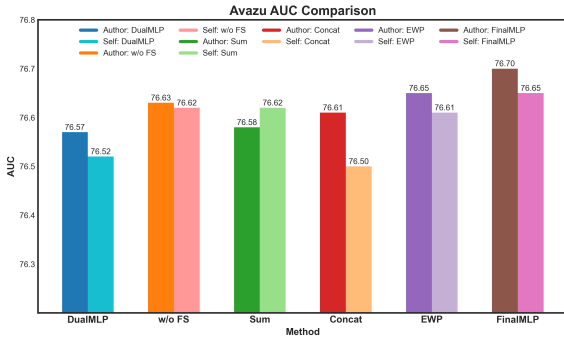
从整体上看，复现的实验结果与原论文中报告的结果展现出高度的一致性。在 Criteo 数据集上，AUC 值和 Logloss 均与论文结果相差无几，显示出实验方法的有效性。在 Avazu 数据集上，尽管 AUC 值与论文结果相近，但 Logloss 略高，这可能暗示模型在该数据集上的表现有待提升。对于 MovieLens 数据集，复现结果显示出略高的 AUC 值，这可能是由于实验条件或数据预处理的差异。在 Frappe 数据集上，AUC 值与论文结果非常接近，而 Logloss 的轻微差异可能需要进一步分析以确定其原因。复现实验在大多数情况下成功地再现了论文中

的结果，这验证了实验的准确性，并为模型性能提供了深入理解。然而，存在的细微差异可能源于实验设置、数据集版本或模型实现的细微差别。未来的研究可以致力于进一步缩小这些差异，并探索可能的改进方向。

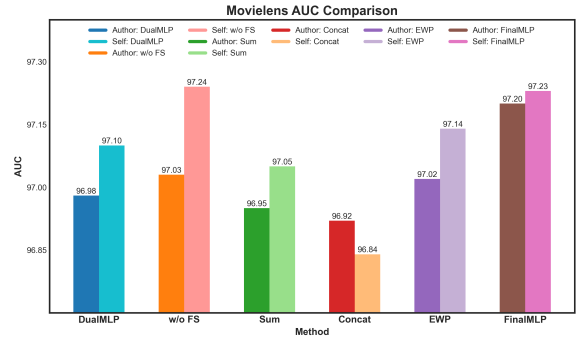
5.3 消融实验结果

消融实验的目的是评估模型中不同组件对整体性能的影响。通过比较 FinalMLP 模型与其变体，包括 DualMLP、w/o FS（无特征选择）、Sum、Concat 和 EWP（其他融合操作），可以确定哪些模块对模型性能至关重要。实验结果表明，移除特征选择模块或用其他融合操作替换双线性融合会导致性能下降。这验证了特征选择和双线性融合模块的有效性。特别是，双线性融合模块比特征选择模块更为关键，因为替换前者会导致更大的性能下降。

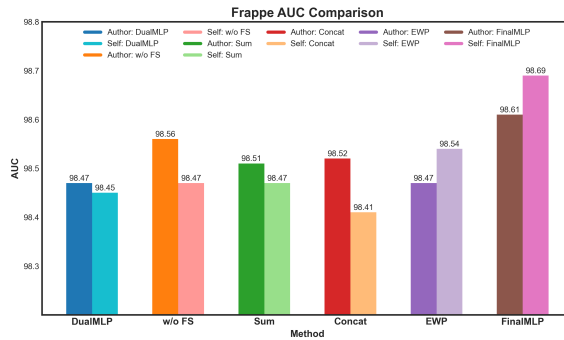
如图 8 所示，对 FinalMLP 模型及其变体在三个不同数据集上的 AUC 性能进行了对比分析。图 (a) 展示了在 Avazu 数据集上的 AUC 比较，图 (b) 展示了在 MovieLens 数据集上的 AUC 比较，而图 (c) 则展示了在 Frappe 数据集上的 AUC 比较。这些柱状图清晰地揭示了特征选择和双线性融合模块在提升模型预测准确性方面的关键作用。



(a) 数据集 Avazu 上的 AUC 对比柱状图



(b) 数据集 MovieLens 上的 AUC 对比柱状图



(c) 数据集 Frappe 上的 AUC 对比柱状图

图 8. 不同数据集上的 AUC 对比分析

5.4 改进实验结果

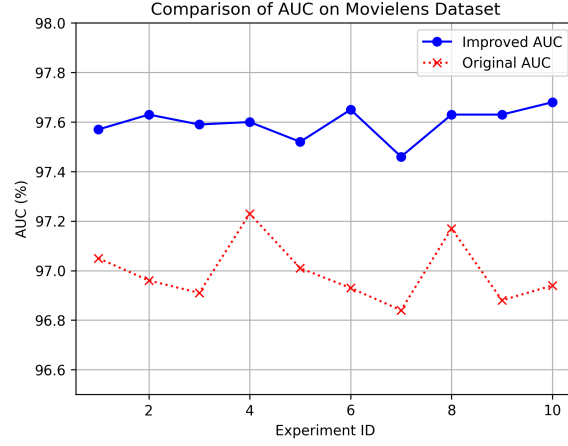
5.4.1 改进实验结果与原实验结果对比

如图 9 所示，展示了在 MovieLens 数据集上进行的 10 次实验中，改进后的 FinalMLP 模型与原始模型的 AUC 性能对比。图 9(a) 详细列出了每次实验的 AUC 值，而图 9(b) 则通过折线图直观地展示了改进模型与原始模型在 AUC 上的表现差异。

在数据集MovieLens上10次实验改进与原结果对比

ID	001	002	003	004	005	006	007	008	009	010	Avg(%)	Std(%)
AUC(%)	97.57	97.63	97.59	97.60	97.52	97.65	97.46	97.63	97.63	<u>97.68</u>	97.53	0.15
	97.05	96.96	96.91	97.23	97.01	96.93	96.84	97.17	96.88	96.94	96.96	0.13

(a) 在 MovieLens 数据集上 10 次实验改进与原结果对比



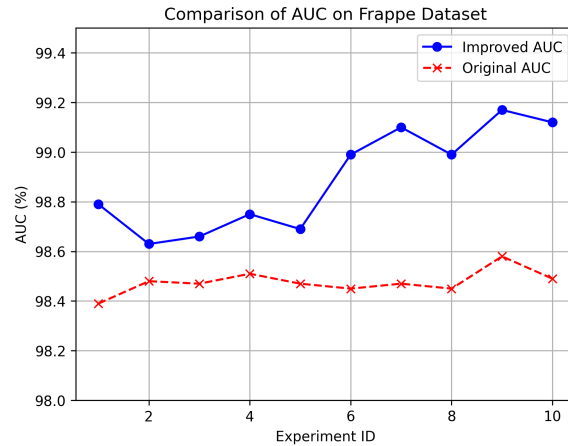
(b) 在 MovieLens 数据集上改进模型与原始模型的 AUC 对比

图 9. MovieLens 数据集上改进模型与原始模型的 AUC 性能对比。

在数据集Frappe上10次实验改进与原结果对比

ID	001	002	003	004	005	006	007	008	009	010	Avg(%)	Std(%)
AUC(%)	98.79	98.63	98.66	98.75	98.69	98.99	99.10	98.99	<u>99.17</u>	99.12	98.86	0.21
	98.39	98.48	98.47	<u>98.51</u>	98.47	98.45	98.47	98.45	98.58	98.49	98.47	0.15

(a) 在 Frappe 数据集上 10 次实验改进与原结果对比



(b) 在 Frappe 数据集上改进模型与原始模型的 AUC 对比

图 10. Frappe 数据集上改进模型与原始模型的 AUC 性能对比。

如图 10 所示，展示了在 Frappe 数据集上进行的 10 次实验中，改进后的 FinalMLP 模型与原始模型的 AUC 性能对比。图 10(a) 详细列出了每次实验的 AUC 值，而图 10(b) 则通

过折线图直观地展示了改进模型与原始模型在 AUC 上的表现差异。

5.5 改进实验结果分析

综合分析表明,本研究中提出的创新特征工程方法在预测准确性方面取得了显著提升。通过对特征进行文本化处理并利用 BERT 模型进行深度语义嵌入,模型能够更有效地捕捉数据中的复杂关系和细微差别。这不仅提高了模型对数据的理解和处理能力,而且在 MovieLens 和 Frappe 数据集上的实验结果也验证了其稳定性和泛化能力。改进后的模型在 AUC 指标上的表现均优于原始模型,平均提升幅度达到了 0.57%,标准差也有所降低,这表明新方法在提高预测性能的同时,也增强了模型的稳定性。这些结果充分证明了特征工程创新方法的有效性,为深度学习在推荐系统中的应用提供了新的视角和可能性。

6 总结与展望

6.1 实验不足之处

在本研究中,我们的特征工程方法在小数据集如 MovieLens 和 Frappe 上运行是可以接受的,尽管运行时间有所增加。然而,在处理 Avazu 和 Criteo 这样的大规模数据集时,我们遇到了显著的挑战。这些大数据集的处理不仅需要显著更多的时间,而且由于内存限制,实验无法成功运行。如图 11 所示,是实验中遇到的报错内容,具体表现为在尝试为特征嵌入分配高达 81.0 GiB 的内存时,出现了 MemoryError 错误。这个错误表明,尽管我们的方法在小数据集上有效,但在处理大规模数据集时,需要更多的计算资源和优化策略。

```
Traceback (most recent call last):
  File "run_expid.py", line 42, in <module>
    build_dataset(feature_encoder, **params)
  File "D:\anaconda3\envs\finmalbert\lib\site-packages\fuxictr\preprocess\build_dataset.py", line 113, in build_dataset
    feature_encoder.fit(train_ddf, **kwargs)
  File "D:\anaconda3\envs\finmalbert\lib\site-packages\fuxictr\preprocess\feature_processor.py", line 130, in fit
    self.fit_sequence_col(col, train_ddf[name].values,
  File "D:\anaconda3\envs\finmalbert\lib\site-packages\fuxictr\preprocess\feature_processor.py", line 268, in fit_sequence_col
    self.processor_dict[name + "::embeddings"] = np.concatenate(embeddings, axis=0)
  File "<__array_function__ internals>", line 200, in concatenate
numpy.core._exceptions.MemoryError: Unable to allocate 81.0 GiB for an array with shape (28300276, 768) and data type float32
```

图 11. 实验中遇到的内存分配错误

6.2 未来研究方向及展望

面对本研究中遇到的挑战,未来的工作将集中在以下几个关键的研究方向:

- **内存优化与高效处理:** 开发策略以降低大规模数据处理的内存需求,可能包括采用更高效的数据结构和算法,以及分批处理和内存映射技术。
- **改进嵌入模型:** 探索对现有嵌入模型的修改和改进,以提高其在大数据集上的性能和效率。这可能涉及到模型架构的调整,或者采用新的训练策略来优化模型的内存占用和计算效率。

- **特征融合技术**: 进一步研究和开发新的特征融合技术, 以更有效地整合不同来源的特征, 提升模型对数据的理解和预测能力。
- **计算资源扩展**: 考虑使用分布式计算或云计算资源, 以克服单机处理能力的限制, 实现对更大数据集的有效处理。

通过这些方向的研究, 期待能够进一步提升模型的性能, 使其在更广泛的应用场景中发挥更大的作用。

参考文献

- [1] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, and Hemal Shah. Wide & deep learning for recommender systems. *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems*, 1(1):7–10, 2016.
- [2] W. Cheng, Y. Shen, and L. Huang. Adaptive factorization network: Learning adaptive-order feature interactions. *The Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*, 34(03):3609–3616, 2020.
- [3] H. Guo, R. Tang, Y. Ye, Z. Li, and X. He. Deepfm: A factorization-machine based neural network for ctr prediction. *International Joint Conference on Artificial Intelligence (IJCAI)*, IJCAI(01):1725–1731, 2017.
- [4] J. Lian, X. Zhou, F. Zhang, and et al. xdeepfm: Combining explicit and implicit feature interactions for recommender systems. *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, 24(01):1754–1763, 2018.
- [5] T. Lin, A. RoyChowdhury, and S. Maji. Bilinear cnn models for fine-grained visual recognition. *IEEE International Conference on Computer Vision (ICCV)*, ICCV(01):1449–1457, 2015.
- [6] J. Ma, Z. Zhao, X. Yi, J. Chen, L. Hong, and E. H. Chi. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, 24(01):1930–1939, 2018.
- [7] S. Rendle. Factorization machines. *Proceedings of the 10th IEEE International Conference on Data Mining (ICDM)*, 10(01):995–1000, 2010.
- [8] W. Song, C. Shi, Z. Xiao, Z. Duan, Y. Xu, M. Zhang, and J. Tang. AutoInt: Automatic feature interaction learning via self-attentive neural networks. *Proceedings of the 28th ACM International Conference on Information and Knowledge Management (CIKM)*, pages 1161–1170, 2019.

- [9] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2017:5998–6008, 2017.
- [10] R. Wang, B. Fu, G. Fu, and M. Wang. Deep & cross network for ad click predictions. *Proceedings of the 11th International Workshop on Data Mining for Online Advertising (ADKDD)*, 11:12:1–12:7, 2017.