

推荐系统

王佳和 张怀松

2024214505, 2024214519

计算机技术, 计算机技术

wang-jh24@mails.tsinghua.edu.cn, thepinezhang@gmail.com,

摘要

推荐系统是现代互联网服务中不可或缺的一部分,尤其在电商、社交网络和内容分发平台中起到关键作用。点击率预测作为推荐系统的核心,直接关系到广告收入和用户满意度。准确的点击率预测可以显著提升个性化推荐的效果,优化用户体验,增加运营效率。对于点击率预测任务,xDeepFM 模型被认为是准确且高效的。我们通过 Criteo 数据集复现 xDeepFM 模型,结合显式和隐式特征交互学习来处理复杂的高阶特征交互。我们的工作涵盖了详尽的文献综述、数据预处理、模型构建和参数优化,使用 AUC 和 Logloss 作为性能评估指标。实验结果表明,尽管与最优配置存在差距,但复现的 xDeepFM 模型在多个指标上展现了较高的预测准确度和改进潜力。报告总结了 xDeepFM 模型在推荐系统中的应用价值,并指出了未来研究的方向,包括进一步优化模型结构和调整超参数。这些工作不仅验证了模型的有效性,还为推荐系统的技术进步提供了有价值的见解。

关键词: 推荐系统, 点击率预测, xDeepFM, 特征交互, 性能评估

摘要

Recommendation systems are an indispensable part of modern Internet services, particularly within e-commerce, social networking, and content distribution platforms. Click-through rate (CTR) pre-

diction, a core component of recommendation systems, is directly linked to advertising revenue and user satisfaction. Accurate CTR predictions can significantly enhance the effectiveness of personalized recommendations, improve user experience, and increase operational efficiency. The xDeepFM model is considered accurate and efficient for CTR prediction tasks. We replicated the xDeepFM model using the Criteo dataset, incorporating both explicit and implicit feature interactions to handle complex high-order feature interactions. Our work included an extensive literature review, data preprocessing, model building, and parameter optimization, using AUC and Logloss as performance indicators. Experimental results show that, despite some gaps with optimal configurations, the replicated xDeepFM model demonstrated high prediction accuracy and potential for improvement on multiple metrics. This report summarizes the application value of the xDeepFM model in recommendation systems and identifies directions for future research, including further optimization of the model structure and adjustment of hyperparameters. This work not only validates the effectiveness of the model but also contributes valuable insights into the technological advancement of recommenda-

tion systems.

Keywords: Recommendation Systems, Click-Through Rate Prediction, xDeepFM, Feature Interaction, Performance Evaluation

1. 引言

如今，移动互联网已经成为人们日常体验中不可或缺的一部分。互联网上不断产生新的数据和信息，导致用户被海量的内容淹没，难以快速找到自己感兴趣的内容。为了缓解互联网产品中的信息过载，满足用户多样化的在线服务需求（如电商、短视频、新闻等），个性化的推荐系统变得越来越重要 [1]。这些系统通过分析用户的历史行为和偏好来预测用户可能感兴趣的项目，从而极大地优化了信息过载问题并提高了用户体验和运营效率 [2]。尤其在电子商务领域，推荐系统不仅帮助消费者发现新产品，还通过个性化推荐提高了购买效率和顾客满意度 [3]。例如，在 2020 年后由于一些限制，随着居家时间的增加，人们越来越倾向于在线购物，这使得电子商务网站必须更有效地利用用户数据来推动销售 [4]。如今，推荐系统已经在各个领域得到了广泛的应用，包括电子商务、社交网络、信息分发，显著地影响着人们的生活 [5]。

点击率预测（CTR）是推荐系统中的一项核心任务 [6, 7]，它直接关联到广告收入和内容推荐的效果。CTR 预测 [8] 在推荐系统中得到广泛应用，其主要目标是估计用户对推荐内容点击的概率，这对于优化广告展示和增加商业收益至关重要。更精准的 CTR 预测可以显著提升用户满意度，因为它能够确保用户更频繁地看到他们感兴趣的内容。此外，CTR 预测帮助平台实现个性化服务，通过分析用户历史行为数据和当前上下文环境，系统能更有效地向用户推荐内容 [9]。在技术层面，CTR 预测模型通过集成多种数据输入（如用户资料、项目属性和场景信息）来实现这一目标。这些模型利用深度学习、注意力机制等高级算法来处理和解析这些数据，从而提高预测的准确性和相关性 [10]。因此，CTR 预测不仅是连接用户与他们感兴趣内容的桥梁，也是

提升广告效率和内容推荐质量的关键技术。这些研究表明，深入理解和精确实现 CTR 预测对于推荐系统的发展至关重要。

本文将对 Wide & Deep、DeepFM 和 xDeepFM 三篇经典点击率预测模型的论文进行文献综述与算法总结，深入分析各模型的核心算法思想和技术架构。同时，本文将对 xDeepFM 模型进行复现，详细介绍复现过程中涉及的评价指标（如 AUC 和 Logloss）及其对模型预测效果的衡量方式，说明所使用的数据集的来源、规模及其预处理方法，以确保数据符合实验要求。在复现流程方面，将完整记录数据预处理、模型构建、参数选择、训练与优化的具体步骤，以保证实验的科学性和可重复性。接着，本文将展示复现的实验结果，通过比较 xDeepFM 和其他模型在不同评价指标上的表现，分析 xDeepFM 的优势和不足，并探讨可能的改进方向。最后，本文总结复现研究的整体效果与发现，为推荐系统中点击率预测模型的研究提供新的见解与实践参考。

2. 文献综述

本部分主要介绍了三篇关键论文，分别是《Wide & Deep Learning for Recommender Systems》[11]、《DeepFM: A Factorization-Machine based Neural Network for CTR Prediction》[12]和《xDeepFM: Combining Explicit and Implicit Feature Interactions for Recommender Systems》[13]。

2.1. Wide & Deep Learning for Recommender Systems

在大规模推荐系统中，传统的广义线性模型（如逻辑回归）虽然简单、可扩展且解释性强，但它主要依赖手工设计的交叉特征（cross-product features）来捕获低阶特征交互。因此，它擅长记忆历史数据中的频繁特征组合，但在遇到未见过的新特征组合时，往往缺乏泛化能力。

另一方面，深度神经网络（DNN）则通过学习嵌入特征向量的低维表示，可以在无需特征工程的情况下自动生成高阶特征交互，从而具有更好的泛化能力。然而，由于 DNN 倾向于过度泛化，在稀疏

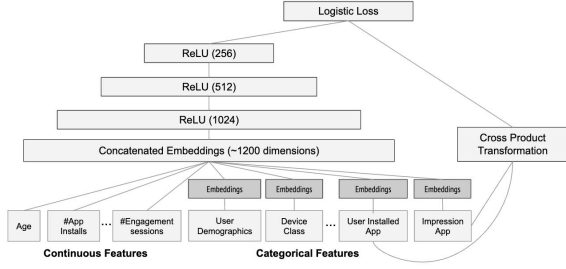


图 1. Wide & Deep Learning 框架图

和高秩（high-rank）数据集上，它可能会推荐与用户兴趣相关性较低的物品。因此，如何在推荐系统中平衡记忆和泛化这两种能力，成为推荐系统面临的核心问题。

为了同时实现记忆和泛化，Cheng 等人提出了 Wide & Deep Learning 框架，将“宽”模型（Wide）与“深”模型（Deep）结合，通过联合训练（joint training）这两部分来提升推荐系统的预测效果。对于 Wide 部分，论文使用广义线性模型，通过手工设计的特征交叉来捕获低阶特征交互，进而实现记忆。这部分擅长捕获频繁出现的特征组合（例如，用户已经安装的 App 和当前展示的 App 的组合），适合那些在历史数据中出现过的交互特征。而对于 Deep 部分，论文使用深度神经网络，通过嵌入和多层非线性层的方式自动学习高阶特征交互，从而实现泛化。与 Wide 部分不同，Deep 部分无需特征工程，适合捕获那些在历史数据中未见过的新特征组合。通过联合训练，Wide 部分和 Deep 部分在训练过程中共享参数和损失函数。联合训练不仅避免了模型在训练时相互独立导致的资源浪费，还可以在优化时相互补充，实现记忆和泛化的平衡。

图 1 展示了 Wide& Deep 模型的整体结构，其主要包括以下几个部分：输入层：模型的输入由用户特征、上下文特征和展示物品的特征组成。这些特征通过不同的预处理方法（如数值特征的标准化和类别特征的嵌入）转换为模型可接受的输入格式。Wide 部分：Wide 部分是一个线性模型，利用特征交叉来捕获低阶特征交互。Wide 部分的输出是一个线性组合，即：

$$y_{\text{wide}} = \mathbf{w}^T [\mathbf{x}, \phi(\mathbf{x})] + b \quad (1)$$

其中， \mathbf{x} 表示输入特征， $\phi(\mathbf{x})$ 表示特征交叉变换， \mathbf{w} 和 b 分别是权重和偏置。Wide 部分对特征交叉的依赖使其能够记忆频繁出现的特征组合，以帮助模型更精确地对历史频率较高的组合进行推荐。Deep 部分是一个多层神经网络，输入为所有特征的低维嵌入表示。每个类别特征首先通过嵌入层转化为低维稠密向量，然后所有嵌入向量和数值特征一起被连接为一个稠密特征向量。这个稠密向量经过多层全连接层和激活函数（如 ReLU）非线性变换，逐步学习更高阶的特征交互，以实现对新特征组合的泛化。Deep 部分的输出即为预测值，用于与 Wide 部分的输出一同进行融合。融合与输出：Wide & Deep 模型通过加权的方式将 Wide 部分和 Deep 部分的输出结合在一起。具体的输出表示为：

$$\hat{y} = \sigma(y_{\text{wide}} + y_{\text{deep}}) \quad (2)$$

其中， σ 是 sigmoid 激活函数，用于二分类的点击率预测任务。联合训练的损失函数是二分类的交叉熵损失，通过反向传播对整个模型的参数进行优化。

通过这一架构，Wide 部分实现了对已见特征组合的记忆，而 Deep 部分则提供了对未见特征组合的泛化能力，从而提升了推荐系统的整体性能。

2.2. DeepFM: A Factorization-Machine based Neural Network for CTR Prediction

Guo 等人则主要探讨了如何在推荐系统中进一步提升点击率（CTR）预测的效果，并解决了 Wide & Deep 模型的一些局限性。推荐系统需要捕捉用户和物品之间复杂的特征交互，但传统的线性模型和因子分解机（Factorization Machines, FM）虽然在捕捉低阶特征交互方面表现良好，却无法有效学习高阶交互。同时，Wide & Deep 模型虽然结合了线性模型和深度神经网络（DNN），但它的 Wide 部分依赖手工设计的特征交叉，费时且难以扩展，尤其是在大规模应用中。

为了解决这些问题，DeepFM 提出了一种无须手工特征工程的端到端学习方法，通过将 FM 和 DNN 集成到一个模型中，同时捕捉低阶和高阶特征

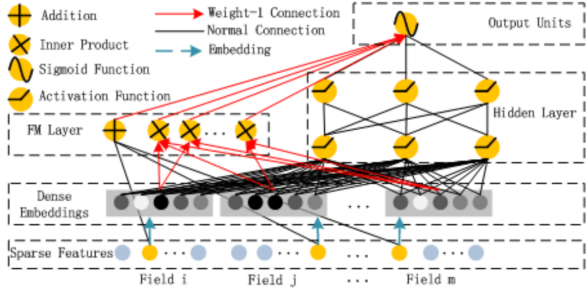


图 2. DeepFM 框架图

交互。具体来说，DeepFM 利用 FM 部分来学习低阶特征交互，特别是二阶交互，通过内积操作建模特征间的相互关系。FM 部分的输出可以表示为：

$$y_{FM} = \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=i+1}^d \langle v_i, v_j \rangle x_i x_j \quad (3)$$

其中， w_i 为特征的线性权重， v_i 为特征的嵌入向量。

而 DNN 部分则用于学习高阶特征交互，它将所有输入特征（包括数值特征和嵌入后的类别特征）连接成一个稠密向量，通过多层全连接网络逐步捕捉更高阶的特征组合。在 DeepFM 中，FM 和 DNN 共享同一个嵌入层，因此整个模型可以端到端训练，不再需要额外的手工特征工程。

DeepFM 模型将 FM 部分和 DNN 部分的输出结合起来，通过 sigmoid 函数生成最终的 CTR 预测值。具体的输出可以表示为：

$$\hat{y} = \sigma(y_{FM} + y_{DNN}) \quad (4)$$

其中， σ 为 sigmoid 激活函数，用于二分类的点击率预测任务

图2展示了 DeepFM 的结构，首先输入层接收类别特征和数值特征。类别特征通过嵌入转化为稠密向量，而数值特征则可以直接输入。FM 部分接收嵌入后的输入，通过计算特征嵌入向量的内积捕捉二阶特征交互，并通过加权求和生成低阶交互输出。DNN 部分则将所有输入特征连接成一个高维稠密向量，经过多层全连接层和激活函数的非线性转换，逐层捕捉更复杂的高阶特征交互。最终，模型将 FM 部分和 DNN 部分的输出进行融合，生成点击率预测值。

DeepFM 在模型的设计上实现了无须手工特征工程的目标，它不但保留了 FM 的低阶特征学习能力，也利用 DNN 有效学习高阶特征交互，从而提高了推荐系统在复杂特征组合下的预测能力。这种架构使得 DeepFM 在多个 CTR 预测任务中表现优越，并且具备较高的训练和推断效率。模型的架构图展示了 DeepFM 的结构，首先输入层接收类别特征和数值特征。类别特征通过嵌入转化为稠密向量，而数值特征则可以直接输入。FM 部分接收嵌入后的输入，通过计算特征嵌入向量的内积捕捉二阶特征交互，并通过加权求和生成低阶交互输出。DNN 部分则将所有输入特征连接成一个高维稠密向量，经过多层全连接层和激活函数的非线性转换，逐层捕捉更复杂的高阶特征交互。最终，模型将 FM 部分和 DNN 部分的输出进行融合，生成点击率预测值。

DeepFM 在模型的设计上实现了无须手工特征工程的目标，它不但保留了 FM 的低阶特征学习能力，也利用 DNN 有效学习高阶特征交互，从而提高了推荐系统在复杂特征组合下的预测能力。这种架构使得 DeepFM 在多个 CTR 预测任务中表现优越，并且具备较高的训练和推断效率。

2.3. xDeepFM: Combining Explicit and Implicit Feature Interactions for Recommender Systems

Lian 等人提出了一种名为 xDeepFM 的模型，进一步扩展了 DeepFM 模型的结构，使其能够显式地学习高阶特征交互。该论文的主要问题是现有方法在特征交互建模方面的局限性：传统的 FM 模型主要捕捉低阶交互，而 DNN 则隐式地捕提高阶交互，但难以明确控制交互的阶数。尽管像 DeepFM 和 Wide & Deep 这样的模型可以通过混合低阶和高阶特征学习来增强推荐效果，但 DNN 隐式捕捉的高阶交互往往是在位级别（bit-wise）的，而不是在向量级别（vector-wise），这可能会导致部分交互信息丢失或噪声干扰。

为了解决这些问题，xDeepFM 提出了一个新的子模块，称为压缩交互网络（Compressed Interaction Network, CIN），该网络专门用于显式建模高阶特征

交互。与传统 DNN 在位级别的交互不同，CIN 在向量级别进行交互，能够捕捉到更具代表性和准确性的特征交互。CIN 的设计灵感来源于卷积神经网络 (CNN) 中的卷积层，通过压缩处理跨层交互，逐层构建出更高阶的特征交互，且计算复杂度没有随之大幅提升。CIN 的核心思想是通过一层层的交互，逐步增加特征交互的阶数，最终构建出一个包含高阶交互的压缩特征表示，能够明确、逐步地控制和增加特征的交互程度。

在 xDeepFM 模型中，CIN 模块负责生成显式高阶交互，而 DNN 部分则继续捕捉隐式高阶交互，这使得 xDeepFM 在模型中同时包含显式和隐式的交互学习能力。具体来说，xDeepFM 的架构包含三部分：线性部分、CIN 部分和 DNN 部分。线性部分用于低阶特征交互的学习，CIN 用于显式的高阶交互建模，而 DNN 则负责捕捉隐式的高阶交互。这种组合方式使得模型在捕捉复杂特征关系时能够更加灵活，既可以处理数据中显著的特征交互，也能通过 DNN 学习到较难检测的隐式关系。

xDeepFM 的输出融合了线性部分、CIN 部分和 DNN 部分的结果。模型首先对每一部分的输出进行加权求和，之后再通过 sigmoid 函数进行激活，以生成最终的预测结果。在 CTR 预测任务中，xDeepFM 不仅显著提高了模型的预测性能，还在多个真实数据集上表现出超越其他主流模型的效果。

通过这种架构设计，xDeepFM 模型在捕捉特征交互上更加全面和高效，不再局限于传统的位级交互，而是能够显式建模和控制高阶交互，并在推荐系统的 CTR 预测任务中展示出卓越的性能。

3. 算法总结

推荐任务 (Recommendation) 的核心是根据用户的历史行为和特征来预测用户对特定物品的偏好或点击率。该任务可形式化地定义为一个概率预测问题，其目标是最大化用户与物品的交互概率。在推荐系统中，给定用户集合 $U = \{u_1, u_2, \dots, u_m\}$ 和物品集合 $I = \{i_1, i_2, \dots, i_n\}$ ，以及已知的用户-物品交互矩阵 R ，其中 $R_{u,i}$ 表示用户 u 对物品 i 的已知偏好或交互（若无交互则为零或未定义），目标是学

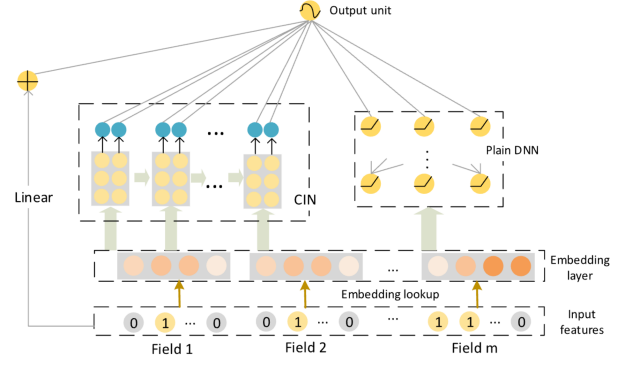


图 3. xDeepFM 框架图

习一个预测函数 $f: U \times I \rightarrow \mathbb{R}$ ，用于预测用户对未交互物品的偏好程度 $\hat{R}_{u,i} = f(u, i)$ ；然后，基于这些预测，为每个用户 u 从未交互的物品中生成一个按预测偏好排序的推荐列表 L_u ，即选择得分最高的前 K 个物品： $L_u = \text{TopK}(\{\hat{R}_{u,i} | i \in I, R_{u,i} \text{ 未知}\})$ 。对于要复现的 xDeepFM 来说，xDeepFM 主要由以下三个模块组成：线性模块、压缩交互网络 (Compressed Interaction Network, CIN) 模块和深度神经网络 (DNN) 模块。

线性模块负责学习一阶（低阶）特征交互，主要用于捕捉直接特征的独立贡献。它是一种广义线性模型，用于学习特征的线性组合权重。在推荐系统中，线性模块通常用于捕捉频繁出现的显著特征，能够有效地建模简单且重要的特征关系。假设输入特征为 \mathbf{x} ，线性模块的输出可表示为：

$$y_{\text{linear}} = \mathbf{w}^T \mathbf{x} + b, \quad (5)$$

其中 \mathbf{w} 为特征权重， b 为偏置项。该模块直接生成推荐预测的初步结果，通过一阶权重学习了各个特征的独立贡献。

CIN 模块则用于显式建模高阶特征交互，逐层增加交互的阶数，能够明确、逐步地控制和增加特征交互的复杂性。CIN 的灵感来源于卷积神经网络 (CNN)，通过逐层交互和特征压缩来构建高阶特征的显式交互。在 CIN 中，高阶交互以向量级 (vector-wise) 进行，即特征嵌入的整个向量参与交互，这使得模型能够更准确地表达特征交互关系，而不只是位级交互。CIN 通过一种特殊的外积运算构建特征

组合，使得特征交互的阶数随着网络深度增加而逐渐提升。CIN 的核心计算方式是通过嵌入矩阵的外积来生成高阶交互特征。假设输入特征的嵌入矩阵为 $X^0 \in \mathbb{R}^{m \times d}$ ，其中 m 是特征的数量， D 是嵌入维度。CIN 在第 k 层的输出矩阵 X^k 由以下方式生成：

$$X_{h,*}^k = \sum_{i=1}^{H_{k-1}} \sum_{j=1}^m W_{ij}^{k,h} (X_{i,*}^{k-1} \circ X_{j,*}^0) \quad (6)$$

其中 $1 \leq h \leq H_k$, $W_{i,j}^{k,h} \in \mathbb{R}^{H_{k-1} \times m}$ 为第 h 个特征向量的参数矩阵， \circ 表示 Hadamard 积（元素乘积），而 X_j^0 是初始嵌入矩阵的第 j 个特征向量。通过该方式，CIN 模块能够在高阶交互中不断生成新的显式特征组合，并通过层次加深逐步增加交互的复杂度。第三个模块 DNN 模块用于捕捉隐式的高阶特征交互，通过非线性转换学习到复杂的特征关系，从而生成更丰富的特征表达。DNN 模块的主要目的是捕捉那些难以通过显式交互提取的特征关系。DNN 通过嵌入层和多层全连接网络实现非线性映射，可以有效学习到特征之间隐含的高阶交互。与 CIN 不同的是，DNN 并不明确控制交互阶数，而是通过多层网络自动捕捉和组合特征信息。DNN 模块首先将输入的特征嵌入表示进行向量化，然后通过多层全连接层逐层生成新的特征表达。假设 e 是连接后的嵌入特征，则 DNN 的第 1 层输出可以表示为：

$$a^{(l)} = \sigma(W^{(l)}a^{(l-1)} + b^{(l)}), \quad (7)$$

其中 $W^{(l)}$ 和 $b^{(l)}$ 为第 l 层的权重矩阵和偏置， σ 是激活函数（如 ReLU）。最终 DNN 的输出用于预测和生成隐式高阶特征交互的表达。

模块融合与输出

xDeepFM 模型通过整合线性模块、CIN 模块和 DNN 模块的输出，生成最终的点击率预测值。模型的整体输出形式可以表示为：

$$\hat{y} = \sigma(w_{linear}^T a + w_{dnn}^T x_{dnn}^k + w_{cin}^T p^+ + b) \quad (8)$$

其中， σ 为 sigmoid 函数， a 为原始特征。 x_{dnn}^k 和 p^+ 分别表示 DNN 和 CIN 的输出。 w^* 和 b 可为

学习参数。对于二分类问题，损失函数为对数损失：

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i) \quad (9)$$

其中 N 是训练集的总数，优化过程是极小化以下目标函数：

$$\mathcal{J} = \mathcal{L} + \lambda_* ||\Theta|| \quad (10)$$

其中 λ_* 代表正则化项， Θ 代表三个部分的参数集合。通过这种融合方式，xDeepFM 能够在不依赖手工特征工程的情况下全面捕捉推荐任务中的低阶和高阶特征交互，并实现显式和隐式特征组合的有效学习。最终，xDeepFM 在 CTR 预测任务中展示出比传统模型更高的预测准确性。

4. 评价指标

本文采用了两种评价指标来评估模型性能。AUC (Area Under the Curve) 是衡量二分类模型表现的重要指标之一，尤其适用于不平衡数据集。AUC 代表 ROC 曲线下的面积，ROC 曲线通过绘制真正例率 (True Positive Rate) 和假正例率 (False Positive Rate) 之间的关系来反映模型在不同阈值下的分类能力。AUC 值越接近 1，表明模型区分正负样本的能力越强；当 AUC 值为 0.5 时，模型的表现相当于随机猜测，表明分类效果较差；而 AUC 值达到 1 则表示模型具有完美的区分能力。在点击率预测 (CTR 预测) 任务中，AUC 是评价模型在预测用户点击行为时整体表现的关键指标。Logloss (对数损失) 则是一种用于衡量模型预测准确性的指标，通过计算预测值与实际值之间的距离来评价模型性能。Logloss 值越小，说明模型预测越接近真实情况，表现越好。其计算公式基于每个样本预测概率的对数形式之和，当模型预测接近真实标签（即对正样本预测接近 1，对负样本预测接近 0）时，Logloss 值较低；而预测偏离真实标签时，Logloss 值则显著增加。Logloss 强调模型的概率预测精度，因此在 CTR 预测任务中，它能够有效评估模型的稳定性和预测结果的置信度。

5. 数据集

本文采用了 Criteo 数据集，这是点击率预测 (CTR) 任务中广泛应用的数据集之一，由 Criteo Labs 提供，包含数百万条广告点击日志记录。Criteo 数据集特别适用于广告推荐系统的研究，因为它具备多样且真实的广告展示环境信息。数据集中的每条记录代表用户与特定广告的一次交互，提供了丰富的特征信息，包括 13 个连续特征和 26 个类别特征。这些特征涉及用户行为、广告内容以及上下文环境等维度，能够帮助模型更准确地识别用户点击广告的倾向。在每条记录中，Criteo 数据集还包含一个二分类标签，标记该条广告记录是否被用户点击 (点击为 1，未点击为 0)，这是模型预测的目标。该数据集中点击标签的类别不平衡现象非常明显，实际点击的比例较低，这反映了真实广告展示中用户点击行为的稀疏性。Criteo 数据集的高维特征和类别不平衡问题为 CTR 预测模型提供了挑战，同时也使其成为评估模型在真实环境下表现的理想数据集。通过在 Criteo 数据集上训练和测试，模型可以更好地学习和捕捉复杂的用户行为模式，从而提升 CTR 预测的准确性和稳定性。这一数据集在推荐系统和广告投放领域的研究中得到了广泛应用，是验证模型性能的标准选择。

6. 复现流程

在本项目中，我们目标是复现 xDeepFM 模型并评估其在点击率预测任务中的性能。本节详细介绍了实验的全过程，包括数据处理、模型配置、训练和测试步骤。下面以使用百分之五数据的实验为例介绍复现流程，对应文件为 ctr0.05.ipynb。

6.1. 数据预处理

使用 Criteo 广告点击数据集，该数据集包含类别特征和数值特征。数据预处理的目的是提高数据质量，确保模型可以有效学习。

1. **数据加载**：使用 pandas 库从 train.txt 文件中以块方式读取数据。选择随机 5% 的数据用于训练，以减少内存占用并加快初步测试的速度。

```
train_data_chunks = pd.read_csv('./oridata/train.txt', names=col_names, sep='\t', chunksize=30000)
data = pd.concat(chunk.sample(frac=train_data_fraction, random_state=42) for chunk in tqdm(train_data_chunks, desc="Loading train data"))
```

2. **缺失值处理**：对于类别特征，将缺失值填充为 '-1'，对于数值特征，缺失值填充为 '0'。

```
data[sparse_features] = data[sparse_features].fillna('-1',)
data[dense_features] = data[dense_features].fillna(0,)
```

3. **特征编码**：对类别特征进行标签编码，对数值特征实施最小-最大归一化，以标准化输入特征。

```
for feat in sparse_features:
    lbe = LabelEncoder()
    data[feat] = lbe.fit_transform(data[feat])
mms = MinMaxScaler(feature_range=(0, 1))
data[dense_features] = mms.fit_transform(data[dense_features])
```

6.2. 模型配置

我们采用了 deepctr_torch 库中的 xDeepFM 模型，该模型集成了线性回归和深度神经网络，能够处理复杂的特征交互。

1. **特征列定义**：确定每个特征的处理方式，包括为每个稀疏特征指定嵌入维度和为每个密集特征指定维度为 1。

```
fixlen_feature_columns = [SparseFeat(feat, vocabulary_size=data[feat].max() + 1, embedding_dim=4)
                           for feat in sparse_features] + [DenseFeat(feat, 1)]
```

```

                                for feat
                                in dense_features]
feature_names = get_feature_names(
    linear_feature_columns +
    dnn_feature_columns)

```

2. **模型实例化**: 配置 xDeepFM 模型, 指定设备 (CPU 或 CUDA), 并设置嵌入的 L2 正则化。

```

model = xDeepFM(linear_feature_columns
    =linear_feature_columns,
    dnn_feature_columns=
    dnn_feature_columns,
    task='binary',
    l2_reg_embedding=1e-5, device=device
)

```

6.3. 模型训练与测试

1. **数据划分**: 将处理后的数据分为 80% 的训练集和 20% 的测试集, 以验证模型性能。

```

train, test = train_test_split(data,
    test_size=0.2, random_state=2020)

```

2. **模型训练**: 配置训练参数, 如批量大小和训练周期, 并在训练集上进行模型训练, 同时设置验证分割以评估过拟合情况。

```

history = model.fit(train_model_input,
    train[target].values, batch_size
    =4096, epochs=10, verbose=1,
    validation_split=0.2)

```

3. **模型测试**: 在测试集上使用模型进行预测, 并计算 LogLoss 和 AUC 指标, 以评估模型的预测性能。

```

pred_ans = model.predict(
    test_model_input, 256)
print("test LogLoss", round(log_loss(
    test[target].values, pred_ans), 4))
print("test AUC", round(roc_auc_score(
    test[target].values, pred_ans), 4))

```

7. 实验结果及分析

实验结果表明, **我们实现的 xDeepFM Ours 模型**在点击率预测任务中的性能仅优于基准模型 (如 LR 和 FM), 但相比其他方法 (如 DNN、DCN 和标准 xDeepFM) 仍有显著差距。这一结果表明, 当前实现的模型在高阶特征交互建模及整体优化上仍有较大的提升空间。

从与传统基准模型比较来看, xDeepFM Ours 展现出了一定的性能优势。例如, 在 100% 数据下, AUC 达到 0.7928, Logloss 降至 0.4572, 明显优于 LR (AUC 0.7577, Logloss 0.4854) 和 FM (AUC 0.7900, Logloss 0.4592)。这表明, xDeepFM Ours 能够在一定程度上捕捉特征交互信息, 从而提升预测性能。然而, 相比更先进的模型如 DNN (AUC 0.7993, Logloss 0.4491) 和标准 xDeepFM (AUC 0.8052, Logloss 0.4418), 我们的模型仍然存在显著性能差距。这种差距说明当前实现未能完全挖掘模型的潜力, 尤其是在处理复杂特征交互时的能力不足。

进一步分析表明, xDeepFM Ours 的性能随着数据使用比例的增加而逐步提升, 但在不同数据规模下仍然表现出一定的局限性。在 5% 数据下, AUC 为 0.7229, Logloss 为 0.6411, 尽管比基准模型略有提升, 但仍远低于标准 xDeepFM 的预期性能。在 50% 数据下, AUC 提升至 0.7520, Logloss 降至 0.5329, 但相较于其他方法 (如 DCN 和 Wide&Deep), 提升幅度依然有限。而在 100% 数据下, 尽管 AUC 和 Logloss 达到最佳值 (AUC 0.7928, Logloss 0.4572), 但依然无法匹敌 DNN 和标准 xDeepFM。这种结果表明, xDeepFM Ours 在面对不同数据规模时的鲁棒性仍需进一步改进。

导致这一现象的原因可能与以下几个方面有关。首先, CIN 模块作为 xDeepFM 的核心组件, 用于显式建模高阶特征交互, 但我们的实现可能在层数、单元数或正则化策略上未达到最优, 导致模型未能充分捕捉特征之间的复杂关系。其次, DNN 部分的设计存在显著不足。我们的模型仅包含 2 层 DNN, 而标准 xDeepFM 通常设计为 3 层 DNN, 导致模型

Model Name	Data Usage (%)	AUC	Logloss	Depth
LR	-	0.7577	0.4854	-, -
FM	-	0.7900	0.4592	-, -
DNN	-	0.7993	0.4491	-, 2
DCN	-	0.8026	0.4467	2, 2
Wide&Deep	-	0.8000	0.4490	-, -
PNN	-	0.8038	0.4927	-, -
DeepFM	-	0.8025	0.4468	-, 2
xDeepFM	-	0.8052	0.4418	3, 2
xDeepFM (Ours)	5	0.7229	0.6411	2, 2
xDeepFM (Ours)	50	0.7520	0.5329	2, 2
xDeepFM (Ours)	100	0.7928	0.7928	2, 2

表 1. Performance comparison of different models on the Criteo dataset. The proposed xDeepFM (Ours) is evaluated with varying data usage.

在隐式特征关系的建模能力上存在不足。此外，隐藏层的结构设计和激活函数选择可能也不够优化，进一步限制了模型对复杂特征的代表能力。最后，超参数的选择也可能对模型性能产生影响。我们在实验中可能未能进行足够细致的超参数调优，例如学习率、嵌入维度和 Dropout 等参数的配置未达到最佳，进而限制了模型的性能发挥。

除了模型结构和超参数调优的问题，数据预处理也可能是影响性能的重要因素。推荐系统中的数据预处理通常对结果有着重要的影响，例如类别特征的编码方式（如 One-Hot 编码或 Embedding）和数值特征的归一化策略。如果这些预处理方法未能充分匹配数据的分布特性，可能会影响模型输入的质量，从而削弱模型对数据的学习能力。此外，实验中的训练资源限制也可能是导致结果偏差的因素之一。例如，硬件资源的限制可能导致我们在训练轮次、批次大小或模型复杂度方面做出妥协，进而影响了模型的最终性能。值得注意的是，我们的模型设计中 DNN 层数较少（2 层），这不仅影响了模型的表达能力，还增加了特征学习的不充分性，从而进一步放大了资源限制带来的问题。

尽管当前实验结果未能达到预期，但我们的模型仍表现出一定的潜力。通过优化 CIN 和 DNN 部分的设计，例如增加 CIN 的层数、优化隐藏层结构

以及引入更适合的激活函数，可以进一步提升模型的特征交互建模能力。同时，针对超参数进行更细致的调优，例如使用网格搜索或随机搜索方法，优化嵌入维度、学习率和正则化参数，可能会显著改善模型性能。此外，改进数据预处理策略，如对类别特征采用更加适配的编码方式，或对数值特征进行更精细的归一化处理，也能够提高模型的学习效果。

总体而言，我们实现的 xDeepFM Ours 模型在 AUC 和 Logloss 指标上仅优于 LR 和 FM，但在其他先进模型（如 DNN、DCN 和标准 xDeepFM）面前仍有明显差距。这表明当前实现尚未完全发挥模型潜力。未来的研究可以集中在模型架构优化、超参数调优和数据预处理改进等方面，以进一步提升性能。尽管如此，实验结果仍然验证了 xDeepFM Ours 在特征交互建模方面的基本能力，其在推荐系统中的应用价值仍值得期待。

8. 结论

本实验的结论是，xDeepFM 模型在点击率预测任务中具有较强的性能优势，尤其是在高阶特征交互的建模上表现突出。通过结合显式的 CIN 模块和隐式的 DNN 模块，xDeepFM 能够有效捕捉到复杂的特征关系，实现更精确的用户偏好预测。尽管我们复现的结果与原论文存在轻微差距，但整体趋势

与原论文一致，验证了 xDeepFM 模型在推荐系统中的有效性和鲁棒性。这表明，将显式和隐式特征交互相结合的设计，对提高推荐系统的预测能力具有重要的实践意义。

9. 训练过程展示与分析

本节分析了模型在 5% 数据下的训练过程，通过绘制 AUC 和 Log Loss 曲线来直观展示模型的训练和验证性能变化趋势。

9.1. AUC 和 Log Loss 曲线

图 4, 图 5展示了模型在 5% 数据下的训练和验证集 AUC、Log Loss 随训练轮次的变化情况。从曲线中可以观察到，模型在训练集上的性能逐步提升，而验证集的性能则在训练后期出现了一定程度的下降，可能表明模型在训练过程中出现了过拟合现象。

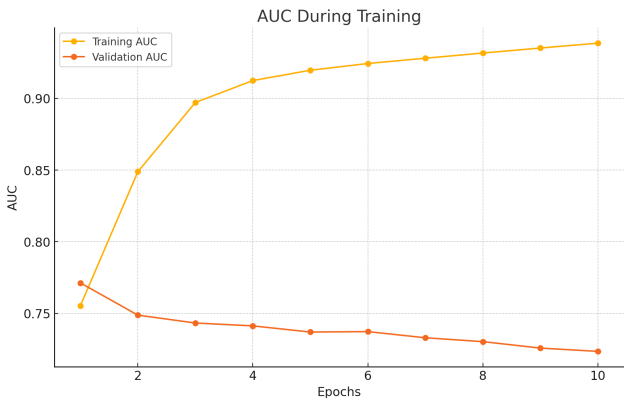


图 4. 模型训练过程的 AUC 曲线图 (5% 数据)

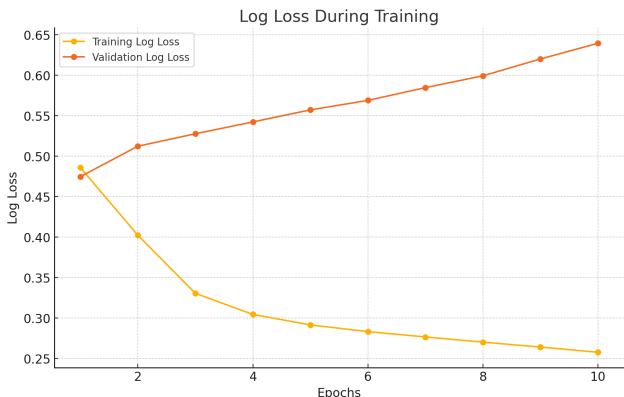


图 5. Wide & 模型训练过程的 Log Loss 曲线图 (5% 数据)

9.2. 训练过程分析

从 AUC 曲线可以观察到，训练集的 AUC 随着训练轮次的增加逐步提升，从初始的 0.7553 稳定增长至最终的 0.9386，这表明模型在训练过程中能够有效地学习特征交互信息并提高预测能力。然而，验证集的 AUC 在前几个轮次中略有下降，从初始的 0.7713 降至最终的 0.7235。这种现象表明模型在训练过程中可能出现了过拟合，即模型在训练数据上的表现逐渐提升，但泛化能力有所下降。

从 Log Loss 曲线可以看出，训练集的 Log Loss 随着轮次的增加显著下降，从初始的 0.4859 降至最终的 0.2578，表明模型逐步收敛并有效减少了分类误差。然而，验证集的 Log Loss 在训练后期逐步升高，从初始的 0.4748 增加至最终的 0.6396，这再次验证了模型可能存在过拟合问题。

综合来看，模型在训练集上的性能持续提升，但在验证集上的性能（AUC 和 Log Loss）在训练的中后期趋于恶化，表明训练时间过长可能导致模型过拟合。此外，验证集 AUC 在第二个训练轮次之后明显下降，表明模型的泛化能力需要进一步优化。

10. 所完成的加分项

加分项内容	对应章节索引
训练过程展示与分析分析	章节9

表 2. 完成的加分项及其在实验报告中的对应章节索引。

如表 2 所示，本次实验中完成的加分项为提供模型训练的可视化结果并对其进行详细分析。具体内容包括模型在训练和验证集上的 AUC 和 Log Loss 曲线，并结合曲线分析模型的训练性能、验证集的表现以及可能存在的过拟合现象。此外，我们还针对模型训练过程中观察到的问题提出了改进建议。相关内容已在章节 9 中详细描述。

11. 成员分工及贡献比

见表 3

任务类别	王佳和负责的部分	张怀松负责的部分	贡献比
文献综述与技术分析	文献综述	技术与模型特点	50%
形式化定义	定义训练目标和损失函数	定义任务的数学模型	50%
论文复现总结	算法模块分析	理论基础总结	50%
模型复现	构建模型架构	参数调优与技术实现	50%
实验设计与执行	设计实验流程	执行实验与数据处理	50%
结果分析与优化	分析结果与提出优化方向	详细优化与调整	50%
撰写报告	撰写结构与内容	完成技术细节与整合结果	50%

表 3. 王佳和, 张怀松在项目中的分工表

12. 心得体会

在此次项目中, 我们共同致力于 xDeepFM 模型的复现和研究, 这个过程既充满挑战也极具启发性。通过深入的文献综述, 我们不仅加深了对推荐系统领域当前最前沿技术的理解, 还对复杂的模型结构和算法有了更为透彻的认识。形式化定义和实验设计的过程强化了我们在理论和实际操作之间的联系, 使我们能够更有效地将理论应用于实际问题解决中。

实验过程中, 我们面对的主要挑战是模型参数的调优和结果的优化。每一步的微调都可能对最终结果产生显著影响, 这要求我们不仅要有耐心, 还需要具备解决问题的创造性思维。实验结果的分析阶段是提升我们批判性思维和解决问题能力的绝佳机会, 我们学会了如何从数据中寻找模式, 识别问题, 并探索可能的改进策略。

在撰写报告的过程中, 我们学习了如何清晰、系统地表达技术细节和研究发现。这不仅提升了我们的写作技能, 也加深了我们对项目的理解和反思。整体来说, 这个项目不仅增强了我们的技术能力和团队合作精神, 更让我们认识到了持续学习和适应快速发展技术领域的重要性。我们期待将这一经验应用于未来的项目, 以达到更高的研究与实践水平。

参考文献

- [1] 朱扬勇 and 孙婧, “推荐系统研究进展,” 计算机科学与探索, vol. 9, no. 5, pp. 513–525, 2015. 2
- [2] Pablo Castells and Dietmar Jannach, “Recommender systems: A primer,” arXiv preprint arXiv:2302.02579, 2023. 2
- [3] J Ben Schafer, Joseph Konstan, and John Riedl, “Recommender systems in e-commerce,” in Proceedings of the 1st ACM conference on Electronic commerce, 1999, pp. 158–166. 2
- [4] Rand Jawad Kadhim Almahmood and Adem Tekerek, “Issues and solutions in deep learning-enabled recommendation systems within the e-commerce field,” Applied Sciences, vol. 12, no. 21, pp. 11256, 2022. 2
- [5] Yakun Li, Jiaomin Liu, and Jiadong Ren, “Social recommendation model based on user interaction in complex social networks,” PloS one, vol. 14, no. 7, pp. e0218957, 2019. 2
- [6] Hengyu Zhang, Chang Meng, Wei Guo, Huifeng Guo, Jieming Zhu, Guangpeng Zhao, Ruiming Tang, and Xiu Li, “Time-aligned exposure-enhanced model for click-through rate prediction,” arXiv preprint arXiv:2308.09966, 2023. 2
- [7] Li Zhang, Weichen Shen, Jianhang Huang, Shijian Li, and Gang Pan, “Field-aware neural factorization machine for click-through rate prediction,” IEEE Access, vol. 7, pp. 75032–75040, 2019. 2

- [8] 王志格, 李汪根, 夏义春, 高坤, 束阳, and 葛英奎, “基于场矩阵分解机和 cnn 的点击率预测模型,” 计算机系统应用, vol. 33, no. 1, pp. 87–98, 2023. 2
 - [9] Wei Zhang, Yahui Han, Baolin Yi, and Zhaoli Zhang, “Click-through rate prediction model integrating user interest and multi-head attention mechanism,” Journal of Big Data, vol. 10, no. 1, pp. 11, 2023. 2
 - [10] Shiqi Li, Zhendong Cui, and Yongquan Pei, “A dual adaptive interaction click-through rate prediction based on attention logarithmic interaction network,” Entropy, vol. 24, no. 12, pp. 1831, 2022. 2
 - [11] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishikesh Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al., “Wide & deep learning for recommender systems,” in Proceedings of the 1st workshop on deep learning for recommender systems, 2016, pp. 7–10. 2
 - [12] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He, “Deepfm: a factorization-machine based neural network for ctr prediction,” arXiv preprint arXiv:1703.04247, 2017. 2
 - [13] Jianxun Lian, Xiaohuan Zhou, Fuzheng Zhang, Zhongxia Chen, Xing Xie, and Guangzhong Sun, “xdeepfm: Combining explicit and implicit feature interactions for recommender systems,” in Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining, 2018, pp. 1754–1763. 2
- * ctr0.05.ipynb: 使用 5% 的数据集运行的 xDeepFM 模型代码和输出。
 - * ctr0.5.ipynb: 使用 50% 的数据集运行的 xDeepFM 模型代码和输出。
 - * ctr1.ipynb: 使用 100% 的数据集运行的 xDeepFM 模型代码和输出。
 - * cin.py: 包含 xDeepFM 模型中压缩交互网络 (CIN) 部分的代码和详细注释。

附录 附录部分包含了提交的所有文件和文件夹的详细列表及其说明，以便清楚地理解每个文件的用途和内容。

文件清单及说明

- **大数据分析文件夹**: 解压大数据分析.rar 后得到的文件夹，包含以下内容：
 - **报告.pdf**: 详细描述了整个项目的研究背景、方法、实验结果和结论。
 - **海报.pptx**: 展示了项目的主要研究成果和关键点，适用于学术报告或展示。
 - **code 文件夹**: 包含所有相关的代码文件，具体包括：