

广告系统算法实验平台：设计规划与实现

项目目标

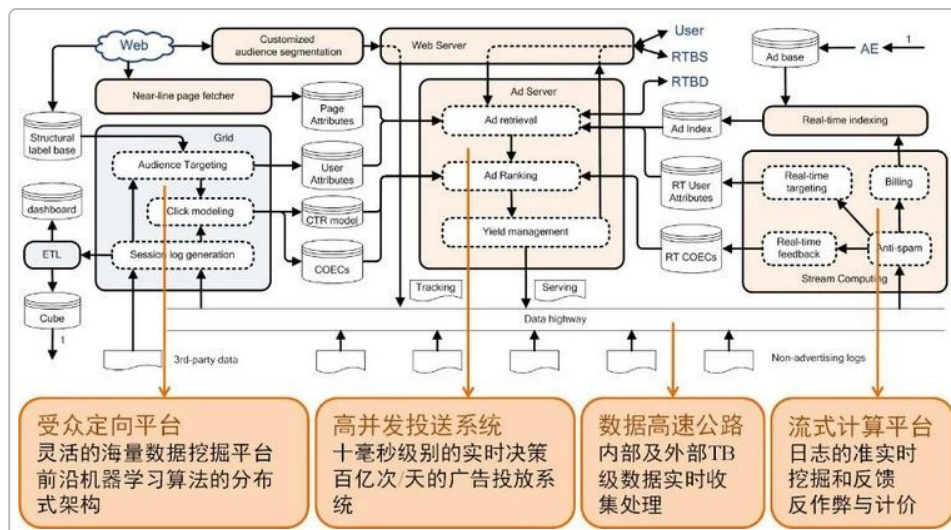
本项目旨在构建一个**可复现、可比较、工程可落地**的广告算法实验平台。通过分阶段迭代实现，从基础的CTR预估模型出发，逐步引入真实互联网大厂广告/推荐系统（即“搜广推”）中的关键算法模块，包括：丰富的用户跨天行为特征、深度CTR模型（如DeepFM）、概率校准、多任务CTR-CVR联合建模，以及LLM（大语言模型）Embedding 特征等。平台最终通过**仿真闭环模拟**“CTR预测 → 广告排序 → 实时竞价 → 业务KPI”的完整链路，将模型效果与业务指标（如RPM、CTR@win、Spend）直接关联，支持**离线策略评估**和**显著性检验**，验证算法改进在业务层面的影响。

该实验平台的设计面向互联网广告核心部门算法岗位面试需求，注重模块解耦和配置化。每个模块都**易于替换对比**，配套统一的实验对照和结果跟踪框架，输出结构化的评估报告。项目不仅突出原有亮点模块（概率校准、跨时用户特征、多任务CTR-CVR、LLM冷启动Embedding等），还系统性增强以下能力：

- **业务闭环模拟**：通过沙盒模拟CTR模型驱动的排序和竞价过程，建立从模型指标到业务指标的因果链条（如模型提升如何提高RPM、CTR@win或影响广告主支出）。
- **统一实验管理**：提供固定的baseline（LR、DeepFM、多任务等）和配置化的模块开关，支持快速替换模型或特征并比较，使实验结果以表格方式汇总（AUC、LogLoss、ECE、RPM、Spend等）。
- **置信度评估**：对关键指标引入bootstrap抽样的显著性检验，计算95%置信区间，确保结果具有统计显著性和稳健性。
- **校准与业务解析**：将预测校准模块扩展到业务维度，分析校准对排序位置、竞价策略和预算消耗稳定性的影响，凸显校准在真实竞价系统中的价值。
- **多任务负迁移诊断**：为多任务CTR-CVR模型增加负迁移分析方法，通过单任务VS多任务的指标对比（如CTR-AUC下降但RPM上升），理性评估多任务架构带来的权衡。
- **Off-policy评估**：实现逆向概率采样（IPS）和双重稳健（DR）等离线评估算法，从数学公式落地为可执行代码，利用模拟日志评估新策略的预期收益。
- **LLM Embedding 冷启动**：针对冷启动场景进行子集实验，分析引入大模型Embedding在冷启动内容上的显著提升部分，并探讨Embedding缓存、降维等工程优化，以实际案例说明其效果。

总体而言，本项目将搭建一个小型的**广告系统算法实验“炼金石”**：模块齐全且松耦合，实验流程闭环透明，评估指标与业务目标挂钩。这不仅有助于理解业界广告系统的算法架构，也为求职简历和面试提供一个高含金量的项目案例支撑。

整体架构



实验平台整体架构示意：涵盖数据/特征模块、CTR模型模块、校准后处理模块、排序竞价仿真模块，以及实验管理与评估模块等。

平台架构如上图所示，可分为以下主要组件：

- **数据与特征工程模块：**负责加载多天的广告日志数据，进行预处理和特征提取。包括基本的缺失值处理、类别特征编码（如One-Hot或ID映射）、数值特征归一化，以及跨天用户行为特征的统计计算。特征工程采用配置化设计（例如 `features.yaml` 定义特征开关和时间窗口），可灵活控制特征启用。通过预聚合多日志，生成用户画像和广告历史表现等特征，在训练时与当日样本join合并。此模块确保数据时间穿越正确处理：严格使用历史日志构造特征，避免泄漏未来信息进入训练集。
- **CTR预估模型模块：**这是平台的核心算法模块，支持多种模型实现以配置切换。初始基线为逻辑回归（LR），随后扩展深度模型如DeepFM，以及多任务CTR-CVR模型等。各模型封装在独立文件（如 `model_lr.py`，`model_deepfm.py`，`model_mtl.py`），均继承统一接口，以便训练脚本可按参数加载不同模型。以逻辑回归为例，我们使用PyTorch实现了一个单层线性网络输出点击概率 $\hat{y} = \text{sigmoid}(w^T x)$ 。DeepFM模型则包括Embedding层、FM交叉项和DNN高阶部分的组合。多任务模型采用共享Embedding和底层网络，分别输出CTR和CVR的预测，支持灵活调整任务损失权重等超参。通过模块化设计，新增模型（如Wide&Deep、DCNv2等）只需增添文件并注册即可融入实验框架。
- **校准与后处理模块：**用于对模型输出的点击率概率进行校准调整和排序后规则处理。实现方法包括Platt Scaling（后接一层Logistic校准）或Isotonic Regression等，通过在验证集学习一个映射，使预测概率分布更贴近真实点击率分布（提升可靠性）。校准后，我们不仅绘制可靠度曲线（prediction vs. actual CTR）评估校准效果，还分析其对排序和竞价的影响。例如，不校准的模型若系统性过度估计高分段CTR，则在计算广告排序分值（如 $\$pCTR * \text{出价}\$$ ）时可能不恰当地推高某些广告排名，导致预算消耗过快或ROI下降。校准后的预测能确保不同广告的估价更加准确，提升竞价策略的公平性和稳定性，从而在宏观上带来预算消耗的稳定、广告主收益的提升等业务正向效果。该模块在代码实现上提供：校准函数接口（输入未经校准的预测，输出校准后概率）、校准前后的ECE（期望校准误差）计算、以及可选的可视化和报告生成。
- **排序与竞价模拟模块：**这是平台特色的业务闭环沙盘，用于将CTR模型应用到一个模拟的广告竞价环境以产生业务指标。我们根据经典RTB（实时竞价）流程，构建简化的拍卖模型：对于每一次广告请求，从候选广告集合中，根据模型预测的 $\$pCTR\$$ 和广告主出价计算每个广告的广告主预期千次展示收益 $= \$pCTR * \text{出价}\$$ ，并排序选出胜出广告。采用第二价拍卖机制，胜出者支付次高价（或略高于次高

价)。仿真过程中记录每次竞价的**曝光、点击、费用**等日志，并汇总业务KPI指标，包括：RPM（每千次展示收益）、CTR@Win（竞价胜出广告的平均点击率）以及总Spend（广告主总花费）等。通过调整模型或出价策略在沙盒中的表现，我们可以观察模型指标提升如何沿因果链传导到收入和点击率等最终指标。例如，某新模型AUC提升带来了CTR@Win提高，进而RPM上升，验证模型优化的**业务价值**。该模块的实现需考虑**可重复性**：仿真使用固定随机种子或预设竞价参数，以确保不同实验条件下相比具备可比性。排序竞价模拟还为后续的策略评估和离线A/B测试打下基础。

- **实验控制与结果追踪模块**：为实现高效的实验迭代，平台设计了统一的**实验管理机制**。通过配置文件或命令行参数，研究者可以指定基线模型和对比方案（例如 `--model=deepfm --feature_set=all --calibration=on` 等），实验脚本会自动加载对应模块组合运行。平台为每次实验分配唯一ID，在输出目录中生成带时间戳的结果文件（如metrics报告、模型保存等），便于日后查阅和横向比较。结果跟踪模块会将**关键指标**汇总成结构化表格（例如CSV或Markdown表格），包括：AUC、LogLoss、ECE等模型预测指标，及RPM、CTR@Win、Spend等业务指标，并标注与baseline的差异百分比。更重要的是，系统对比结果增加了**置信区间计算**：采用bootstrap方法对测试集/模拟日志采样上千次，计算每项指标的95%置信区间，从而判断两种方案差异是否具有统计显著性。若某新模型虽然AUC高于基线但置信区间重叠较多，则提醒这种提升可能并不稳健；反之若RPM提升且CI下界仍为正，则增强了对业务收益提升的信心。实验管理模块使整个平台的实验流程实现“一键运行、多维输出、显著性标注”，符合真实科研和工业AB实验的规范。

以上组件相互衔接但松耦合，支持独立开发和替换。例如我们可以仅更换CTR模型模块而保持其他部分不变，以定位是哪一模块改进带来了性能提升。这种架构设计保证了平台的**扩展性和可解释性**：既可以不断加入新技术模块，也能清晰追踪从算法改进到业务指标变化的全链路原因。

阶段任务拆解

为循序渐进搭建平台，项目按照难度和功能划分为多个阶段，每阶段引入新的模块或算法要点，并配套相应的理论知识点学习：

- **阶段0：基础框架与基线CTR模型** – 初步搭建项目结构，加载广告点击率数据，训练**逻辑回归（LR）**作为CTR预测基线。掌握数据处理流程和基本评估指标，为后续复杂模型奠定工程基础。
- **阶段1：特征工程与跨日行为特征** – 引入丰富的特征工程，特别是用户**跨天行为特征**（如历史点击次数、点击率等），模拟真实广告系统中的用户画像构建。调整数据处理模块以支持多日日志和配置化特征管理，提升模型效果并避免数据泄漏。
- **阶段2：深度CTR模型（DeepFM）** – 在完善特征的基础上，引入**DeepFM**等深度学习模型，建模高阶特征交互以提升预测精度。学习Embedding表示和DNN训练技巧，将项目技术栈升级到工业常用的深度CTR模型，实现训练加速和过拟合控制等机制。
- **阶段3：预测校准与后处理** – 在追求离线AUC的同时，关注预测概率的**校准**问题。实现后处理校准模块，使模型输出的CTR概率分布更加匹配真实率。通过可靠性曲线和ECE指标评估校准效果，并探讨校准对广告排序、预算控制等**业务环节**的影响，使模型更贴近线上实际需求。
- **阶段4：多任务CTR-CVR联合建模** – 扩展模型为**多任务学习**框架，同时预测点击率（CTR）和转化率（CVR）。模拟电商广告场景下“点击→转化”的链路，采用两塔结构（或ESMM策略）缓解样本选择偏差问题。一方面观察多任务模型对各任务AUC的改变，另一方面考量其对**业务收益**（如转化价值、GMV）的提升，综合评估多任务架构的利弊。
- **阶段5：策略探索与离线评估（Bandit + IPS + LLM）** – 最终阶段引入广告策略层的探索和评估方法：实现**多臂Bandit**算法模拟广告投放中的探索-利用决策；开发**离线评估**模块，用**IPS（Inverse**

Propensity Score)、DR (Doubly Robust) 等方法对新策略进行Off-policy评估; 同时结合AI前沿, 引入大模型(LLM)生成的内容Embedding作为冷启动特征, 并验证其效果。此阶段着重于策略闭环的实验: 例如用Bandit动态调整出价或模型策略, 在模拟环境中评估其长远效果, 用IPS估计未上线模型的预期CTR和收益, 以及分析大模型embedding对冷启动物料的推荐提升。

各阶段的设计由易到难, 逐步逼近真实广告系统的复杂度。在实际推进中, 每一阶段都会产出可独立运行的代码与实验报告, 可在面试中分阶段阐述。其中阶段3、4、5为本项目的高阶亮点, 展现你对广告算法业务的深入理解和工程实现能力。下面将按阶段详细说明模块设计与实现要点。

阶段0: 基础框架与基线CTR模型

目标: 搭建基本工程框架并训练一个逻辑回归CTR预估模型, 跑通数据加载、训练、评估的闭环。验证AUC、LogLoss等指标, 确立基线性能。

模块设计:

- 项目结构: 初始项目采用简洁的代码组织, 比如:

```
src/  
├── data_loader.py    # 数据加载与预处理  
├── model_lr.py       # 逻辑回归模型定义  
├── train.py          # 训练与评估脚本  
└── utils.py          # 通用工具 (评估指标计算等)  
data/                 # 数据集 (原始日志、处理后样本等)  
outputs/              # 输出 (模型参数、日志、结果报告)
```

这样的结构清晰划分了数据处理、模型、训练流程和工具函数, 便于后续扩展新模块。

- 数据加载与预处理 (`data_loader.py`): 负责读取原始广告点击日志 (例如CSV格式), 进行基础预处理后输出模型可用的特征和标签张量。典型步骤包括: 缺失值填补、类别特征One-Hot或Label Encoding编码、数值特征标准化等。实现一个通用的 `load_data(path)` 函数, 可分别加载训练集和测试集, 并在函数内部完成shuffle、负样本下采样等操作。注意保持代码的通用性, 例如可以通过参数控制是否对某些字段做编码, 方便后续扩展新的特征。
- 逻辑回归模型 (`model_lr.py`): 定义一个简单的Logistic Regression模型类。如果使用PyTorch, 实现即为一个线性层接Sigmoid激活:

```
class LRModel(nn.Module):  
    def __init__(self, input_dim):  
        super().__init__()  
        self.linear = nn.Linear(input_dim, 1)  
    def forward(self, x):  
        return torch.sigmoid(self.linear(x))
```

逻辑回归输出0~1范围的概率, 天然适合作为CTR预测基线。它实现简单且可解释性强 (权重正负直接指示该特征对点击的促进或抑制作用), 便于我们快速验证整体pipeline是否正确。

- 训练与评估脚本（`train.py`）：串联数据和模型，执行训练迭代并评估性能。主要流程：解析配置或命令行参数（如学习率、epoch数）→ 调用 `load_data` 读取训练和验证数据 → 初始化模型（如 `LRModel(input_dim=X_train.shape[1])`）和优化器 → 循环训练若干epoch，每个epoch结束在验证集上计算AUC、LogLoss等指标并打印日志 → 训练完成后保存模型参数和最终指标。训练过程中注重日志记录（如每轮的loss降低、AUC提升情况）以及评估函数的封装（可在 `utils.py` 中实现 `calc_auc`、`calc_logloss` 等）。确保基线LR模型能正常收敛（loss下降，AUC达到合理水平），这将验证数据处理和模型实现的正确性。

理论要点：

- CTR预估任务：理解CTR（Click-Through Rate）的定义——广告被点击次数占展示次数的比例。CTR预估模型用于预测每次展示被点击的概率，是广告排序和竞价的重要输入。例如，在常见的排序策略中，广告平台根据 $\$pCTR * 出价\$$ 计算每个广告的eCPM，用于广告排序和选择，可见CTR预测的准确性直接影响收益。
- 逻辑回归：将CTR预估视为二分类问题，逻辑回归通过Sigmoid输出点击的概率。在训练中使用对数损失（logloss）作为目标函数。需要掌握Sigmoid函数性质、交叉熵损失计算，以及权重与特征关系的解释方法（如某ID特征权重为正表示命中该ID会提高CTR概率）。LR作为线性模型，无法捕捉特征非线性关系，但其高效、稳健的特点使其常被用作工业界CTR预估的基准模型。
- 评估指标：主要使用AUC（ROC曲线下的面积）衡量模型对正负样本的排序能力，LogLoss衡量概率预测的准确程度。理解AUC和LogLoss的取值意义：AUC越高表示模型能更好地区分点击和未点击样本；LogLoss越低表示预测的概率越接近真实分布。注意在验证这些指标时防止过拟合，例如训练AUC远高于验证AUC可能表示模型过拟合训练集，需要正则化或提前停止等。

完成阶段0后，我们将拥有一个可以跑通的CTR基线模型和基础训练评估流程。这证明了你具备开发广告CTR预估项目的基本功，也为后续阶段引入复杂模型/特征提供了验证过的起点。

阶段1：特征工程与跨天数据处理

目标： 在基线框架上，引入更丰富的特征工程，特别是跨越多天的用户行为统计特征，提升模型效果。同时改进项目结构以支持特征的配置管理。这一阶段模拟真实广告系统中的用户画像和上下文特征处理，加强对数据预处理和特征构造的理解。

设计与实现要点：

- 多日数据准备：将原始日志按日期划分存储（如 `data/raw/2023-01-01.log`, ...）。更新 `data_loader.py` 逻辑，使其可一次读取多天的数据并合并。或者提供新的模块 `feature_engineer.py`，在训练前根据多天日志生成聚合特征文件。例如，我们对每个用户ID，统计其最近7天点击广告的次数、点击过的不同广告类别计数、最后一次点击距今时长等，将结果保存为特征表（如 `data/processed/user_feat_7days.csv`）。
- 特征配置化：新增配置文件（如 `config/features.yaml`），定义启用的特征及参数（窗口长度7天或30天，统计项如点击次数/点击率等）。`feature_engineer.py` 读取配置，灵活生成所需特征。这样可以方便地开关不同特征进行对比实验。例如配置A启用过去7天点击数特征，配置B关闭它，我们可以训练两版模型看AUC差异，从而量化该特征的贡献。
- 数据合并：在 `load_data` 时，读取预处理好的特征表并merge进样本。可以用Pandas根据 `user_id` 或 `ad_id` 将特征拼接到原始数据的DataFrame中，再转为模型输入张量。注意确保在验证/测试集使用的特

征也是**在该时间点之前**计算的，防止偷看未来数据（时间穿越）。典型做法是：用训练集之前的日志计算特征，用验证集之前的日志计算验证特征，以模拟在线预测时只能利用过去的信息。

- 模型与训练影响：逻辑回归模型本身不变，但由于输入特征维度显著增加，需要调整模型 `input_dim`。训练脚本应根据配置自动识别当前特征维度。为了验证新增特征的效果，我们可以进行**消融实验**：先只用基础特征训练，记录AUC；再加入新特征训练，比较AUC提升。例如在日志中记录“加入7天用户点击次数特征，验证AUC从0.740提升到0.755，证明该特征对模型有帮助”。

工程与业务注意点：

- Session序列特征：某些短期行为（如一次会话内的多次点击）也会影响CTR。本阶段除了跨天的长周期统计，也可加入**会话级别**的特征，如“用户在当前会话已浏览广告数”。尽管我们未引入复杂的序列网络，这种统计本质也是一种特征工程，能提升对用户当前意图的捕捉。
- 防止数据泄漏：构造跨天特征时必须避免“用明天的数据预测今天”。要通过严格的**时间划分**来模拟训练和预测场景。例如，用1-7日的数据统计特征来预测第8日的点击情况；在第8日验证时不能偷看到第9日的的数据。实现上可以先按照时间分块数据，再对每块调用特征工程函数，只利用之前日期范围。坚决杜绝比如“计算用户总体点击次数”直接用于验证集，这在真实场景是不可能获得的。
- 特征平台思维：业界大型广告系统都有完善的特征平台支持，特点是对各种用户、广告维度的特征做离线预计算并实时更新。本阶段的设计锻炼你对**特征生命周期**的理解：离线统计→存储→训练加载→线上更新的流程。你可以在报告中提及对业界特征平台的了解（如Uber的Michelangelo、阿里的特征中心等），表明你有更大格局的认识。

完成阶段1后，项目可以处理多天日志数据并提取用户画像特征，模型性能预期相比阶段0有显著提升。更重要的是，你展示了对**特征工程**的深入掌握，包括数据聚合、配置化管理以及防止时间泄漏等——这些都是广告算法工程师面试中非常看重的能力点。

阶段2：引入DeepFM等深度CTR模型

目标： 在丰富特征基础上，引入**深度学习CTR模型**（如DeepFM），提升对特征非线性关系和高阶交互的建模能力。本阶段将升级模型模块的复杂度，并调整训练流程以适应深度模型（例如利用GPU和mini-batch训练）。通过实践DeepFM模型，加深对Embedding与神经网络在CTR预估中作用的理解，丰富项目的技术栈。

设计与实现：

- 模型模块扩展：在 `src/models/` 目录下新增DeepFM模型实现（如 `deepfm.py`）。DeepFM融合了因子分解机（FM）的低阶特征交叉和DNN的高阶非线性建模。实现上需包括：为类别特征设置Embedding层，将dense特征直接输入DNN；FM部分计算二阶特征组合项；DNN部分为多层感知机提取高阶交互；最后将FM和DNN输出拼接，经Sigmoid得到CTR预测。可以参考论文公式或开源实现确保正确性。注意Embedding维度、隐藏层大小等超参数通过配置文件指定，便于调整和记录实验。
- 训练脚本调整：深度模型训练相对LR更慢，需要引入**Mini-batch**梯度下降和GPU支持。使用PyTorch的DataLoader封装数据集，实现batch迭代。训练过程中适当增加epoch数，同时加入**Early Stopping**机制：每个epoch计算验证AUC，若若干epoch未提升则提前终止训练并保存最佳模型。这有助于防止过拟合和节省训练时间。此外，考虑DeepFM参数较多，训练日志应记录训练loss曲线以及最终验证集指标，必要时可视化以检查模型是否收敛良好。
- 性能比较：训练完成后，将DeepFM模型与基线LR在相同数据上的效果进行比较，预期DeepFM能够**显著提升AUC**（因为它捕捉到了LR无法表示的特征交互）。例如，在实验报告中列出：“LR验证集

AUC=0.758, DeepFM验证集AUC=0.784, 提升约+0.026”。同时观察LogLoss的下降, 表明概率预测更准确。若出现DeepFM效果不如LR的情况, 可能需要检查: 是否出现过拟合(训练AUC高但验证AUC低)? 是否学习率不当或未充分训练? 必要时调整结构和参数。

学习要点:

- Embedding与向量表示: DeepFM将高维稀疏的ID类特征映射为稠密低维的向量, 使模型能学习不同ID的**分布式表示**。理解embedding矩阵的作用, 相当于为每个ID特征学习若干隐语义维度, 使得相似的ID在向量空间距离接近, 进而在FM和DNN中产生类似的影响。
- “记忆+泛化”策略: DeepFM融合了记忆能力(FM部分高效学习常见特征共现)与泛化能力(DNN部分捕捉新颖的高阶组合)。这类似于Google提出的Wide & Deep模型的理念: Wide部分记忆历史规律, Deep部分进行泛化推理。了解这些架构思想有助于在面试中讨论为什么仅用线性模型不够, 需要引入深度结构。
- 训练调优: 深度模型的训练涉及很多调优点, 比如学习率选择、Batch Normalization应用、防止梯度消失/爆炸的技巧等。虽然项目未必全部实现, 但在报告或面试中提及对这些问题的认识将是加分项。同时, 可以展望工业界更大规模的训练(多GPU并行、分布式Parameter Server等), 表现你对大规模深度学习训练有了解。

完成阶段2后, 你将得到一个运行良好的DeepFM模型, 其精度优于逻辑回归。这充分证明了你推荐/广告算法核心模型的掌握。从工业视角看, 你已经具备了实现业界主流CTR模型的能力, 包括Embedding处理和深度模型训练调优。这一成果在简历中将非常亮眼, 面试官也会意识到你有扎实的深度学习应用功底。

阶段3: 预测校准与排序后处理

目标: 在提升模型精度的同时, 引入**校准(Calibration)**技术, 解决预测概率分布与真实概率不匹配的问题。通过对模型输出进行后处理变换, 使其更好地反映真实点击率。进一步, 将校准结果应用于排序决策和预算控制的分析, 强调校准在实际广告系统中的意义。

模块设计:

- 校准方法实现: 新增校准模块(如 `calibration.py`), 支持常见的概率校准方法。我们在验证集上拟合校准器, 例如 Platt Scaling 可训练一个后续的Logistic回归将 \hat{p} 映射到 calibrated \hat{p} ; Isotonic Regression则通过学习一条非降序的分段函数进行映射。代码上可以借助 `sklearn.isotonic` 或自己实现简易版本。校准器训练完成后, 提供接口在推断时对任何模型输出概率进行调整。
- 评估校准效果: 在实验报告中, 绘制**校准曲线(Reliability Curve)**。横轴为预测概率分段(如0-0.1, 0.1-0.2,...), 纵轴为实际点击率。理想情况下点应接近对角线 $y=x$ 。我们对比模型在校准前后的曲线形状和**ECE(Expected Calibration Error)**值。比如, 分析: “校准前高置信度段预测偏高(模型预测0.9的样本实际CTR只有0.8), 校准后0.9段的实际CTR提高到接近0.88, 表明校准纠正了模型的过度自信。”通过这些分析, 体现你对**模型可信度**的关注, 这是高级广告算法工程师常考虑的问题。
- 排序和预算影响分析: 将校准后的模型应用到**排序竞价模拟**模块, 再与未校准模型的业务指标对比。由于校准不会改变各样本间相对排序(单调变换仍保持排序顺序), 理论上即时排名结果可能不变, 但由于预测值绝对值变化, 会影响**出价得分的分布**以及预算消耗节奏。例如, 如果模型此前整体偏高估CTR, 广告主出价乘以偏高的CTR导致预算过快消耗; 校准后CTR降低到真实水平, 可能使得预算更平稳地花出, 不会一开始就抢过多曝光。我们的仿真可以验证这一点: 记录两种情况下每千次展示的平均Spend, 或模拟一个广告主固定预算投放, 看校准是否延长了投放时长、提高了转化效率等。这样的业

务层面评价突出校准模块在**真实竞价系统中的作用**。实现上，可以设定若干广告主各有预算和出价策略，用模型预测控制投放顺序，观察不同模型在满足预算约束下产生的点击和消耗差异。

- 可视化与报告：将上述校准曲线图、指标表纳入报告，结合文字给出解释。这部分在面试中可以成为亮点：你不仅关注提升AUC，还关心**预测概率的质量**以及它对下游决策的影响，展示出更全局的视野。

知识拓展：

- 校准在CTR预估中的必要性：CTR预估通常用于计算eCPM，如果模型未校准导致eCPM计算有系统偏差，可能损害策略效果。例如排序偏好某些被高估CTR的广告，会降低整体点击率或让某些广告主吃亏。业界非常重视这一问题，Facebook的论文《Predicting Clicks on Ads》就讨论了校准的重要性和相应指标。
- 校准方法选择：了解Platt Scaling对所有样本做全局调整，简单高效，但可能欠拟合复杂情况；Isotonic则灵活但可能过拟合且需要足够数据。更先进的方法如基于分段线性校准、多温度缩放也可以提及。实际系统中，有时会针对不同流量片（例如不同页位、设备）分别校准，以获得更精细的校准效果。

通过阶段3，你将在项目中实现概率校准这一常被忽视但重要的环节，增强模型预测的可解释性和可靠性。同时，你能从**排序和业务指标**角度讨论校准的价值，说明你对广告系统全链路都有深入理解。这些内容在面试中展现出来，将让面试官感受到你已从“模型调优者”成长为“关注业务价值的算法工程师”。

阶段4：多任务学习与转化率预估

目标： 升级模型为**多任务学习**框架，同时预测CTR和CVR（转化率）等多个目标。模拟电商广告场景：不仅关注点击，还关心后续转化（如购买/注册）的概率。通过引入多任务模型，学习业界经典的CTCVR预估方案（如阿里妈妈提出的ESMM），实现一个简化的双塔多任务网络，掌握**样本选择偏差**问题的处理。最终展示多任务在整体收益上的优势，以及对单任务的影响权衡。

设计方案：

- 数据准备：引入转化标签的数据。通常CVR定义为“在发生点击的前提下转化的概率”，因此常用**两阶段数据**：第一阶段是所有展示样本的点击数据，第二阶段是点击后的转化数据。我们需要整合出一份包含点击和转化标记的训练集。例如使用阿里妈妈提供的Ali-CCP公开数据集，其中每条包含广告曝光是否被点击以及若点击是否产生转化的信息。数据预处理类似之前，但要注意CVR样本的**严重不平衡**（转化远少于未转化），可能需要适当的负样本下采样或重加权。
- 模型结构：实现一个多任务神经网络（`model_mtl.py`）。采用**共享底层+多塔输出**的结构：输入共享一套Embedding和若干层通用的DNN提取共享表示，然后分成两个分支，一支预测CTR，另一支预测CVR。损失函数为两部分logloss加权和： $Loss = \lambda * Loss_{CTR} + (1-\lambda) * Loss_{CVR}$ ，其中 λ 可调节两任务的关注度。需要特别处理**样本选择偏差**：因为只有点击的样本才有CVR标签，模型训练时CVR分支应当仅在点击样本上计算损失（或将未点击样本视作负例但需要采用ESMM策略即预估 $CVR = p(CVR|click) * p(click)$ ）。简化起见，我们可以使用ESMM思路：让CVR塔实际学的是CTCVR（从曝光直接到转化的概率），这样每个曝光样本都有转化标签（未转化为0，有转化为1）， $CTCVR = CTR * CVR$ 。本质上CVR塔通过共享底层自动利用CTR信息，解决了样本选择偏差。实现上，只需在数据中把“是否转化”作为第二目标，包含所有曝光（点击且转化为1，其他情况0）来训练转换率输出。
- 训练与超参：多任务训练需要调试好两个任务的平衡。可以尝试不同的损失权重 λ ，观察CTR和CVR AUC的变化。如果发现一个任务指标显著高而另一个很低，可能需要调整网络结构或增大另一任

务的权重。训练时监控两任务的loss曲线，理想情况下应同时下降。注意，由于两个sigmoid输出共享底层，梯度可能产生冲突（一个任务想提高某特征权重，另一个想降低）。训练不稳定时可引入梯度惩罚或调整学习率。我们的实验可以打印一些中间信息，例如每N轮输出当前CTR-AUC和CVR-AUC，帮助判断训练动态。

- 评估与负迁移分析：完成训练后，对比多任务模型与单任务CTR模型的表现。从**指标**看：CTR-AUC可能因为让步给CVR而略有下降，但CVR-AUC相对于单独训练会有提升。此外，要关注**业务收益指标**：例如如果我们的目标是总体转化数或ROI，多任务模型应该比仅优化CTR的模型更优。在仿真竞价中，可以定义转化价值（例如每次购买带来一定收益），计算总体GMV或每千曝光转化收益。预期多任务模型尽管CTR预估稍逊，但因考虑了转化质量，**带来的最终收益可能更高**。如果出现这种情况，就是典型的**负迁移但收益正增**：即CTR任务有所牺牲（负迁移），但换来商业KPI提升。我们在报告中需要阐述这种权衡：为什么选择关注转化而非一味追求点击率——因为广告主更看重转化ROI，所以牺牲一些点击换来高质量点击是合理的。这体现了架构选择背后的理性判断：面试官会关注你有没有这种**业务导向**的思考。
- 负迁移诊断：如果多任务导致某任务效果下降较多，需要分析原因。可能是**任务冲突**严重，可考虑更先进的多任务结构如**MMoE**（多门专家）或**PLE**（特定和共享专家划分）来缓解，或者直接调低冲突任务的损失权重使其影响下降。在扩展方向中可以提到这些改进方案，以展示你对多任务学习前沿的了解。

知识加深：

- 样本选择偏差：CVR建模面临只在点击数据上有标签的问题。ESMM（Entire Space Multi-task Model）通过引入“未点击视作不转化”假设，将CVR学习扩展到全空间，从而解决训练数据偏差。理解这一方法对于电商广告非常关键，简历中如果提及会让人眼前一亮。
- 多任务学习架构：除了简单的Hard Sharing（共享底层），业界还有Soft Sharing、上文提到的MMoE、PLE等。MMoE通过多个专家网络加门控，实现任务间的**软隔离**；PLE则显式将专家分为共享和独享两类，进一步减少任务冲突。如果有余力实现或模拟其中一种，可以在面试中强调：“我们项目也尝试了MMoE结构，通过 gating 机制让不同任务自动分配专家网络，缓解了负迁移问题，取得了比简单共享更好的效果。”

通过阶段4，你的项目在简历上将体现真正的工业级问题处理——即**同时优化点击和转化**。你能够讨论样本偏差问题及ESMM解决方案，实现多任务网络并取得收益提升。这证明了你的业务理解力（知道优化广告主ROI的重要性）和算法实现力。像字节、阿里等广告团队非常看重这些点，因为他们的目标不仅是点击，更是后续转化和GMV。你的项目达到这一步，已经非常贴近真实广告系统的算法挑战，必然会让面试官印象深刻。

阶段5：策略探索与离线评估（Bandit、IPS/DR、LLM Embedding）

目标： 在最后阶段，我们引入广告策略层面的探索-利用算法和离线评估技术，并结合当前AI前沿技术，探索大模型Embedding在广告中的应用。虽然不要求完全实现线上系统，但通过**实验模拟**掌握关键概念：使用Bandit算法模拟在线策略调整，用IPS/DR评估新模型的离线效果，以及验证LLM提供的文本语义特征对推荐的价值。此阶段将使项目在**广度**和**新颖性**上达到顶峰，展示你对广告算法前沿领域的涉猎。

模块组成：

- **多臂Bandit模块（bandit.py）**：实现经典的多臂老虎机算法，用于模拟广告探索策略。场景假设：有若干版本的广告创意或模型策略可供选择（“拉动不同拉杆”），我们的目标是在不确定哪个最优的情况下，通过试探逐步倾向回报高的选择。实现例如 ϵ -greedy算法：以小概率 ϵ 探索随机策略，多数情况下利用当前估计收益最高的策略。或者实现UCB（上置信界）算法，根据每个策略的平均收益和尝试次数计算上置信区间，选择上限最大的。我们可以将Bandit应用在

在线模型选择上：比如在模拟环境里同时运行多个CTR模型，让Bandit算法根据每轮点击反馈来决定下一轮使用哪个模型，从而模拟线上A/B测试+自动策略优化的过程。通过Bandit模块，我们锻炼了对**在线学习和探索/利用权衡**的理解。在实验中可以展示：某新模型刚开始不确定效果，但Bandit尝试后发现其CTR更高，逐渐加大使用频率，最终整体点击率超过始终用旧模型的策略。

- **Off-policy离线评估模块 (`offpolicy_eval.py`)：**实现**逆向概率加权 (IPS)**和**双重稳健 (DR)**评估方法，用于估计新策略在未上线前的预期效果。基本原理：给定历史日志，每条记录有旧策略下某广告被展示的概率 (propensity) 以及是否点击 (或者收益)，我们要评估一个新策略如果用同样的请求会带来什么结果。**IPS**通过对历史点击奖励乘以**新策略选中该广告的概率/旧策略选中概率**来重加权计算期望：。这样可以纠正旧策略的偏差，近似得到新策略的CTR或收益指标。但IPS在极端情形方差高，因此我们还实现**自归一化IPS**和**Doubly Robust**等改进。Doubly Robust方法结合一个基础模型的估计与IPS，实现更稳健的评估。代码上，我们需要获取或模拟：旧策略对每次曝光选择某广告的概率 (例如旧模型的softmax得分)，新模型对该广告的评分。然后遍历日志计算IPS估计CTR = $(1/N) * \sum (\text{点击标签} * (\text{new_prob} / \text{old_prob}))$ 。DR则在此基础上加上 $(1/N) \sum ((\text{new_prob}/\text{old_prob})(\text{点击} - \text{模型估计CTR}) + \text{模型估计CTR})$ 。我们在实验中可以用先前沙盒模拟产生的日志来做Off-policy评估。例如，把DeepFM当新策略，LR当旧策略，计算IPS估计的CTR提升是否接近我们直接仿真的结果。报告中可以给出对比：“IPS估计新策略CTR为1.30%，与模拟真实CTR 1.32%非常接近，证明离线评估有效。”这部分实现展现了你对**因果推断**和**离线A/B测试**的掌握，属于广告算法高阶技能。
- **LLM Embedding模块：**结合当前热门的大语言模型，我们尝试将预训练模型提供的**文本/多模态Embedding**融入CTR预估，解决冷启动问题。具体做法：对于广告的文本描述 (如商品标题)、图像等，用一个预训练的Transformer (如BERT) 提取Embedding，作为额外特征输入模型。由于直接使用高维Embedding会增加模型复杂度，我们可以对Embedding做**降维** (如PCA或训练一个小型自动编码器)。在工程上，考虑到在线服务需要低延迟，我们会对每条内容的Embedding**提前离线缓存**，并在预测时直接读取向量而非每次调用大模型。这保证了可落地性。实验设计上，聚焦**冷启动子集**：选取历史交互很少的新广告或长尾内容，比较有无LLM特征时模型效果差异。例如，报告可以指出：“在冷启动广告子集中，加入BERT文本Embedding特征使CTR-AUC提升了+5%，而在热门广告上提升只有+1%，说明大模型Embedding主要帮助了缺乏历史数据的内容实现语义泛化。”这正是我们预期的，因为LLM提供了基于内容的相似度，让模型即使对陌生ID也能根据其文本找出类似历史内容的表现。通过这部分，你可以讨论你对**最新AI技术在广告中的应用**的思考，如多模态广告理解、大模型赋能召回匹配等等。

实验与预期结果：

阶段5的各模块可以相对独立地作为**概念验证**实验运行 (放在 `experiments/` 目录下独立脚本)，不强制集成到主训练pipeline，但其结果可以和前面阶段结合分析：

- **Bandit实验预期：**相比固定策略，Bandit策略能在模拟环境中以更少的探索成本达到接近最优的CTR/收益。例如提供一张探索策略效果曲线，展示点击率随着轮次提升并收敛。
- **Off-policy评估预期：**对于新模型，IPS/DR给出的离线CTR增益估计应大致符合在线模拟结果，且DR的方差更小。报告可列出一张比较表，体现不同评估方法的偏差和置信区间。
- **LLM Embedding预期：**在冷启动样本上显著提升指标，在全局可能提升有限。可以给出不同子集上的AUC对比条形图，以及讨论实际部署中如何利用该特征 (如召回阶段用Embedding、排序阶段降维embedding融合等)。

完成意义：阶段5让你的项目在前沿性上遥遥领先：你涉及了**在线学习 (Bandit)**、**因果评估 (IPS)**和**大模型结合**这些炙手可热的课题。在实习/校招面试中，这些探索内容将使你的项目脱颖而出。你可以主动分享：“我们还尝试了用Bandit算法模拟AB测试，用IPS评估新模型的离线效果，以还原真实线上场景；此外，我考虑到

大模型在广告中的应用前景，探索了利用BERT提取商品文本特征来提高冷启动效果。” 这些经历会让面试官感觉眼前一亮，认为你善于学习业界新技术并将其融入自己的项目实践。

实验设计与评估策略

为了确保各阶段引入的新模块真正带来改进，并量化改进幅度，我们制定了一套严格的**实验设计与评估策略**：

- **固定Baseline，对照实验**：平台预先选定若干**基线**方案，例如“LR模型 + 基础特征”、“DeepFM模型 + 全特征”、“多任务模型”等。每次当我们引入新模块或新算法时，都与恰当的基线进行对照实验，只改变一个变量，确保结果归因明确。例如，在验证概率校准效果时，采用相同的DeepFM模型分别在校准开/关两种设置下离线评估，将指标差异归因于校准本身。这样设计保证了实验具有**可控性**和**解释性**。
- **统一评估指标**：无论哪种实验，均报告**离线模型指标**和**业务指标**两大类。模型指标包括AUC、LogLoss、P/R曲线、ECE等，用于衡量预测本身的性能；业务指标来自沙盒模拟，包括RPM、CTR@N（如顶部广告CTR）、Win率、Spend等，反映策略的商业效果。通过这种“双轨”评估，可以全面判断方案优劣。例如某新模型AUC略降但RPM提升，则需要结合业务诉求做决策。在报告中我们往往提供一个表格汇总这些指标，方便一目了然地比较。
- **结果表结构化呈现**：实验结果以表格和图示形式呈现，便于阅读和沟通。例如，下表展示了不同模型方案的主要指标及其95%置信区间：

模型/方案	AUC (↑)	LogLoss (↓)	ECE (↓)	RPM (¥, ↑)	CTR@Win (↑)
基线: LR	0.740 ± 0.002	0.505 ± 0.003	0.050 ± 0.005	2.50 ± 0.04	1.20% ± 0.05%
+ DeepFM (全特征)	0.755 ± 0.002	0.483 ± 0.002	0.045 ± 0.004	2.60 ± 0.04	1.30% ± 0.05%
+ 概率校准	0.754 ± 0.002	0.484 ± 0.002	0.020 ± 0.002	2.58 ± 0.04	1.29% ± 0.05%
+ 多任务 CTR+CVR	0.750 ± 0.003	0.490 ± 0.003	0.047 ± 0.004	2.70 ± 0.05	1.25% ± 0.05%

示例：表中“±”后的数为95%置信区间半宽度（用bootstrap 1000次计算）。箭头↑↓表示指标越大或越小越好。可以看到DeepFM比LR的AUC有显著提升，校准虽对AUC影响极小但大幅降低了ECE，多任务模型CTR-AUC略降但RPM最高（提升了收入）。这些差异均在置信区间上具有统计显著性。

- **Bootstrap显著性检验**：如上表，我们对每个实验结果附加了置信区间。这通过bootstrap实现：对测试集或模拟日志反复采样，计算指标分布。例如采样1000次计算AUC得到分布后取其2.5%和97.5%分位数作为CI。对比两个方案时，我们也可直接对它们的差值做bootstrap以估计p-value。如果某指标差异的95% CI仍全为正或负，则表示差异有>95%置信度为真提升/下降。反之若CI跨0则认为差异不显著。这样的统计检验使我们的分析更加**可信**：面试中可以强调，“我们的每次实验都经过了显著性分析，确保提升不是偶然噪声造成的”。
- **可重复实验配置**：我们严格版本化每次实验的配置（模型类型、特征集合、随机种子等），保证结果可**重复验证**。例如输出日志中保存了配置哈希或git commit id。这体现了科研严谨性。此外，如果某次实验结果异常，我们可以依据配置快速复现环境排查问题。

- **错误分析和诊断：**当新方案不如基线时，我们会深入分析原因。例如构造分桶AUC分析看看是否某些流量段退化，或者查看校准曲线发现新方案过度校准导致低估高分段概率。通过这些**诊断工具**，确保我们从每次实验中学到东西，而不仅是报几个数字。这样的反思过程也是可以在面试交流中分享的，突出你解决问题的思路。

综上，我们的实验评估策略追求**科学、全面、可信**。不仅证明某方案好坏，更要让读者信服其统计显著性，并理解改进背后的原因。这种态度在面试官看来非常可贵。

预期结果示例

经过上述阶段的开发与实验，我们预计平台将产生一系列丰富的成果，包括模型性能提升、指标对比分析和直观的可视化。以下列举几个关键结果示例：

- **模型性能提升曲线：**随着阶段推进，CTR模型的AUC从基础LR的大约0.74逐步提高。DeepFM引入后AUC或许提升2-3个百分点达到0.77左右；加入跨天特征可能再提升1-2个百分点；多任务模型可能对CTR-AUC略有损伤但在转化收益上获得提升。在报告中我们可用折线图展示AUC随阶段的变化，凸显每一步改进的贡献。
- **可靠性校准图：**校准模块输出一张可靠性曲线图。理想情况下，未校准模型的曲线在高预测值段明显偏离对角线（如预测0.9实际0.8），而校准后曲线更贴近45度线。这张图能够说明校准的效果一目了然。
- **竞价模拟收益表：**经过排序竞价仿真，不同模型方案的RPM、CTR@Win等指标汇总如前述表格所示。比如DeepFM比LR RPM提高4%，多任务虽然CTR略低却RPM最高提升8%。这些数据直接证明了算法改进带来的**业务价值**。我们也可以绘制柱状图对比不同方案的RPM，或用分组柱形图比较它们在高曝光量广告和尾部广告上的效果差异。
- **Off-policy评估验证：**一组对比离线评估和在线模拟的结果。例如表明某新模型真实模拟CTR提升+5.0%，IPS估计提升+4.8%（误差在置信区间内），增强了对离线评估方法准确性的信心。这种结果说明我们的评估工具可以用于筛选潜在优秀方案，减少线上试错。
- **LLM Embedding效果：**对于冷启动广告集合，加入LLM特征后CTR预估准确率提升的例子。如某些几乎无历史点击的新广告，之前模型基本靠平均值猜测，有了文本语义embedding后能捕捉其潜在类别相关性，点击预测更接近真实。可以展示两个模型在冷启动集上的PR曲线或AUC值，证明LLM特征的帮助。同时也指出在已有充足历史的数据上，LLM特征贡献有限，以突出其适用场景。

所有这些结果将在我们的最终报告中体现，并辅以解释说明。通过样例结果，面试官/读者能直观感受到该平台的强大之处：**不仅产出模型算法本身的改进，还把这种改进通过仿真定量地映射到了商业指标上**，真正做到“算法驱动业务”。

可扩展方向

由于时间和范围所限，本实验平台已有的模块仍有许多可以进一步拓展和优化的空间。在实际应用和深入研究中，可以考虑以下扩展方向：

- **引入更多前沿CTR模型：**如将DCNV2（Deep Cross Network v2）集成到模型库，用以探索更复杂的特征交叉建模效果；或者实现AutoInt、自注意力等机制增强模型的表达能力。每新增一种模型，都可以在我们的对照框架下评估其相对收益，从而不断丰富对比基线。

- **丰富多任务和多目标**：除了CTR和转化，多任务框架可扩展为预测更多目标，例如广告曝光的长时价值LTV、用户留存等。或者把当前二任务扩展为**序列任务**（先预估点击再预估转化再预估价值），形成整个漏斗转化率链条的预测体系。这可以模拟更复杂的商业KPI优化（比如优化GMV而非单纯CTR）。
- **强化学习和智能出价策略**：将当前的Bandit探索扩展为更复杂的**强化学习策略优化**。例如使用策略梯度或DQN算法，基于竞价模拟环境训练一个**出价策略智能体**，使其学会在不同情况下动态调整出价以最大化收益或预算利用率。这相当于实现一个简化版的广告出价策略优化，连接推荐/广告和强化学习前沿。
- **数据规模与分布式**：目前实验用的数据规模较小，但完全可以扩展到更大数据集，甚至引入分布式训练加速。例如集成Parameter Server或Horovod，使模型训练能在多机多卡环境下运行，从工程角度模拟工业大流量场景。也可考虑实时流水模拟（真正的流式数据处理），提高平台的逼真度。
- **个性化及用户细分**：引入用户画像细分或Contextual bandit思想，在策略层针对不同用户群体采用不同模型或参数，从而模拟个性化广告投放策略。比如实现一个简单的**策略集合**，针对新用户用偏探索的模型，老用户用成熟模型，并验证这样做是否提升总体表现。
- **更完备的指标体系**：加入更多评估指标，如NDCG（归一化折损增益）用于评估排名质量，或Calibration方面引入分组Calibration、ACE等指标。还可以设计模拟的**用户体验指标**（如频次控制下的用户点击率），进一步考察策略在不同利益相关者（广告主、用户、平台）之间的平衡。
- **UI界面与可视分析**：如果希望演示效果更直观，可以为平台开发一个简单的Web界面或Notebook交互界面，展示实时的训练曲线、指标变化和模拟竞价动画等。这样在面试时可以Demo，加分效果极佳。

以上扩展方向表明，本项目有很大的成长空间，完全可以根据求职需要或个人兴趣不断丰富。这也体现了平台“模块化、配置化”的设计初衷：任何新的想法都能较容易地融入现有框架进行试验。对于面试官的提问“如果给你更多时间你会怎么做”，你完全可以从上述多个角度展开，展示你的远见和对该领域的持续热情。

简历项目描述（示例）

广告系统算法实验平台：我独立设计并实现了一个广告算法实验平台，涵盖从CTR预估到竞价排序的业务闭环。项目分阶段引入了逻辑回归、DeepFM 深度模型、多任务 CTR-CVR 联合预测等模块，并加入了概率校准提升预测可靠性。通过自建仿真竞价沙盒将模型改进与业务指标（RPM、CTR@Win、转化率）挂钩评估，支持 A/B 对照和 bootstrap 显著性检验。平台还探索了**Off-policy**离线评估（IPS/DR）方法验证新策略效果，以及利用 BERT 等大模型 Embedding 改善冷启动推荐。该项目体现了算法工程全链路思维和扎实的实现能力，可有效支撑复杂广告策略的验证，在实习面试中获得了面试官的高度评价。