

# 个性化新闻标题生成模型训练思路与优化建议

---

## 1. 用户编码器训练 (User Encoder)

**目标：**学习用户兴趣表征，为后续个性化生成提供用户向量。

**主要流程：**

- **数据准备：**
  - 输入：新闻类别、标题、正文的索引化表示 (`news_vert`, `news_title`, `news_body`)，用户点击历史 (`TrainUsers`)、点击样本 (`TrainSamples`)。
  - 数据由预处理脚本生成。
- **模型结构：**
  - 采用 NRMS/NAML 等结构，分别对新闻和用户进行编码。
  - 新闻编码器将新闻内容转为向量，用户编码器聚合用户历史点击新闻的向量。
- **训练目标：**
  - 以点击预测为目标，优化交叉熵损失，提升用户兴趣建模能力。
  - 训练完成后，保存新闻和用户的向量 (`news_scoring`, `global_user_embed`)，用于生成模型。
- **评估指标：**
  - AUC、MRR、nDCG、CTR 等推荐系统常用指标。

---

## 2. 个性化生成器训练 (Personalized Generator)

**目标：**根据新闻正文和用户兴趣，生成个性化新闻标题。

**主要流程：**

- **数据准备：**
  - 输入：新闻正文 (`sources`)、标题 (`target_inputs`, `target_outputs`)、用户兴趣向量 (`global_user_embed` 或用户点击历史编码)。
- **模型结构：**
  - 基于 Seq2Seq (LSTM/Transformer 编码器+解码器)，可选 Pointer-Generator 机制。
  - 用户兴趣向量作为条件输入，影响解码过程。
- **训练流程：**
  1. **预训练：**用标准 Seq2Seq 目标 (最大似然/交叉熵) 训练生成器，提升基本生成能力。
  2. **个性化训练：**引入用户兴趣向量，采用强化学习 (如 A2C) 优化个性化目标 (如覆盖率、流畅性、个性化得分等)。
    - 奖励函数综合了覆盖率 (ROUGE)、个性化相关性、流畅性等。

- 用户兴趣向量由用户编码器输出，新闻正文向量由新闻编码器输出。

- **评估与测试：**

- 采用 ROUGE 等指标评估生成标题的质量和个性化程度。
- 

### 3. 可能的优化方法

#### 1. 用户兴趣建模优化

- 引入更多行为特征（如时间、位置、设备等）。
- 尝试更复杂的用户建模结构（如 Transformer-based 用户建模）。

#### 2. 生成器优化

- 使用更强大的预训练语言模型（如 BERT、T5 等）作为编码器或解码器。
- 引入多任务学习（如同时优化点击率和生成质量）。
- 增加多样性奖励，避免生成模板化标题。

#### 3. 训练策略优化

- 采用 Curriculum Learning，先训练基础生成，再逐步引入个性化和强化学习目标。
- 使用更细致的奖励函数设计，区分不同类型的个性化需求。

#### 4. 数据增强与负采样

- 对用户历史、新闻正文进行数据增强，提升模型泛化能力。
- 负采样策略优化，提升训练效率和效果。

#### 5. 推理阶段优化

- 引入多样化解码（如 Top-k、Top-p 采样）提升生成多样性。
  - 结合用户实时反馈动态调整生成策略。
- 

### 总结：

该模型通过“用户兴趣建模+个性化生成”两阶段训练，结合推荐与生成思想，实现了个性化新闻标题生成。优化空间主要在用户兴趣建模、生成器结构、奖励函数设计和训练策略等方面。