

大模型抽取知识点的知识图谱及推荐系统

7/25

在之前demo的代码基础上，对全部商品（26k+）进行商品知识点抽取。

遇到如下的问题：

虽然已经限制了知识点的类型，但未限制知识点的内容。目前对于一类知识点，如Functionality知识点下，大模型对同一含义存在多种表述。需要将其统一为一定数目的知识点内容，以能够在图谱中正确地连接。

在demo上，由于商品数量较少（仅200+），对于知识点下的内容，采用了BERTopic+HDBSCAN的方式进行聚类，通过调整聚类的超参数，可以得到一定数量的知识点内容簇。经人工审查，聚类效果较好。

在全部数据上，商品总数很大。同样的聚类方式，同一个超参数在不同的知识点下表现不同。存在部分聚类效果很差（将上千个词语聚为20+簇，但同一簇差异过大、或聚类为几百簇，类别数量过多等等）

考虑可能是以下的原因：

1. BERTopic模型能力不足，得到的embedding表达能力不够
2. HDBSCAN聚类方法无监督，难以得到预期数量的聚类个数

此外，如果对每个类、每个商品生成知识点后，都采用如上方式聚类，会导致有新商品进入时，需要完全更新聚类结果（有可能产生完全不存在的新簇、或改变原来各个簇的结果），与整体设想的“增量更新”优势相悖。

希望讨论改进方向：

1. 重新选择embedding模型/聚类方法，采用更强的模型/有监督的聚类方法生成知识点，如BGE模型+K-means等
 - a. 进展：尝试了基于BGE模型的embedding，效果相对BERTopic较优，但仍然存在部分类似的问题

b. 进展：K-means等聚类方法效果较差，目前选择K=10、20时，均会有被强制分为一类的差异较大的知识点

2. 和限定知识点类型类似地，限定各类下知识点内容，如对更多（比如全部商品的5%）商品进行不加限制地生成，然后对这些商品的知识点采用目前的方式聚类，通过调参得到每个类上一定数目的知识点类型&内容，在全量生成时，限制为/投影为已经确定的知识点类型&内容

a. 进展：正在尝试中。目前选择了1000个商品

7/17

目前已经完成了对于demo（选取Amazon VideoGames数据集，随机20个user的全部交互记录，以及对应商品）的

1. 商品知识点抽取
2. 用户兴趣抽取
3. 基于大模型抽取知识图谱的图推荐

的全流程实验。（仅作为代码可跑通的测试，未记录最终推荐结果等指标）

其中，在知识点抽取方面，针对大模型生成结果做了如下的限制：

1. 限制商品知识点类型为10类（方法：先不限制地生成知识点类型，对所有类型进行聚类，选取最有代表性的9类+Others类 重新生成）

```
ALLOWED_TYPES = {  
    "Functionality",  
    "Compatibility",  
    "Components",  
    "User Demographics",  
    "Usage Scenarios",  
    "Game Features",  
    "Technical Specifications",  
    "Content Additions",  
    "Performance Metrics",  
    "Others"  
}
```

2. 限制用户兴趣内容属于商品知识点内容（方法：先不限制地生成用户兴趣内容，再将内容投影到距离最近的知识点内容上）

↑ 聚类使用的embedding来自 BERTopic模型，使用HDBSCAN无监督聚类。

（“最有代表性”指的是 聚类结果合理&聚类各簇覆盖较广，也就是 有较多的商品都有属于这类的知识点，且同聚类知识点含义相近、不同类知识点含义不同）

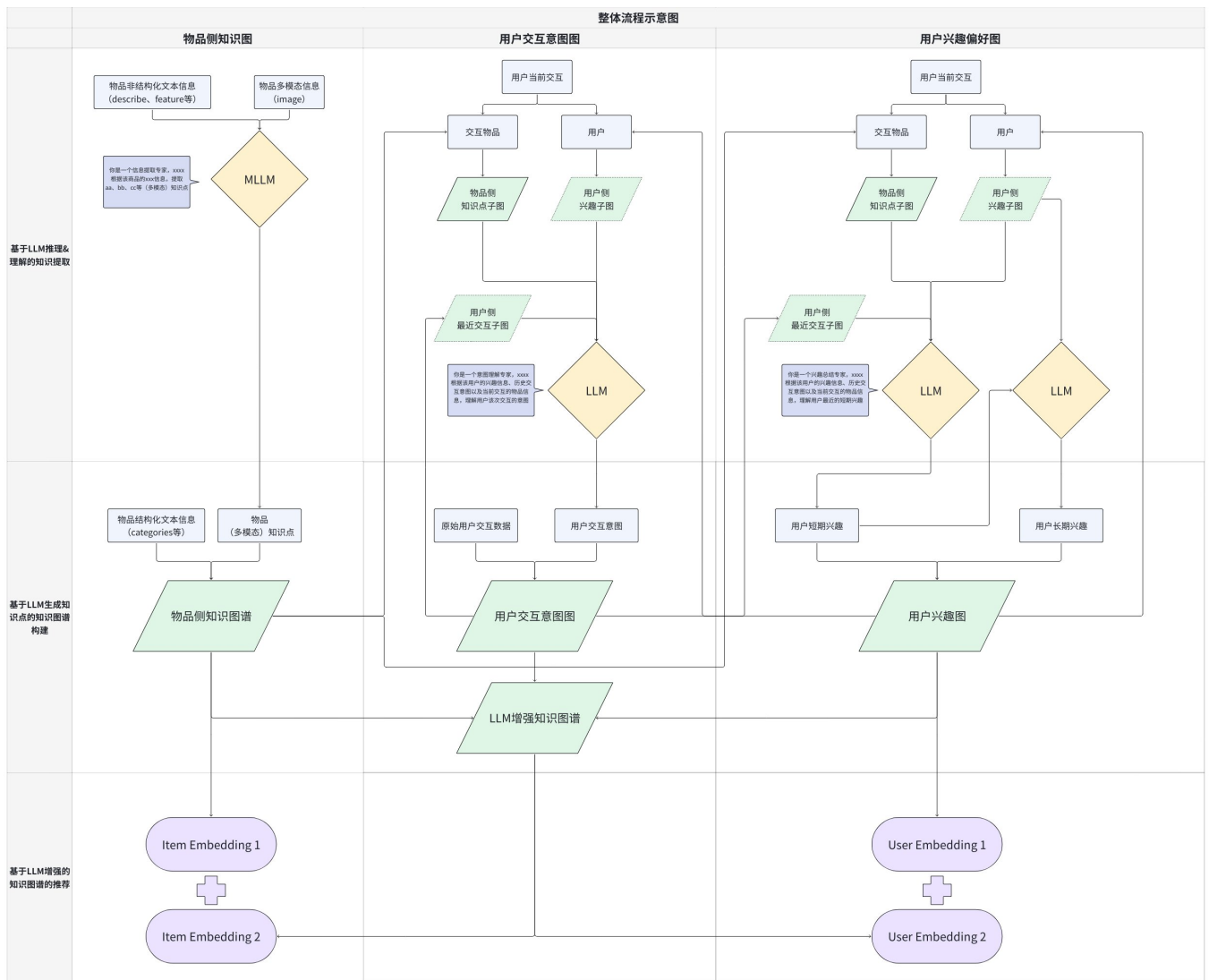
使用大模型生成用户兴趣时，以如下的方式使用了Graph RAG：

先基于生成的商品知识点构建商品侧知识图谱，对于用户的每次交互，使用交互的物品在商品侧知识图谱中查询对应的子图、使用用户在用户侧兴趣图中查询对应的子图。将两个子图以文本的形式作为Prompt的一部分输入LLM，进行基于Graph RAG的用户兴趣抽取。

计划对全部商品（26k+）做同样工作。预期下周五可以获得商品侧知识图谱，并开始构建用户兴趣图谱

7/3

整体流程示意



基于Graph Rag的用户兴趣更新&交互意图理解

7/1

在之前的尝试的基础上，更新了想法，主要有：1.不仅抽取物品侧的多模态知识，也抽取用户侧的“兴趣偏好”以及u-i交互的“交互类型”，构建三张基于LLM抽取知识的图，并综合为一个知识图谱。2.得到知识图之后，基于知识图的特性，设计两阶段图学习策略，进行基于图的推荐

方法构想汇报总结：

从LLM抽取的知识点图谱到结构可解释的推荐建模

1. 背景与问题出发点

当前主流的 LLM4Rec 工作大多以两种方式使用大模型：

- **Embedding 提取器**：让 LLM 输出用于表示的向量，再送入下游推荐模型；
- **直接生成推荐结果**：基于prompt直接输出物品序列。

但这些方式忽略了LLM在语言理解、推理和知识抽取方面的能力，且对推荐过程的可解释性帮助有限。

2. 我们的方法出发点

我们希望利用LLM的推理和语言抽象能力，构建出包含用户兴趣、物品属性、交互动机等信息的**结构化知识图谱**，并进一步用于推荐建模。

但不是让LLM参与整个推荐流程，而是只在**离线阶段用于图谱构建**，后续推荐模型完全使用**轻量图神经网络结构**完成，保持训练和部署效率。

3. 图谱构建阶段（离线，使用LLM）

构建的图谱结构由三部分组成：

1. **用户侧图**：用户与其兴趣标签（长期/短期）之间的边，如 “UserA → 长期兴趣 → 剧情片”；
2. **物品侧图**：物品与其属性知识点之间的边，如 “MovieX → 拍摄风格 → 蒙太奇”；
3. **用户-物品交互图**：保留用户行为类型（如 “喜欢” “不喜欢” “点击但未完成观看”），并明确区分**正负边**。

我们对 LLM 输出的知识点做语义聚类归一，使用 schema 限定边类型，并加入**置信度判定模块**，筛选出结构清晰、高可信度的子图用于下游建模。

4. 推荐模型设计阶段（在线，使用小模型）

为了更有效地利用我们构建的结构化图谱，我们采用**双阶段图学习策略**：

► 第一阶段：三张图分别建模

- 用户侧图：学习基于兴趣标签聚合的用户表示；
- 物品侧图：学习基于属性标签聚合的物品表示；
- 用户-物品交互图：用于监督学习，尤其是对**正负交互进行区分建模**。

► 第二阶段：统一融合图建模

- 构建一张包含所有实体（user/item/knowledge point）的融合图；
- 使用图神经网络进行进一步传播和建模；
- 利用知识点连接的路径提升 user-item 之间的表示交互。

5. 模型增强设计

为了更好地学习结构中的语义信息，我们设计了以下机制：

✓ Relation-aware attention

- 对不同类型边赋予不同的注意力权重，正向兴趣传播增强，负向行为传播抑制或反向作用；
- 同时保留路径解释能力，例如“因为你不喜欢 A，所以未推荐 B”。

✓ 对比学习目标

- 显式构造 (user, pos_item, neg_item) 三元组；
- 推动模型拉近兴趣一致项、推远负反馈项。

✓ 图结构嵌入损失

- 对于正边，鼓励节点embedding靠近；
- 对于负边，鼓励拉远；
- 对于知识点，聚合其连接的user/item形成聚类中心。

6. 可解释性设计

- 每个推荐结果都可以基于路径进行解释：

例如：“推荐《少年时代》是因为你近期关注‘成长题材’，该片属于‘成长题材’和‘纪录片风格’。”

- 也可以输出用户兴趣标签占比，构成用户画像。

7. 方法优势与创新点

- 利用 LLM 的强推理能力抽取**结构化标签型知识点**，替代传统embedding；
- 构建出**多结构分离 + 联合融合**的图谱建模框架；
- 引入**边类型感知机制和对比损失**，提升用户兴趣建模精度和方向性；
- 保留推荐过程的**结构路径可解释性**，实现“看得懂”的推荐理由。

8. 当前进展与下一步计划

- 已完成知识点抽取和结构图谱设计方案；
- 下一步准备构建初步图谱、搭建三图建模和对比学习框架；
- 后续将结合路径可解释性做case study，并评估融合图带来的性能增益。