

# 思路

## 大致思路

1. 创建feature-label数据集
2. 利用监督学习方法，训练entity-recognition“分类器”
3. 把分类器整合到聊天里
4. 测试聊天

### 数据集特征

- Bell商家
  - 三个标签：feature，status，location
  - 少数商家有recommendation
  - 只有一个层级：餐厅
- 校内餐厅
  - 分层级：大食堂-小餐厅
  - 有些商家又是独立的，比如博雅轩
- 校内设施
  - 分层级：大楼层-小设施

## Intent

1. 询问

## Entity

1. 大一级设施：B1，A3，D6.....
2. 小一级设施
  1. 具体的某家餐厅
  2. 具体的某个地点，比如图书馆
3. 各种学校事务操作方法
  1. 选课，请假，等等（一个how to就是一个entity）
4. 各种生活小问题
  5. 一个how to

## 整理完毕的思路

非常棒，你的设计思路整体上非常合理且现代化，采用了LLM辅助开发的流程，可行性很高。这套方案覆盖了从数据准备到模型训练再到最终响应的完整闭环。

下面我将针对你提出的每一步进行分析，指出其中的**合理之处**以及一些**可以优化或需要特别注意**的地方。

---

## 方案评估与优化建议

### 1. 爬取信息，构建JSON数据库

- **合理之处** 👍：这是构建检索式问答系统的根基。使用JSON格式非常灵活，便于读取和检索。将信息分类（餐厅、教学楼、常见问题）也是很好的实践。
- **可优化之处** 💡：
  - **结构深度**: 设计JSON结构时，尽量让其具有良好的层级关系。例如，一个“教学楼”对象下可以嵌套“楼层分布”、“开放时间”、“可用自习室列表”等子对象，这样便于精确检索。
  - **别名与关键词**: 为每个实体（如“A3图书馆”）增加一个 `aliases` 或 `keywords` 字段，包含所有可能的别称（如“新图”、“A3”、“三馆”）。这在实体识别和后续检索时非常有用。
  - **非结构化数据**: “学习常见问题”属于非结构化数据(FAQ)。你可以将其整理成“question-answer”对的列表。后续可以采用更高级的向量检索（语义搜索）来匹配用户问题，而不仅仅是关键词检索。

### 2. 汇总实体，让大模型生成可能的Intent

- **合理之处** 👍：这是一个非常聪明的“自下而上”的方法。让LLM基于你已有的实体数据来“脑补”用户可能会问什么，可以快速生成大量高度相关的意图，避免了人工冥思苦想。
- **可优化之处** 💡：
  - **补充通用意图**: LLM基于你的数据可能只会生成与数据直接相关的意图（如 `query_location`）。但一个完整的Chatbot还需要处理很多通用对话意图，例如：**问候 (greeting)**、**感谢 (thank\_you)**、**告别 (goodbye)**、**寻求帮助/询问能力 (ask\_capabilities)**、**无法回答/兜底 (fallback)**、**闲聊 (chitchat)** 等。你需要手动补充这些通用意图。
  - **意图粒度**: 注意LLM可能生成过于宽泛或过于细致的意图。例如，`query_A3_library_hours` 和 `query_gym_hours` 最好合并为统一的 `query_facility_hours` 意图，将具体设施名作为实体。你需要对LLM生成的结果进行人工筛选和归纳。

### 3. 让LLM生成Feature-Label数据用于监督学习

- **合理之处** 👍：这是整个方案中最能体现效率提升的地方。相比纯人工标注，使用LLM生成训练数据可以极大地缩短开发周期。
- **可优化之处** 💡：

- **数据质量和多样性是关键:** LLM生成的数据可能缺乏多样性，句式单一。你需要设计多样化的Prompt，要求LLM生成不同语气（正式、口语化）、不同句式（陈述句、疑问句）、甚至包含常见错别字或简化表达的样本。
- **!! 人工审核与修正 (Human-in-the-loop):** 这是**至关重要**的一步！LLM生成的数据绝不能直接使用。你必须进行抽样检查或完整地过一遍，修正其中错误的标注。一个包含少量高质量、经人工修正的数据集，远比一个包含大量低质量自动生成数据的数据集要好。**Garbage in, garbage out.**
- **负样本:** 不仅要生成能匹配到意图的“正样本”，也要生成一些**不属于任何意图**的“负样本”(即out-of-scope)，并将它们统一标注为 oos 或 None 意图。这能让你的模型学会拒绝回答超出其能力范围的问题。

## 4. 分别训练实体识别和意图识别的模型

- **合理之处** 🍌: 将意图分类和实体识别作为两个独立的任务是NLU（自然语言理解）的标准做法，便于独立优化和评估。
- **可优化之处** 💡:
  - **考虑使用框架:** 与其完全从零开始写训练脚本，不如考虑使用现成的对话AI框架，如Rasa。Rasa的设计哲学与你的方案高度契合，它内置了处理意图和实体的流水线（NLU Pipeline），你可以方便地插入自己训练的模型或使用其预置模型。它能帮你管理训练数据、训练流程和模型版本，事半功倍。
  - **联合训练:** 虽然是两个任务，但一些先进的模型架构（如基于Transformer的DIETClassifier）可以同时意图和实体进行联合训练，有时能取得更好的效果，因为两个任务可以互相提供有用的信息。Rasa就支持这种模型。

## 5. 设计一个回应函数

- **合理之处** 🍌: 函数的输入参数（实体、意图、情感）设计得很好，考虑到了核心信息和个性化表达。
  - **可优化之处** 💡:
    - **对话管理 (Dialogue Management):** 这是你当前方案中隐藏的最复杂的部分。你的回应函数不仅要能“单轮”回答，还要能处理更复杂的对话逻辑。
      - **澄清式提问:** 如果意图是 query\_facility\_hours，但实体 facility\_name 缺失，函数不应直接回答“不知道”，而应反问：“请问您想查询哪个设施的开放时间呢？”
      - **上下文记忆:** 如果用户先问“A3图书馆在哪？”，你回答后，用户接着问“那它几点开门？”，你的函数需要能“记住”上文的实体“A3图书馆”。
      - **处理歧义:** 如果用户问“机房在哪？”，而学校有多个机房，函数应列出选项让用户选择，而不是随意返回一个。
    - **响应模板:** 不要硬编码(hardcode)回答。可以设计一个响应模板库，根据不同的意图和实体组合，动态地填充信息生成自然语言，并根据情感参数调整语气词（如“很高兴为您服务！” vs “别着急，我帮您查到了...”）。
-

# 总结与核心建议

你的方案非常出色，有很强的前瞻性。核心建议是：

- 1. **重视数据质量**: 在第1步和第3步，投入更多精力在JSON的结构设计和LLM生成数据的审核修正上。
- 2. **引入对话管理**: 在第5步，将简单的“回应函数”升级为具备上下文记忆和澄清反问能力的“对话管理器”。
- 3. **拥抱开源框架**: 强烈建议研究一下 **Rasa** 框架。它能为你当前的方案提供一个健壮的骨架，让你更专注于业务逻辑和数据，而不是重复造轮子。

这是一个非常有前景的方案，祝您项目顺利！

## 信息检索

以entity为主键重构知识库

`{}` JSON

```
1  "starbucks_campus_a": {
2      "name": "星巴克 (A区店)",
3      "type": "restaurant",
4      "location": "校园A区，图书馆一楼旁边",
5      "hours": { "weekdays": "08:00 - 20:00", "weekends": "09:00 - 19:00" },
6      "features": "提供在线点单和自习座位。"
7  }
```