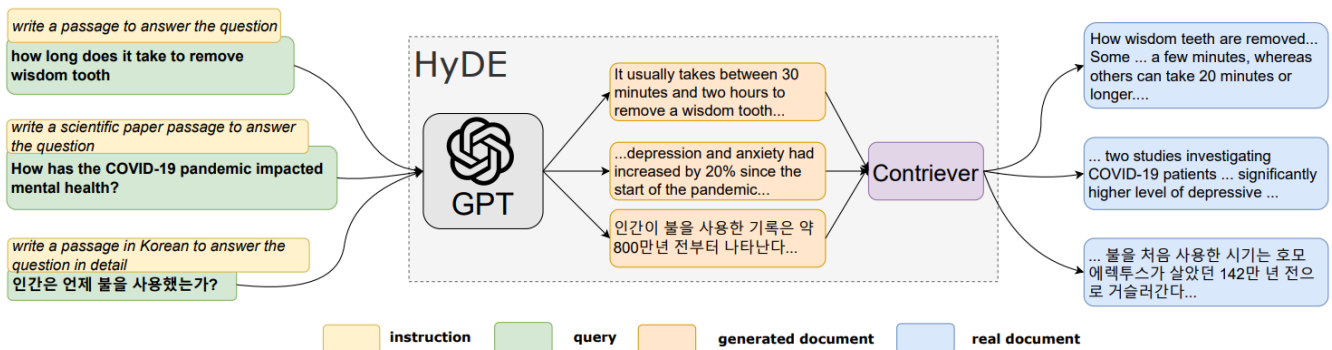


Langchain 0.2 RAG 快速入门 #6

HyDE

Hypothetical Document Embeddings

HyDE不是基于原始问题生成查询，而是专注于为给定的查询生成假设性文档。生成这类假设性文档背后的直觉是，它们的嵌入向量可以用来在语料库嵌入空间中识别一个邻域，基于向量相似性检索出相似的真正文文档。在这种情况下，RAG将能够根据假设性文档检索到更多相关的文档，从而准确回答用户查询。



生成假设文档

```
1 from langchain.prompts import ChatPromptTemplate
2
3 hyde_prompt = ChatPromptTemplate.from_template(
4     """
5     请撰写一段科学论文的内容来回答以下问题：
6     问题：{question}\n
7     段落：
8     """
9 )
10
11 generate_doc_chain = (
12     {'question': RunnablePassthrough()}
13     | hyde_prompt
14     | glm4_air_model
15     | StrOutputParser()
16 )
17 question = "工作满意度指数有哪些？"
18 res = generate_doc_chain.invoke(question)
```

```
19
20 retrieval_chain = generate_doc_chain | retriever
21 retrieved_docs = retrieval_chain.invoke({"question": question})
22 retrieved_docs
```

RAG chain

```
1 template = """Answer the following question based on the provided context:
2
3     {context}
4
5     Question: {question}
6     """
7
8     prompt = ChatPromptTemplate.from_template(template)
9
10    final_rag_chain = (
11        prompt
12        | glm4_air_model
13        | StrOutputParser()
14    )
15
16    res = final_rag_chain.invoke({"context": retrieved_docs, "question":
17        question})
17    print(res)
```

完整代码

```
1 import os
2 from uuid import uuid4
3
4
5 # 从我本地的一个文件，导入我设置的ZHIPU_AK
6 from ai_docs_cfg import ZHIPU_AK, MY_LANGCHAIN_API_KEY
7
8 from langchain_community.document_loaders import PyPDFLoader, DirectoryLoader
9 from langchain.text_splitter import RecursiveCharacterTextSplitter
10
11 from langchain_community.vectorstores import Chroma
12 from langchain_core.output_parsers import StrOutputParser
13 from langchain_core.runnables import RunnablePassthrough
14 from langchain_openai import ChatOpenAI
15 from langchain import hub
16 from langchain.load import loads, dumps
```

```

17 from typing import List
18
19 # Langsmith 配置, 不用可注掉
20 unique_id = uuid4().hex[0:8]
21 os.environ["LANGCHAIN_PROJECT"] = f" RAG - {unique_id}"
22 # os.environ["LANGCHAIN_TRACING_V2"] = 'true'
23 os.environ["LANGCHAIN_API_KEY"] = MY_LANGCHAIN_API_KEY
24
25 if __name__ == '__main__':
26 # 初始化 模型
27 glm4_air_model = ChatOpenAI(
28     model_name="glm-4-0520",
29     openai_api_base="https://open.bigmodel.cn/api/paas/v4",
30     openai_api_key=ZHIPU_AK,
31 )
32
33 from langchain_community.embeddings.huggingface import HuggingFaceEmbeddings
34
35 # 模型所在路径
36 embedding_model_name = 'D:\\LLM\\bce_models\\bce-embedding-base_v1'
37
38 # 使用 GPU
39 embedding_model_kwargs = {'device': 'cuda:0'}
40
41 # 使用 CPU
42 # embedding_model_kwargs = {'device': 'cpu'}
43
44 # 配置
45 embedding_encode_kwargs = {'batch_size': 32, 'normalize_embeddings': True, }
46
47 # 初始化
48 embed_model = HuggingFaceEmbeddings(
49     model_name=embedding_model_name,
50     model_kwargs=embedding_model_kwargs,
51     encode_kwargs=embedding_encode_kwargs
52 )
53
54 data_path = 'D:\\LLM\\my_projects\\new_langchain_test\\AI_docs\\data'
55 vdb_path =
56     'D:\\LLM\\my_projects\\new_langchain_test\\AI_docs\\data\\vectorstore'
57 loader = DirectoryLoader(data_path, glob="*.pdf", loader_cls=PyPDFLoader)
58 documents = loader.load()
59
60 # Split text into chunks
61 text_splitter = RecursiveCharacterTextSplitter(chunk_size=500,
62     chunk_overlap=20)
63 text_chunks = text_splitter.split_documents(documents)

```

```

63
64 # 创建向量数据库
65 vectorstore = Chroma(persist_directory=vdb_path,
66                       embedding_function=embed_model)
67
68 # 索引数据
69 #res = vectorstore.add_documents(text_chunks)
70
71 # 检索器
72 retriever = vectorstore.as_retriever(search_kwargs={'k': 5})
73
74
75 from langchain.prompts import ChatPromptTemplate
76
77 hyde_prompt = ChatPromptTemplate.from_template(
78     """
79     请撰写一段科学论文的内容来回答以下问题:
80     问题: {question}\n
81     段落:
82     """
83 )
84
85 generate_doc_chain = (
86     {'question': RunnablePassthrough()}
87     | hyde_prompt
88     | glm4_air_model
89     | StrOutputParser()
90 )
91
92 question = "工作满意度指数有哪些？"
93 res = generate_doc_chain.invoke(question)
94
95 retrieval_chain = generate_doc_chain | retriever
96 retrieved_docs = retrieval_chain.invoke({"question": question})
97
98 retrieved_docs
99
100 template = """Answer the following question based on the provided context:
101 {context}
102 Question: {question}
103 """
104
105 prompt = ChatPromptTemplate.from_template(template)
106
107 final_rag_chain = (
108     prompt
109     | glm4_air_model
110     | StrOutputParser()
111 )

```

```
110
111 res = final_rag_chain.invoke({"context": retrieved_docs, "question":
    question})
112 print(res)
113
```