

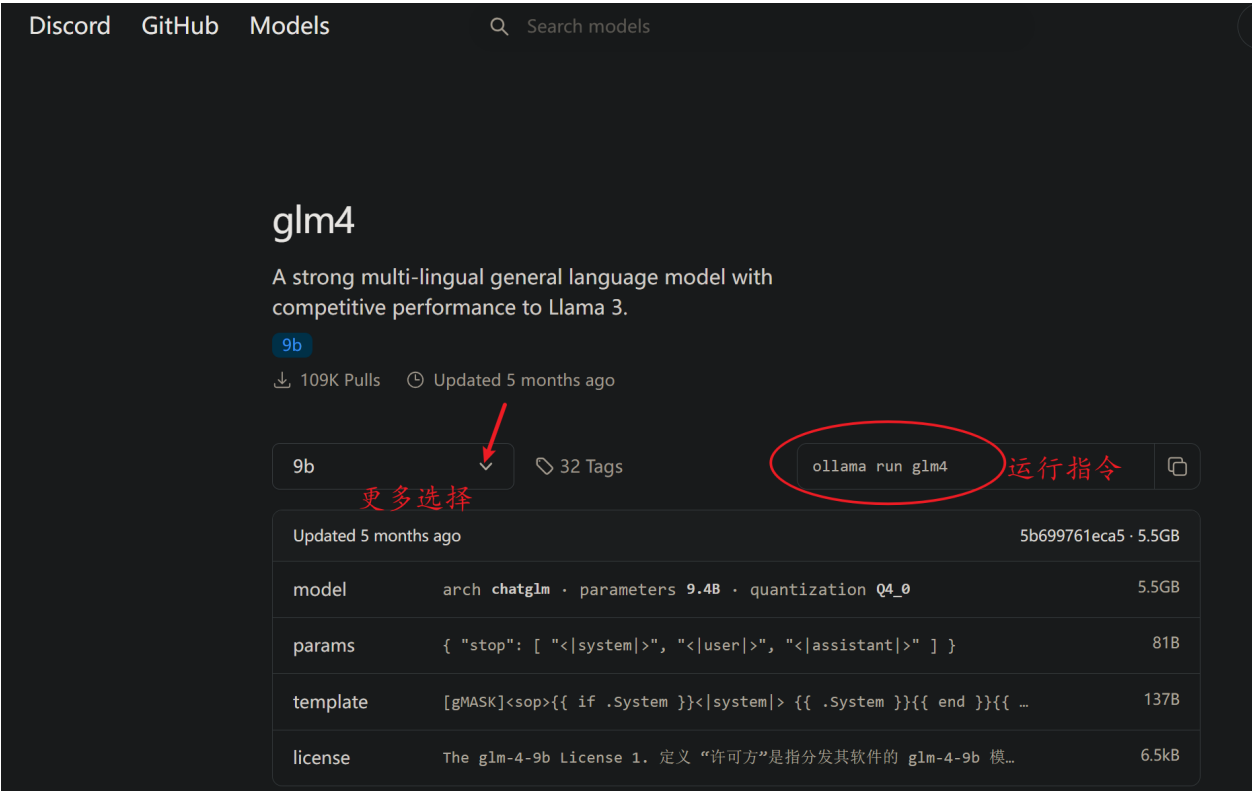
基于ollama的ragas评估指南

一、ollma环境部署

ollma对Windows环境支持比较好，安装运行也简单。当然在linux下会更灵活。安装后需要合理配置一些参数。

二、ollma模型下载

到[模型库](#)搜索关键词。找到目标模型后，可通过ollama pull 命令下载。注意至少一个语言模型和一个嵌入模型。



下载完成后，可以通过ollama list 命令显示当前可用的完整模型名称，这是在调用时必须提供的参数。

三、本地API调用

可以参考[官方文档](#)通过http请求进行最简单的实验，确保能正常工作。

ollama对openai风格接口[已经部分兼容](#)。部署好直接用openai库就可以调用。但为了便于ragas使用，我们用langchain_openai提供的包装进行测试。

```

1 from langchain_openai import OpenAIEmbeddings, ChatOpenAI
2 from langchain_core.messages import HumanMessage, SystemMessage
3
4 # 创建OpenAI客户端实例，指向本地服务
5 llm = ChatOpenAI( model="glm4:latest",
6     base_url='http://127.0.0.1:11434/v1/',
7     api_key="not-needed" # API key可以是任意值，因为没有认证
8 )
9 messages = [
10     SystemMessage(content="你是一个AI助手，善于解答各种问题"),
11     HumanMessage(content="美联储进入降息周期后，对世界金融市场各个品种有哪些影响？")
12 ]
13 response = llm.invoke(messages)
14 print(response.content)
15 one = OpenAIEmbeddings(model="aerok/xiaobu-embedding-v2:latest", deployment="aerok/xiaobu-embedding-v2:latest",
16     api_key='0', base_url='http://127.0.0.1:11434/v1/', )
17 out = one.embed_query(text="今天天气真热啊！")
18 print(out)

```

由于ollama对嵌入模型的openai风格接口仅支持输入文本而不支持输入token，所以需要修改langchain_openai关于Embedding的源码。主要是 `_get_len_safe_embeddings` 和 `_aget_len_safe_embeddings` 两个方法。尤其是后者为异步调用方式，很重要。将源码中两个函数重命名或者注释掉，只保留其中最有价值的代码段就可以工作了。具体如下：

```

1 def _get_len_safe_embeddings(
2     self, texts: List[str], *, engine: str, chunk_size: Optional[int] = None
3 ) -> List[List[float]]:
4     batched_embeddings: List[List[float]] = []
5     for i, text in enumerate(texts, 1):
6         valid = text if len(text) <= 510 else text[-510:]
7         response = self.client.create(
8             input= valid, **self._invocation_params
9         )
10         if not isinstance(response, dict):
11             response = response.model_dump()
12         batched_embeddings.extend(r["embedding"] for r in response["data"])
13
14     return batched_embeddings

```

```

14
15     async def _aget_len_safe_embeddings(
16         self, texts: List[str], *, engine: str, chunk_size: Optional[int] = Non
17         e
18     ) -> List[List[float]]:
19         batched_embeddings: List[List[float]] = []
20         for i, text in enumerate(texts,1):
21             valid = text if len(text) <= 510 else text[-510:]
22             response = await self.async_client.create(
23                 input= valid, **self._invocation_params
24             )
25             if not isinstance(response, dict):
26                 response = response.model_dump()
27             batched_embeddings.extend(r["embedding"] for r in response["data"])
28
29         return batched_embeddings

```

四、在ragas下的使用

具体代码如下：

```

1  import json
2  from langchain_openai import ChatOpenAI, OpenAIEmbeddings
3  from ragas import evaluate, RunConfig, EvaluationDataset
4  from ragas.llms import LangchainLLMWrapper
5  from ragas.embeddings import LangchainEmbeddingsWrapper
6  from ragas.metrics import (ContextPrecision, ContextRecall,
7
8                               Faithfulness, AnswerRelevancy, SemanticSimilarity)
9
10 llm = ChatOpenAI( model="glm4:latest",
11                  base_url='http://127.0.0.1:11434/v1/',
12                  api_key="not-needed" )
13 evaluator_llm = LangchainLLMWrapper(llm)
14 embed_model = OpenAIEmbeddings(model="aerok/xiaobu-embedding-v2:latest",
15                                 deployment="aerok/xiaobu-embedding-v2:latest",
16                                 api_key='0', base_url='http://127.0.0.1:11434/v
17 1/', )
18 evaluator_embeddings = LangchainEmbeddingsWrapper(embed_model)
19 metrics = [ Faithfulness(llm=evaluator_llm),
20             AnswerRelevancy(llm=evaluator_llm, embeddings=evaluator_embeddings),

```

```

19     ContextPrecision(llm=evaluator_llm),
20     ContextRecall(llm=evaluator_llm), ]
21
22 if __name__ == '__main__':
23     # 读取json文件中的数据，最好是字典的列表。
24     with (open('other/qainfo.json', 'r', encoding='utf-8')) as fr:
25         qainfo = json.load(fr)['results']
26         new = []
27         for k in qainfo:
28             one = {'user_input': k['query'], 'response': k['response'],
29                   'reference': k['gt_answer'],
30                   'retrieved_contexts': [one['text'] for one in k["retrieved_c
ontext"]]
31             }
32             new.append(one)
33     eval_dataset = EvaluationDataset.from_list(new[50:60]) # 选择需要评估的范
围。
34     print(eval_dataset) # 答应评估数据集的基本信息
35     results = evaluate(dataset=eval_dataset, metrics=metrics, batch_size=1,
36                       run_config=RunConfig(timeout=900))
37     print(results)
38 # {'context_recall': 0.860, 'factual_correctness': 0.498, 'faithfulness': 0.65
9}
39     df = results.to_pandas()
40     df.to_csv("result.csv", index=False)

```

其中数据集是参考ragas文档介绍准备的json格式。