

Langchain + ChatGPT

Geração Aumentada de Recuperação (RAG)

Fundamentos e Contextualização

- Biblioteca | Framework
- Simplificar o desenvolvimento de IA
 - Correntes
 - Links
- Componentes do LangChain
 - Interface LLM
 - Modelos de prompt
 - Módulos de recuperação
 - Agentes
 - Memória



Geração Aumentada de Recuperação [RAG]

Aplicações

- Agentes Inteligentes
- Sistemas de Perguntas e Respostas (Q&A)
- Chatbots Personalizados
- Aplicações de Resumo (Summarization)
- Automação de Fluxos de Trabalho
- Geração de Dados Sintéticos
- **Geração Aumentada de Recuperação (RAG, Retrieval-Augmented Generation)**

Motivação:

- Os modelos de LLM não conhecem seus dados
- Os aplicativos de IA devem aproveitar dados personalizados para serem eficazes

O que é:

- Abordagem de arquitetura tenta melhorar a eficácia das aplicações de grandes modelo de linguagem (LLM) usando dados personalizados. Isso é feito recuperando dados/documentos relevantes para uma pergunta ou tarefa e fornecendo-os como contexto para o LLM.

Geração Aumentada de Recuperação [RAG]

Que problemas resolve:

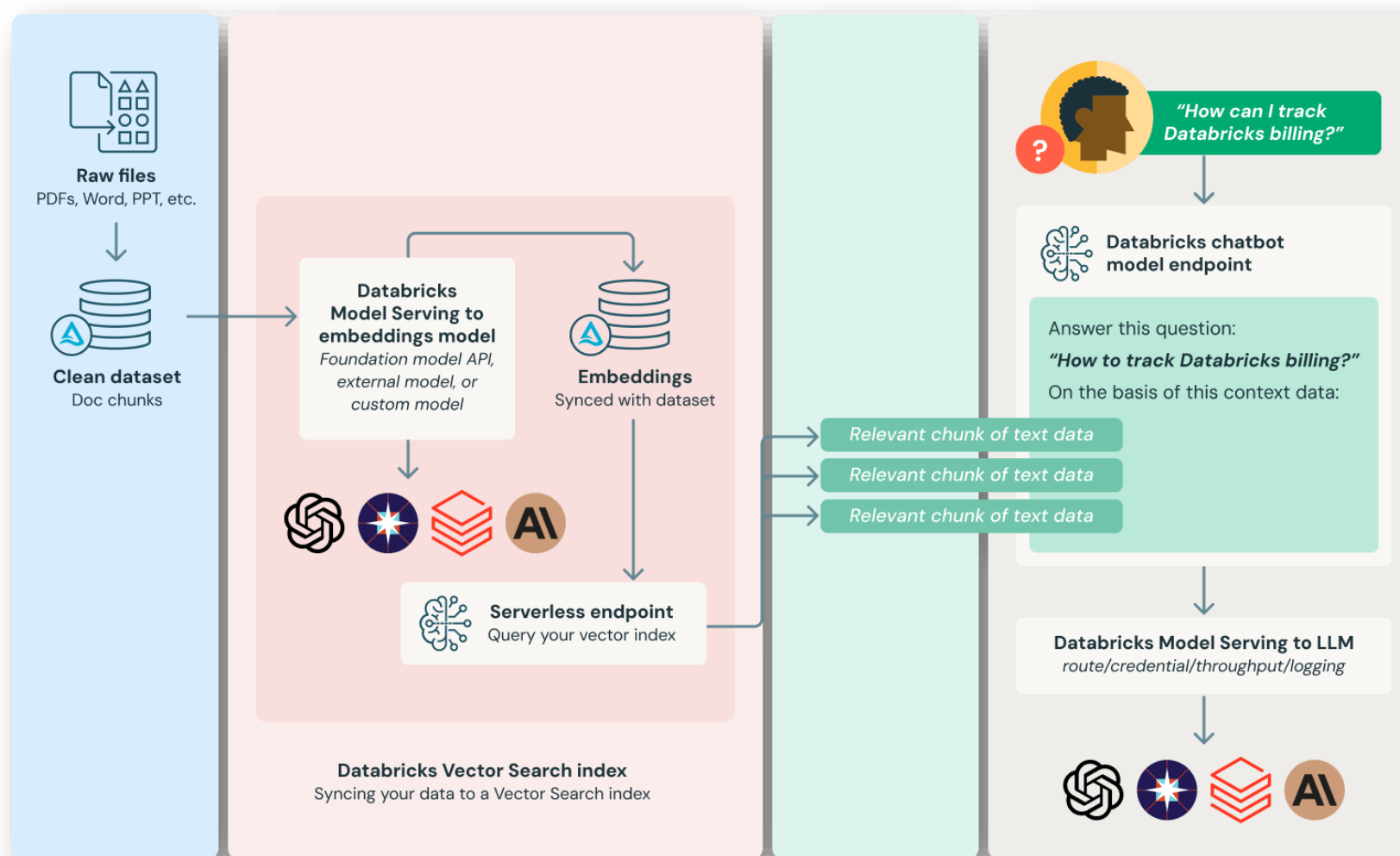
- Os modelos de LLM não conhecem seus dados
- Os aplicativos de IA devem aproveitar dados personalizados para serem eficazes
- Limitações de Conhecimento Estático
- Redução de Alucinações

Geração Aumentada de Recuperação [RAG]

Quais são os casos de uso da RAG :

- Assistentes de suporte ao cliente
- Análise de documentos corporativos
- Gestão de conhecimento empresarial
- Educação e treinamento
- Monitoramento e manutenção preditiva
- Melhorias em processos logísticos
- Q&A para novos funcionarios

Arquitetura



Geração Aumentada de Recuperação [RAG]

Código | Armazenando os dados

```
1 import argparse
2 import os
3 import shutil
4 from langchain_community.document_loaders import PyPDFDirectoryLoader
5 from langchain_text_splitters import RecursiveCharacterTextSplitter
6 from langchain.schema.document import Document
7 from get_embedding_function import get_embedding_function
8 from langchain_chroma import Chroma
9 from dotenv import load_dotenv
10 import openai
11
12 # Carrega as variáveis de ambiente. Assume que o projeto contém um arquivo .env com as chaves da API.
13 load_dotenv()
14 # ---- Define a chave da API da OpenAI
15 # Altere o nome da variável de ambiente de "OPENAI_API_KEY" para o nome especificado no
16 # seu arquivo .env.
17 openai.api_key = os.environ["OPENAI_API_KEY"]
18
19 CHROMA_PATH = "chroma"
20 DATA_PATH = "data"
21
22
23 def main():
24     # Verifica se o banco de dados deve ser limpo (usando o parâmetro --reset).
25     parser = argparse.ArgumentParser()
26     parser.add_argument(
27         "--reset", action="store_true", help="Reinicia o banco de dados."
28     )
29     args = parser.parse_args()
30     if args.reset:
31         print("💡 Limpando o banco de dados")
32         clear_database()
33
34     # Cria (ou atualiza) o banco de dados.
35     documents = load_documents()
36     chunks = split_documents(documents)
37     add_to_chroma(chunks)
```

```
1 # from langchain_community.embeddings.ollama import OllamaEmbeddings
2 # from langchain_community.embeddings.bedrock import BedrockEmbeddings
3 from langchain_openai import OpenAIEmbeddings
4 from dotenv import load_dotenv
5 import openai
6 import os
7
8 # Carrega as variáveis de ambiente. Assume que o projeto contém um arquivo .env com as chaves da API.
9 load_dotenv()
10 # ---- Define a chave da API da OpenAI
11 # Altere o nome da variável de ambiente de "OPENAI_API_KEY" para o nome especificado no
12 # seu arquivo .env.
13 openai.api_key = os.environ["OPENAI_API_KEY"]
14
15
16 def get_embedding_function():
17     # ---- Modelo da Amazon
18     # embeddings = BedrockEmbeddings(credentials_profile_name="default", region_name="us-east-1")
19
20     # ---- Modelos locais
21     # embeddings = OllamaEmbeddings(model="nomic-embed-text")
22
23     embeddings = OpenAIEmbeddings(model="text-embedding-3-large")
24     return embeddings
25
```


Código | Armazenando os dados

```
1 def add_to_chroma(chunks: List[Document]):
2     # Carrega o banco de dados existente.
3     db = Chroma(
4         persist_directory=CHROMA_PATH, embedding_function=get_embedding_function()
5     )
6
7     # Calcula os IDs das páginas.
8     chunks_with_ids = calculate_chunk_ids(chunks)
9
10    # Adiciona ou atualiza os documentos.
11    existing_items = db.get(include=[]) # IDs são sempre incluídos por padrão
12    existing_ids = set(existing_items["ids"])
13    print(f"Número de documentos existentes no banco de dados: {len(existing_ids)}")
14
15    # Adiciona apenas documentos que não existem no banco de dados.
16    new_chunks = []
17    for chunk in chunks_with_ids:
18        if chunk.metadata["id"] not in existing_ids:
19            new_chunks.append(chunk)
20
21    if len(new_chunks):
22        print(f"👉 Adicionando novos documentos: {len(new_chunks)}")
23        new_chunk_ids = [chunk.metadata["id"] for chunk in new_chunks]
24        db.add_documents(new_chunks, ids=new_chunk_ids)
25        db.persist()
26    else:
27        print("✅ Nenhum novo documento para adicionar")
```

```
1 def calculate_chunk_ids(chunks):
2     # Isso criará IDs como "data/monopoly.pdf:6:2"
3     # Fonte da Página : Número da Página : Índice do Fragmento
4
5     last_page_id = None
6     current_chunk_index = 0
7
8     for chunk in chunks:
9         source = chunk.metadata.get("source")
10        page = chunk.metadata.get("page")
11        current_page_id = f"{source}:{page}"
12
13        # Se o ID da página for o mesmo que o último, incrementa o índice.
14        if current_page_id == last_page_id:
15            current_chunk_index += 1
16        else:
17            current_chunk_index = 0
18
19        # Calcula o ID do fragmento.
20        chunk_id = f"{current_page_id}:{current_chunk_index}"
21        last_page_id = current_page_id
22
23        # Adiciona o ID aos metadados do fragmento.
24        chunk.metadata["id"] = chunk_id
25
26    return chunks
```

Código | Consulta dos dados

```
1 import argparse
2 from langchain_chroma import Chroma
3 from langchain.prompts import ChatPromptTemplate
4
5 # from langchain_community.llms.ollama import Ollama
6 from langchain_openai import ChatOpenAI
7
8 from get_embedding_function import get_embedding_function
9
10 CHROMA_PATH = "chroma"
11
12 PROMPT_TEMPLATE = ""
13 Responda à pergunta com base apenas no seguinte contexto:
14
15 {context}
16
17 ---
18
19 Responda à pergunta com base no contexto acima: {question}
20 ""
21
22
23 def main():
24     # Cria a interface de linha de comando (CLI).
25     parser = argparse.ArgumentParser()
26     parser.add_argument("query_text", type=str, help="O texto da consulta.")
27     args = parser.parse_args()
28     query_text = args.query_text
29     query_rag(query_text)
```

```
1 def query_rag(query_text: str):
2     # Prepara o banco de dados.
3     embedding_function = get_embedding_function()
4     db = Chroma(persist_directory=CHROMA_PATH, embedding_function=embedding_function)
5
6     # Pesquisa no banco de dados.
7     results = db.similarity_search_with_score(query_text, k=5)
8
9     context_text = "\n\n---\n\n".join([doc.page_content for doc, _score in results])
10    prompt_template = ChatPromptTemplate.from_template(PROMPT_TEMPLATE)
11    prompt = prompt_template.format(context=context_text, question=query_text)
12    # print(prompt, "\n\n---\n\n")
13
14    model = ChatOpenAI(model="gpt-4o")
15    response_text = model.invoke(prompt)
16
17    sources = [doc.metadata.get("id", None) for doc, _score in results]
18    formatted_response = f"Resposta: {response_text.content}\nFontes: {sources}"
19    print(formatted_response)
20    return response_text
21
22
23 if __name__ == "__main__":
24     main()
25
```

Resultado



```
1 >>> python query_data.py "De que inteligencias artificiais esse trabalho fala?"
2
3 Resposta: O trabalho mencionado fala sobre inteligências artificiais generativas,
4 como o ChatGPT e MidJourney, além de abordar conceitos como Aprendizado de Máquina
5 (Machine Learning) e Deep Learning. Além disso, menciona as Redes Neurais
6 Generativas Adversativas (GANs).
7
8 ---
9
10 Fontes: ['data/IMPACTO DA INTELIGENCIA ARTIFICIAL NA ENGENHARIA.pdf:27:1',
11          'data/IMPACTO DA INTELIGENCIA ARTIFICIAL NA ENGENHARIA.pdf:11:1',
12          'data/IMPACTO DA INTELIGENCIA ARTIFICIAL NA ENGENHARIA.pdf:47:0',
13          'data/IMPACTO DA INTELIGENCIA ARTIFICIAL NA ENGENHARIA.pdf:6:0',
14          'data/IMPACTO DA INTELIGENCIA ARTIFICIAL NA ENGENHARIA.pdf:21:0']
```