

RAG System with ChatGPT Documentation

Introduction

This **Python** code implements a Retrieval-Augmented Generation (RAG) system using OpenAI's language models (GPT-4o-mini in this case). The RAG setup allows a language model to retrieve relevant documents from a vector database before generating responses, which improves the accuracy and relevance of answers to user queries.

This code is a reimplementation of the solution discussed in this tutorial: [RAG With Llama 3.1 8B, Ollama, and Langchain: Tutorial](#). If you are interested in building an RAG system with **Llama** (an open-sourced LLM that can run locally on your device), please check that tutorial instead.

Prerequisites

- **OpenAI API Key:** This script uses the OpenAI API for processing the data and generating responses. Set up your API key as an environment variable: `OPENAI_API_KEY`.

Dependencies

To use this code, you will need to install the following libraries:

- `langchain`
- `openai`
- `pickle`
- `os`
- `time`

You can install these in the terminal by executing the command:

```
pip install langchain langchain-chroma langchain-community langchain-core  
langchain-openai langchain-ollama openai
```

Code Walkthrough

Import Libraries

```
1  import os
2  import pickle # Import pickle for serialization
3  import time # Import time for measuring execution time
4
5  from langchain.prompts import PromptTemplate
6  from langchain.text_splitter import RecursiveCharacterTextSplitter
7  from langchain_chroma import Chroma # Updated import
8  from langchain_community.document_loaders import PyPDFLoader, WebBaseLoader
9  from langchain_core.output_parsers import StrOutputParser
10 from langchain_ollama import ChatOllama
11 from langchain_openai import OpenAIEmbeddings # Updated import
12 from openai import OpenAI
```

This section imports all the required libraries:

- `langchain` module handle prompt templates, text splitting, embeddings, and document loaders.
- `openai` library provides an interface to interact with the OpenAI API.
- `pickle` is used to serialize (save) data to disk.

Note that this code assumed the original data format was PDF. Therefore, if other types of file need to be processed, you will need to import the corresponding file loader from `langchain_community.document_loaders`.

Setting Environment Variables

```
14  client = OpenAI(
15      api_key=os.getenv("OPENAI_API_KEY"),
16  )
```

- **OpenAI client**: The `OpenAI` client initializes with your API key for interacting with OpenAI's API.

Define Paths for Serialized Data

```
18 # Define paths for serialized data
19 SPLIT_DOCS_PATH = 'split_documents_pdf.pkl' # Path to save split documents
20 VECTOR_STORE_DIR = 'chroma_vectorstore_PDF' # Directory to save Chroma data
```

These paths specify where serialized data, including split documents and vector store data, will be stored.

Step 1: Loading and Splitting Documents

Loading from PDF Files

```
29 file_paths = [
30     "documents/physics.pdf"
31 ]
32
33 def load_and_split_documents_PDF(file_paths, split_docs_path):
34     if os.path.exists(split_docs_path):
35         print("Loading existing split documents...")
36         with open(split_docs_path, 'rb') as f:
37             doc_splits = pickle.load(f)
38     else:
39         print("Loading and splitting documents...")
40         # Initialize loaders for each file and load documents
41         loaders = [PyPDFLoader(file_path) for file_path in file_paths]
42         docs = [loader.load() for loader in loaders]
43         docs_list = [item for sublist in docs for item in sublist] # Flatten the list
44
45         text_splitter = RecursiveCharacterTextSplitter.from_tiktoken_encoder(
46             chunk_size=250, chunk_overlap=30
47         )
48         doc_splits = text_splitter.split_documents(docs_list)
49
50         # Save split documents to disk
51         with open(split_docs_path, 'wb') as f:
52             pickle.dump(doc_splits, f)
53         print(f"Split documents saved to {split_docs_path}.")
54     return doc_splits
55
56 # Load or split and save documents
57 doc_splits = load_and_split_documents_PDF(file_paths, SPLIT_DOCS_PATH)
```

This function loads PDF documents, splits them into chunks, and saves them for reuse.

Step 2: Create or Load Vector Store with Embeddings

```
54 def create_or_load_vectorstore(vector_store_dir, doc_splits):
55     if os.path.exists(vector_store_dir):
56         print("Loading existing Chroma vector store...")
57         embedding = OpenAIEmbeddings(openai_api_key=os.getenv("OPENAI_API_KEY"))
58         vectorstore = Chroma(
59             persist_directory=vector_store_dir,
60             embedding_function=embedding
61         )
62     else:
63         print("Creating new Chroma vector store...")
64         embedding = OpenAIEmbeddings(openai_api_key=os.getenv("OPENAI_API_KEY"))
65         vectorstore = Chroma.from_documents(
66             documents=doc_splits,
67             embedding=embedding,
68             persist_directory=vector_store_dir
69         )
70     # Save the vector store to disk
71     print(f"Chroma vector store saved to {vector_store_dir}.")
72     return vectorstore
73
74 vectorstore = create_or_load_vectorstore(VECTOR_STORE_DIR, doc_splits)
75 retriever = vectorstore.as_retriever(search_kwargs={"k": 4}) # Retrieve top 4 relevant documents
```

This function either loads an existing vector store or creates a new one from document chunks. Embeddings are generated using OpenAIEmbeddings, and the store is saved for reuse.

Step 3: Define Prompt Templates

```
77 # Step 3: Define Prompt Template for the Language Model
78 prompt_template = PromptTemplate(
79     template="""You are an assistant for question-answering tasks.
80     Use the following documents to answer the question. Use your own knowledge when documents do not have the answer.
81
82     NEVER share the source of your response.
83     Do not include document in your response.
84
85     Question: {question}
86     Documents: {documents}
87     Answer:
88     """,
89     input_variables=["question", "documents"]
90 )
91
92 prompt_template Og = PromptTemplate(
93     template="""You are an assistant for question-answering tasks.
94
95     Question: {question}
96     Answer:
97     """,
98     input_variables=["question"]
99 )
```

Two prompt templates are defined: one for the RAG setup and one without retrieval for comparison purposes.

Step 4: Define Chat Function with OpenAI's API

```
101 # Step 4: Define a function to interact with OpenAI's ChatGPT API
102 def ask_chatgpt(question, documents, model="gpt-4o-mini"):
103     # Format the prompt with question and documents
104     prompt = prompt_template.format(question=question, documents=documents)
105
106     response = client.chat.completions.create(
107         model=model,
108         messages=[{"role": "user", "content": prompt}],
109         temperature=0
110     )
111     return response.choices[0].message.content
```

This function formats the question and document context for RAG-based queries and uses OpenAI's API to get a response.

Step 5: RAG Application Class and OG Application Class

```
113 # Step 5: Define Application Classes
114 class RAGApplication:
115     def __init__(self, retriever):
116         self.retriever = retriever
117
118     def run(self, question):
119         # Retrieve relevant documents using the 'invoke' method
120         documents = self.retriever.invoke(question)
121         # Extract content from retrieved documents
122         doc_texts = "\n".join([doc.page_content for doc in documents])
123         # Get the answer from the language model
124         answer = ask_chatgpt(question, doc_texts)
125         return answer
126
127 # Original LLM pipeline (without RAG) for comparison
128 class BasicLLMApplication:
129     def __init__(self, model="gpt-4o-mini"):
130         self.model = model
131
132     def run(self, question):
133         prompt = prompt_template_og.format(question=question)
134         response = client.chat.completions.create(
135             model=self.model,
136             messages=[{"role": "user", "content": prompt}],
137             temperature=0
138         )
139         return response.choices[0].message.content
```

The RAG class initializes with a retriever, retrieves documents based on the question, and formats them for the language model's answer generation.

The basic LLM class simply formats the prompt and ask the language model to generate answer based on the model's own knowledge

Step 6: Main Execution Loop

```
141 # Step 6: Test the Application with Interactive Loop
142 if __name__ == "__main__":
143     try:
144         # Initialize the RAG application
145         rag_application = RAGApplication(retriever)
146
147         # Initialize the original LLM without RAG
148         basic_llm_application = BasicLLMApplication()
149
150         # Welcome message and instructions
151         print("Welcome to the RAG Question-Answering System!")
152         print("Type your questions below. To exit, type 'exit' or 'quit'.\\n")
153
154         while True:
155             # Prompt user for input
156             question = input("Your question: ")
157
158             # Handle exit commands
159             if question.lower() in ['exit', 'quit']:
160                 print("Exiting the application. Goodbye!")
161                 break
162             elif question.strip() == "":
163                 print("Please enter a valid question.\\n")
164                 continue
165
166             # Measure Basic LLM response time
167             start_basic = time.perf_counter()
168             basic_answer = basic_llm_application.run(question)
169             end_basic = time.perf_counter()
170             basic_time = end_basic - start_basic
171
172             # Measure RAG response time
173             start_rag = time.perf_counter()
174             rag_answer = rag_application.run(question)
175             end_rag = time.perf_counter()
176             rag_time = end_rag - start_rag
177
178
179             # Display the answers with elapsed times
180             print("\\nAnswer with RAG:")
181             print(rag_answer)
182             print(f"Time taken: {rag_time:.2f} seconds\\n")
183
184             print("Original Model Answer:")
185             print(basic_answer)
186             print(f"Time taken: {basic_time:.2f} seconds\\n")
187
188         except KeyboardInterrupt:
189             # Handle Ctrl+C gracefully
190             print("\\nInterrupted by user. Exiting the application. Goodbye!")
191         except Exception as e:
192             print(f"An error occurred during execution: {e}")
```

This interactive loop takes user questions, retrieves responses from both RAG and basic LLM setups, and displays the answers and time taken.

Conclusion

This script provides a framework for creating an RAG system with ChatGPT, enabling better quality answers by integrating document retrieval into language model queries. The most important part of the code is Step2, where we used OpenAI's embedding API to create a vector store (embeddings) from our data chunks. Embeddings are the key for LLM to generate high quality answers in a short amount of time.