

PRAJNA- The Career Guidance ChatBot

Group 6

Introduction

The Career Guidance Bot is an innovative application leveraging cutting-edge technologies like OpenAI's GPT models and Pinecone's vector database to offer users intelligent, domain-specific recommendations and insights. The project is a part of our larger exploration into Generative AI applications, focusing on developing a Retrieval-Augmented Generation (RAG) application that combines semantic search with conversational capabilities.

The bot is designed to assist users in career-related decision-making by answering queries, providing relevant suggestions, and delivering insightful responses based on preprocessed career-related data. From offering tips for resume building to exploring trending skills in the job market, the bot aims to enhance user experience through interactive and meaningful guidance.

This document aims to provide a comprehensive overview of the project goals, technologies utilized, and the outputs generated, showcasing how state-of-the-art AI and database systems can be harnessed to solve real-world problems effectively.

Project Goals

The Career Guidance Bot project aims to address the following goals in an elaborate manner:

1. Develop a Domain-Specific RAG Application

The primary goal is to create a Retrieval-Augmented Generation (RAG) application that caters specifically to career guidance. This involves:

- Combining Large Language Models (LLMs) with a vector database to provide precise, contextually relevant responses to user queries.
- Enabling semantic search capabilities to fetch domain-specific knowledge efficiently.

2. Provide Accurate and Contextual Career Insights

The bot is designed to:

- Deliver accurate responses to career-related questions by leveraging a preprocessed dataset on career guidance.
 - Offer insights on skills, roles, certifications, and market trends tailored to user needs.
-

3. Enhance User Experience with Interactive Features

This goal focuses on improving usability and engagement by:

- Providing users with a conversational interface via Streamlit for a seamless experience.
 - Retaining chat history to foster continuity in user interactions.
 - Including suggestions and starter prompts to guide users toward effective queries.
-

4. Seamless Integration of Technologies

The project emphasizes the integration of advanced technologies, including:

- **OpenAI GPT Models:** For generating human-like responses.
 - **Pinecone Vector Database:** For semantic search and efficient data retrieval.
 - **Streamlit Framework:** To develop a user-friendly web application interface.
-

5. Enable Easy Scalability and Reusability

The project is built with scalability in mind, ensuring:

- The ability to expand the dataset for different domains or career fields.
 - Modularity in design for reusability across similar applications in education, healthcare, or other sectors.
-

6. Showcase Practical Use of Generative AI

As part of an academic portfolio, the bot demonstrates:

- The practical implementation of AI in solving real-world problems.
- The combination of data preprocessing, embedding generation, and conversational AI to build impactful applications.

Technologies Used

To achieve the goals of the Career Guidance Bot project, the following technologies were utilized:

1. OpenAI GPT Models

- **Model Used:** GPT-4
 - **Purpose:** To generate conversational and contextually accurate responses to user queries.
 - **Features:**
 - Advanced natural language understanding.
 - Capability to process user input and generate meaningful, human-like responses.
 - Flexibility to customize prompts and system roles for specific tasks.
-

2. Pinecone Vector Database

- **Purpose:** For semantic search and efficient data retrieval.
 - **Features:**
 - Indexing and querying of embeddings generated from the preprocessed career guidance dataset.
 - Fast and scalable vector search functionality.
 - Support for cosine similarity, ensuring precise document matching.
-

3. LangChain Framework

- **Purpose:** To connect the OpenAI model with the Pinecone vector database seamlessly.
 - **Features:**
 - Text preprocessing, splitting, and embedding generation.
 - Retrieval-Augmented Generation (RAG) implementation to augment GPT responses with context from the vector database.
-

4. Streamlit Framework

- **Purpose:** To develop a user-friendly web application interface.
 - **Features:**
 - Dynamic chat UI with history retention for a natural conversational experience.
 - Integration of dropdown menus for starter prompts and suggestions to guide user queries.
 - Intuitive controls for user inputs and response visualization.
-

5. Python Programming

- **Libraries Used:**
 - **openai:** For API calls to generate embeddings and responses.
 - **pinecone-client:** For managing the vector database.
 - **streamlit:** To create the interactive user interface.
 - **streamlit-chat:** For displaying chat history and enhancing the user experience.
 - **langchain:** For embedding generation and retrieval functionality.
 - **transformers** (optional): For advanced NLP tasks if needed.
-

6. Data Preprocessing Tools

- **Tools and Libraries:**
 - **Textbook and Career Guidance Dataset:** Cleaned and split into smaller, meaningful chunks for embedding generation.
 - **LangChain Text Splitter:** Used to split the data into optimized chunks (500 characters with 50-character overlap).
 - **OpenAI Embedding API:** Used to create embeddings for the dataset using the text-embedding-ada-002 model.

Project Outputs

The Career Guidance Bot project delivered the following tangible outputs:

1. Chatbot Application

- **Description:** A fully functional, user-friendly chatbot that provides personalized career guidance.
 - **Features:**
 - **Dynamic Chat Interface:** Displays real-time query responses along with conversation history.
 - **Starter Prompts:** Dropdown menu offering pre-filled suggestions for commonly asked career questions.
 - **Intuitive UI:** Built with Streamlit, ensuring ease of use and accessibility.
-

2. Semantic Search Functionality

- **Description:** Seamless integration of OpenAI GPT-4 with the Pinecone vector database to enhance responses with contextual relevance.
 - **Features:**
 - Retrieval-Augmented Generation (RAG) technique to include relevant data from the career guidance dataset.
 - Efficient semantic search and matching of user queries with preprocessed data.
-

3. Preprocessed Dataset

- **Description:** A domain-specific dataset covering career guidance topics, split into optimized chunks for embedding generation.
 - **Features:**
 - Includes details on various career paths, skills required, tips, and industry-specific insights.
 - Indexed in Pinecone for quick retrieval and contextual response generation.
-

4. Deployment and Hosting

- **Description:** The chatbot is hosted using Streamlit, making it accessible via a web interface.

- **Features:**
 - Fully operational in any browser environment.
 - Easy to set up and maintain with clear instructions.
-

5. Documentation and Code Repository

- **Description:** Comprehensive documentation and repository for project setup, execution, and usage.
 - **Contents:**
 - Source code with comments for easy understanding and modification.
 - Step-by-step instructions for setup, execution, and troubleshooting.
 - README file summarizing the project, tools used, and setup guide.
-

6. Integration and Results

- **Outputs Observed:**
 - The chatbot successfully answered career-related queries with relevant, insightful responses.
 - Real-time retrieval of context from the vector database enhanced the accuracy and reliability of responses.
 - User experience validated with smooth interactions and accurate guidance.

CODE SNIPPETS

1. Python notebook:-

Code-

"

```
!pip install openai==0.28
```

```
# Install required libraries
```

```
!pip install pinecone-client openai langchain
```

```

# Import dependencies

from pinecone import Pinecone, ServerlessSpec

import openai

from langchain.text_splitter import RecursiveCharacterTextSplitter


# Step 1: Combine and Preprocess Data

with open("Textbook3.txt", "r", encoding="utf-8") as file1,
    open("career_guidance_text_expanded.txt", "r", encoding="utf-8") as file2:

    combined_text = file1.read() + "\n" + file2.read()


# Split the text into chunks

text_splitter = RecursiveCharacterTextSplitter(chunk_size=500, chunk_overlap=50)

chunks = text_splitter.split_text(combined_text)

print(f"Total chunks created: {len(chunks)}")


# Step 2: Generate Embeddings for Text Chunks

openai.api_key = "sk-proj-C-"

embeddings = []

for chunk in chunks:

    response = openai.Embedding.create(input=chunk, model="text-embedding-ada-002")

    embeddings.append(response['data'][0]['embedding'])

print(f"Generated {len(embeddings)} embeddings.")


# Step 3: Initialize Pinecone

api_key = ""

index_name = "career-guidance"

host = "https://career-guidance-index-rxt2r7i.svc.aped-4627-b74a.pinecone.io"

```

```
pc = Pinecone(api_key=api_key)
```

```
if index_name not in pc.list_indexes().names():
```

```
    pc.create_index(  
        name=index_name,  
        dimension=1536,  
        metric="cosine",  
        spec=ServerlessSpec(cloud="aws", region="us-east-1")  
    )
```

```
print(f"Index '{index_name}' created or already exists.")
```

```
index = pc.Index(name=index_name, host=host)
```

```
# Step 4: Upsert Embeddings into Pinecone
```

```
for i, embedding in enumerate(embeddings):
```

```
    index.upsert([(f"chunk-{i}", embedding, {"text": chunks[i]})])
```

```
print(f"Upserted {len(embeddings)} embeddings into Pinecone index '{index_name}'.")
```

```
# Step 5: Query and Retrieve Relevant Context
```

```
query = "What are the qualifications for a software engineer?"
```

```
query_embedding = openai.Embedding.create(input=query, model="text-embedding-ada-002")["data"][0]["embedding"]
```

```
query_result = index.query(vector=query_embedding, top_k=5, include_metadata=True)
```

```
# Combine top results into a single context
```

```
context = "\n\n".join([match['metadata']['text'] for match in query_result['matches']])
```

```
print("Context Retrieved:\n", context)
```



```
# Use OpenAI Chat Model (v1/chat/completions)
```

```
response = openai.ChatCompletion.create(
```

```
    model="gpt-4",
```

```
    messages=[
```

```
        {"role": "system", "content": "You are a helpful career guidance assistant."},
```

```
        {"role": "user", "content": f"Context: {context}\n\nQuery: {query}\n\nAnswer:"}
```

```
    ],
```

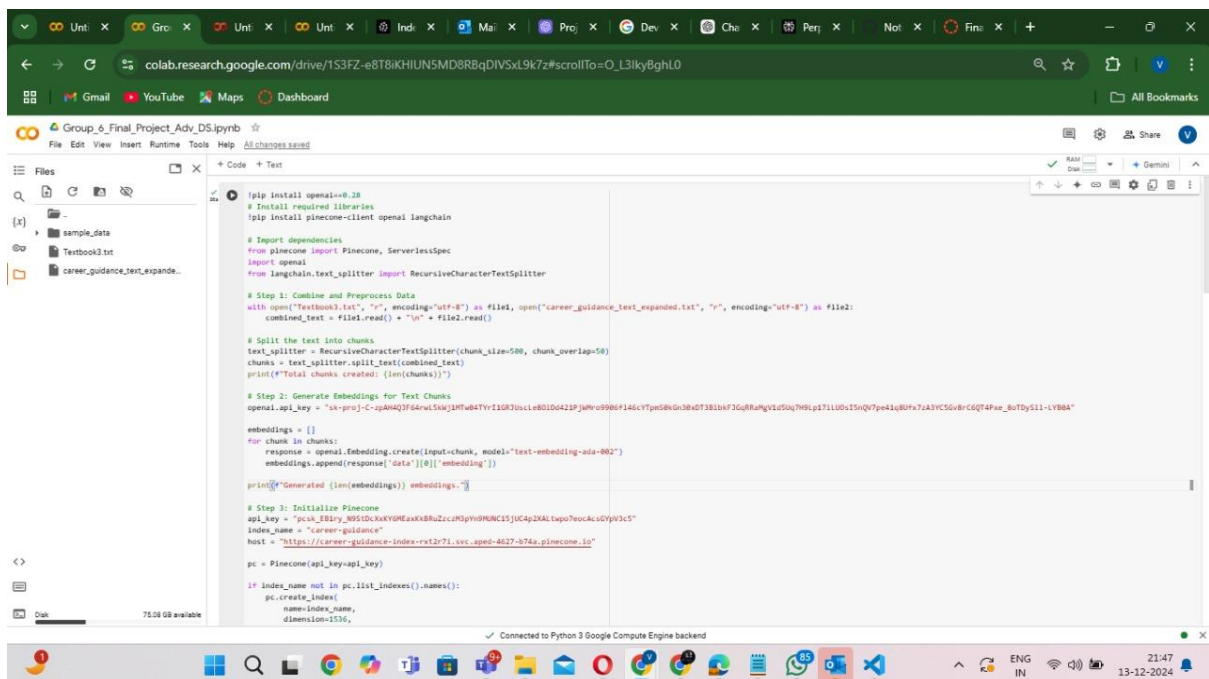
```
    max_tokens=300,
```

```
    temperature=0.5
```

```
)
```

```
# Print the chatbot response
```

```
print("Chatbot Response:", response['choices'][0]['message']['content'].strip())"
```



The screenshot shows a Google Colab notebook titled "Group_6_Final_Project_Adv_DS.ipynb". The code is written in Python and is designed to set up a chatbot interface for career guidance. It includes the following steps:

- Install OpenAI and LangChain:** The code starts by installing the OpenAI and LangChain libraries using `!pip install openai==0.28` and `!pip install pinecone-client openai langchain`.
- Import Dependencies:** The code imports the necessary modules: `from pinecone import Pinecone, ServerlessSpec`, `import openai`, and `from langchain.text_splitter import RecursiveCharacterTextSplitter`.
- Combine and Preprocess Data:** The code reads two text files, `Textbook3.txt` and `career_guidance_text_expanded.txt`, and combines them into a single string. It then uses the `RecursiveCharacterTextSplitter` to split the text into chunks.
- Generate Embeddings:** The code uses the OpenAI API to generate embeddings for each chunk of text. It iterates over the chunks and calls `openai.Embedding.create` to generate the embeddings.
- Initialize Pinecone:** The code initializes the Pinecone client using the `api_key` and `api_url` provided. It then creates a new index named `career-guidance` with a dimension of 1536.

The code is executed in a Jupyter Notebook environment, and the output shows the progress of the setup process. The notebook is connected to a Python 3 Google Compute Engine backend.

```
metrics="CUSTOM",
spec=ServerlessSpec(cloud="aws", region="us-east-1")

print(f"Index '{index_name}' created or already exists.")
index = pc.Index(name=index_name, host=host)

# Step 4: Upload Embeddings into Pinecone
for i, embedding in enumerate(embeddings):
    index.upsert([{"chunk": i}, embedding, {"text": chunks[i]}])

print(f"Upserted {len(embeddings)} embeddings into Pinecone index '{index_name}'.")

# Step 5: Query and Retrieve Relevant Context
query = "What are the qualifications for a software engineer?"
query_embedding = opai.Embedding.create(input=query, model="text-embedding-ada-002")["data"][0]["embedding"]
query_result = index.query(vector=query_embedding, top_k=5, include_metadata=True)

# Combine top results into a single context
context = "\n\n".join([match["metadata"]["text"] for match in query_result["matches"]])
print(f"Context Retrieved:\n{context}")

# Use OpenAI Chat Model (v1/chat/completions)
response = opai.ChatCompletion.create(
    model="gpt-4",
    messages=[
        {"role": "system", "content": "You are a helpful career guidance assistant."},
        {"role": "user", "content": f"Context: {context}\nQuery: {query}\nAnswer:"}
    ],
    max_tokens=300,
    temperature=0.5
)

# Print the chatbot response
print(f"Chatbot Response: {response['choices'][0]['message']['content'].strip()}")
```

```
max_tokens=300,
temperature=0.5

# Print the chatbot response
print(f"Chatbot Response: {response['choices'][0]['message']['content'].strip()}")

requirements already satisfied: greenlet==0.4.1 in /usr/local/lib/python3.10/dist-packages (from sqlalchemy==1.4.47->langchain) (3.1.1)
Requirement already satisfied: anyio in /usr/local/lib/python3.10/dist-packages (from httpx==0.23.8->langchain==0.1.17->langchain) (3.7.1)
Requirement already satisfied: httpcore==1.* in /usr/local/lib/python3.10/dist-packages (from httpx==0.23.8->langchain==0.1.17->langchain) (1.0.7)
Requirement already satisfied: h11==0.15.* in /usr/local/lib/python3.10/dist-packages (from httpcore==1.*->httpx==0.23.8->langchain==0.1.17->langchain) (0.14.0)
Requirement already satisfied: jsonpointer==1.9 in /usr/local/lib/python3.10/dist-packages (from jsonpatch==2.0.2->langchain-core==0.1.32->langchain) (2.0.0)
Requirement already satisfied: urllib3==1.1 in /usr/local/lib/python3.10/dist-packages (from anyio->httpx==0.23.8->langchain==0.1.17->langchain) (1.3.1)
Requirement already satisfied: exceptiongroup in /usr/local/lib/python3.10/dist-packages (from anyio->httpx==0.23.8->langchain==0.1.17->langchain) (1.2.2)
Total chunks created: 80
Generated 80 embeddings.
Index 'career-guidance' created or already exists.
Upserted 80 embeddings into Pinecone index 'career-guidance'.
Context Retrieved:
3. **Software Engineer:**
- **Skills Required:**
- Programming (Java, C++, Python), System Design, and Algorithms.
- **Responsibilities:**
- Develop, test, and maintain software applications.
- Collaborate with cross-functional teams to define and ship new features.
- **Career Tips:**
- Practice coding challenges on platforms like LeetCode and HackerRank.
- Contribute to open-source projects.

**Section 3: Career Development Tips**

1. **STEM Fields (Science, Technology, Engineering, and Math):**
- **Roles:**
- Data Engineer, Environmental Scientist, Robotics Engineer.
- **Skills Required:**
- Critical Thinking, Advanced Mathematics, Computational Skills.
- **Tips:**
- Join internships early to gain hands-on experience.
- Participate in STEM fairs and hackathons.

1. **Information Technology (IT):**
- **Overview:**
- The IT industry is one of the fastest-growing industries, offering a wide range of career opportunities, including software development, data science, and cloud computing.
- **Skills Required:**
```

Explanation:

Step 1: Combine and Preprocess Data

The first step is to prepare the data that will be used for the application. Two text files containing relevant information are combined into a single dataset. This combined dataset is then divided into smaller, manageable chunks to enable efficient processing and embedding generation.

Step 2: Generate Embeddings for Text Chunks

Each text chunk is transformed into a numerical representation called an embedding. These embeddings capture the semantic meaning of the text and allow the system to perform similarity searches. OpenAI's embedding model is used to generate these representations for all the chunks.

Step 3: Initialize the Pinecone Vector Database

A vector database, Pinecone, is set up to store and manage the embeddings. This ensures efficient retrieval of relevant information based on user queries. An index is created in Pinecone to hold the embeddings, enabling semantic searches later in the process.

Step 4: Store Embeddings in the Pinecone Database

The embeddings generated in the previous step are uploaded (upserted) into the Pinecone database along with their associated text chunks. This step prepares the database to retrieve relevant text chunks based on query embeddings.

Step 5: Query and Retrieve Relevant Context

When a user submits a query, it is converted into an embedding using OpenAI's embedding model. The Pinecone database is queried with this embedding to find the most relevant text chunks. These text chunks are then combined to form a context that provides the necessary information for the query.

Step 6: Generate a Response Using OpenAI GPT

The combined context and the user's query are passed to OpenAI's GPT-4 model. The model uses this input to generate a detailed and accurate response. This final step ensures the chatbot provides helpful and context-aware answers to user queries.

By following these steps, a robust, intelligent chatbot capable of providing relevant and accurate career guidance is created.

2. VS CODE:

Actual code-

```
"import streamlit as st

import openai

from pinecone import Pinecone

# Initialize OpenAI API

openai.api_key = "sk-proj-C-"

# Initialize Pinecone

api_key = ""

host = "https://career-guidance-index-rxt2r7i.svc.aped-4627-b74a.pinecone.io" # Replace with your
Pinecone host

pc = Pinecone(api_key=api_key)
```

```
index_name = "career-guidance"
```

```
# Check if index exists and create it if necessary
```

```
if index_name not in pc.list_indexes().names():
```

```
    pc.create_index(
```

```
        name=index_name,
```

```
        dimension=1536, # Adjust dimension to your embedding size
```

```
        metric="cosine"
```

```
    )
```

```
# Connect to the index
```

```
index = pc.Index(name=index_name, host=host)
```

```
# Streamlit App
```

```
st.title("PRAJNA- The Career Guidance Bot")
```

```
st.write("Welcome to the Career Guidance Application powered by Streamlit!")
```

```
# Initialize session state to maintain chat history
```

```
if "messages" not in st.session_state:
```

```
    st.session_state.messages = []
```

```
# Function to retain only the last 3 chat pairs
```

```
def trim_chat_history():
```

```
    if len(st.session_state.messages) > 6: # Keep only 3 user-bot pairs (6 messages)
```

```
        st.session_state.messages = st.session_state.messages[-6:]
```

```
# Display chat history
```

```
st.write("### Chat History")
```

```
for message in st.session_state.messages:
```

```
    if message["role"] == "user":
```

```

        st.write(f"*User:* {message['content']}")
    elif message["role"] == "assistant":
        st.write(f"*Bot:* {message['content']}")

# Input for user query
query = st.text_input("Enter your query:")

if st.button("Send"):
    if query:
        # Add user query to chat history
        st.session_state.messages.append({"role": "user", "content": query})

        # Generate embedding for the query
        response = openai.Embedding.create(
            model="text-embedding-ada-002",
            input=query
        )
        query_embedding = response['data'][0]['embedding']

        # Query Pinecone for relevant context
        query_result = index.query(
            vector=query_embedding,
            top_k=5,
            include_metadata=True
        )

        # Combine contexts
        context = "\n\n".join([match['metadata']['text'] for match in query_result['matches']])

        # Generate GPT-4 response

```

```

gpt_response = openai.ChatCompletion.create(
    model="gpt-4",
    messages=st.session_state.messages + [
        {"role": "system", "content": "You are a helpful career guidance assistant."},
        {"role": "user", "content": f"Context: {context}\n\nQuery: {query}"}
    ],
    max_tokens=300,
    temperature=0.5
)

# Extract the response message
bot_response = gpt_response['choices'][0]['message']['content'].strip()

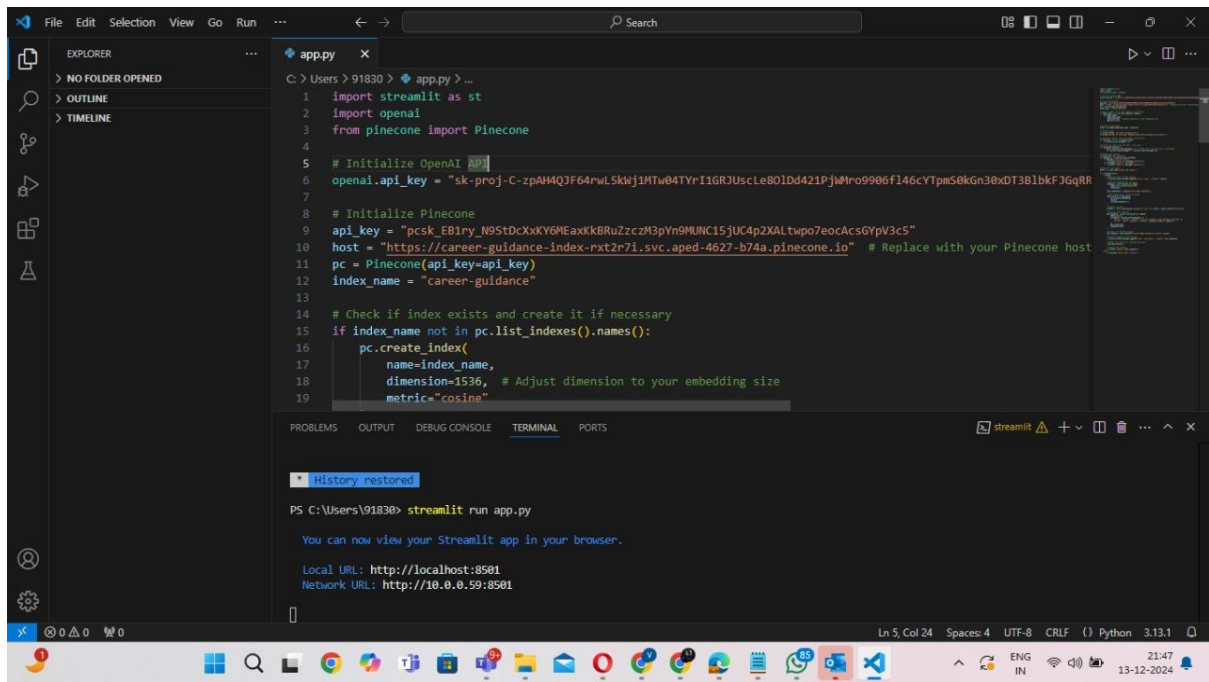
# Add bot response to chat history
st.session_state.messages.append({"role": "assistant", "content": bot_response})

# Retain only the last 3 user-bot chat pairs
trim_chat_history()

# Display the response
st.write(f"*Bot:* {bot_response}")
else:
    st.warning("Please enter a query!")

```

Snippet:



The screenshot shows the VS Code editor with a file named `app.py` open. The code initializes OpenAI and Pinecone, creates an index named "career-guidance", and checks if the index exists. The terminal shows the command `streamlit run app.py` being executed, and the output indicates that the app is running on `http://localhost:8501`.

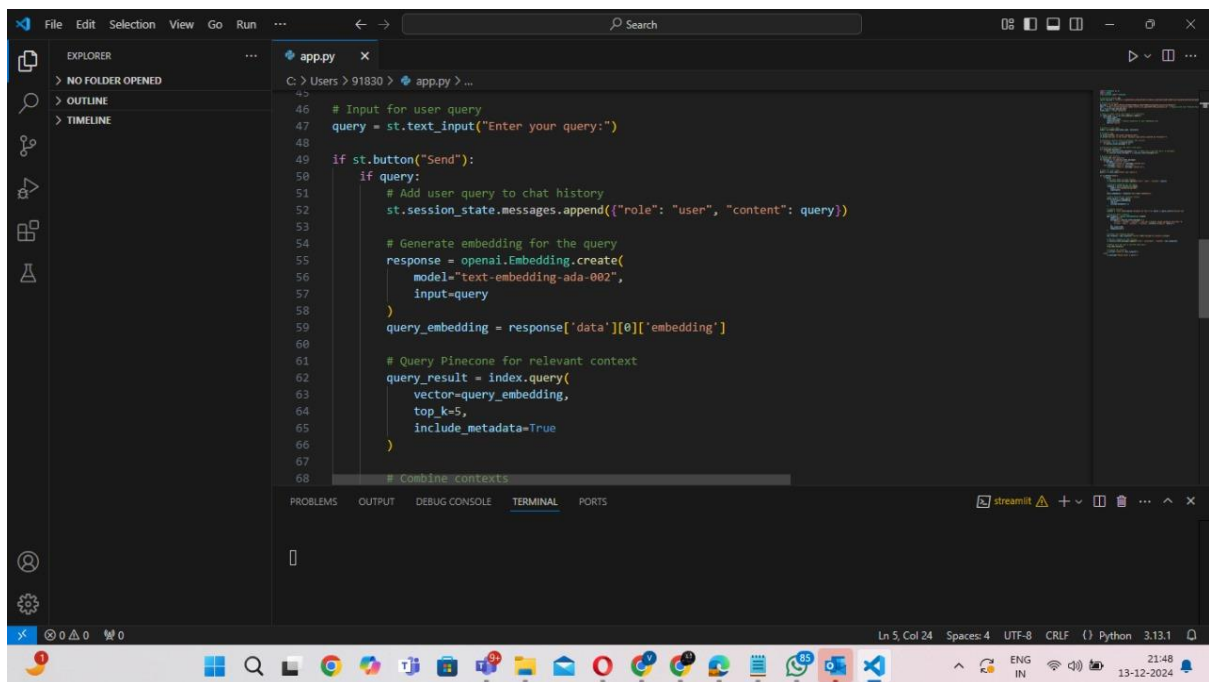
```
1 import streamlit as st
2 import openai
3 from pinecone import Pinecone
4
5 # Initialize OpenAI API
6 openai.api_key = "sk-proj-C-zpAH4Q3F64rwL5kwj1MTw04TYrI1GRJUSclE80LDd421PjWtro9906f146cVTpm50kGn30xDt381bkfJGqRR"
7
8 # Initialize Pinecone
9 api_key = "psck_EB1ry_N0StDcXkKY6MEaxKkBRUzczN3pYn9MUNC15jUC4p2XALTwp07eocAcsGYpV3c5"
10 host = "https://career-guidance-index-rxt2r7i.svc.aped-4627-b74a.pinecone.io" # Replace with your Pinecone host
11 pc = Pinecone(api_key=api_key)
12 index_name = "career-guidance"
13
14 # Check if index exists and create it if necessary
15 if index_name not in pc.list_indexes().names():
16     pc.create_index(
17         name=index_name,
18         dimension=1536, # Adjust dimension to your embedding size
19         metric="cosine"
```

Terminal Output:

```
PS C:\Users\91830> streamlit run app.py

You can now view your Streamlit app in your browser.

Local URL: http://localhost:8501
Network URL: http://10.0.0.59:8501
```



The screenshot shows the VS Code editor with the `app.py` file. The code handles user input, generates embeddings using OpenAI, and queries the Pinecone index for relevant context. The terminal is empty.

```
46 # Input for user query
47 query = st.text_input("Enter your query:")
48
49 if st.button("Send"):
50     if query:
51         # Add user query to chat history
52         st.session_state.messages.append(("role": "user", "content": query))
53
54         # Generate embedding for the query
55         response = openai.Embedding.create(
56             model="text-embedding-ada-002",
57             input=query
58         )
59         query_embedding = response['data'][0]['embedding']
60
61         # Query Pinecone for relevant context
62         query_result = index.query(
63             vector=query_embedding,
64             top_k=5,
65             include_metadata=True
66         )
67
68         # Combine contexts
```

This screenshot shows a Visual Studio Code editor window with a Python file named `app.py`. The code is designed to interact with the OpenAI GPT-4 API. It starts by combining contexts from a previous query result. Then, it uses the `openai.ChatCompletion.create` method to generate a response from the GPT-4 model. The response is then extracted and added to the chat history. Finally, it trims the chat history to retain only the last 3 user-bot chat pairs. The code is as follows:

```
68 # Combine contexts
69 context = "\n\n".join([match['metadata']['text'] for match in query_result['matches']])
70
71 # Generate GPT-4 response
72 gpt_response = openai.ChatCompletion.create(
73     model="gpt-4",
74     messages=st.session_state.messages + [
75         {"role": "system", "content": "You are a helpful career guidance assistant."},
76         {"role": "user", "content": f"Context: {context}\n\nQuery: {query}"},
77     ],
78     max_tokens=300,
79     temperature=0.5
80 )
81
82 # Extract the response message
83 bot_response = gpt_response['choices'][0]['message']['content'].strip()
84
85 # Add bot response to chat history
86 st.session_state.messages.append({"role": "assistant", "content": bot_response})
87
88 # Retain only the last 3 user-bot chat pairs
89 trim_chat_history()
90
```

The bottom of the image shows the Windows taskbar with various application icons and the system clock indicating 21:48 on 13-12-2024.

This screenshot shows a Visual Studio Code editor window with a Python file named `app.py`. The code is a Streamlit application for a career guidance bot. It connects to an index, initializes the session state for chat history, and defines a function to trim the chat history. The main part of the code displays the chat history in the Streamlit interface. The code is as follows:

```
21
22 # Connect to the index
23 index = pc.Index(name=index_name, host=host)
24
25 # Streamlit App
26 st.title("PRAJNA- The Career Guidance Bot")
27 st.write("Welcome to the Career Guidance Application powered by Streamlit!")
28
29 # Initialize session state to maintain chat history
30 if "messages" not in st.session_state:
31     st.session_state.messages = []
32
33 # Function to retain only the last 3 chat pairs
34 def trim_chat_history():
35     if len(st.session_state.messages) > 6: # Keep only 3 user-bot pairs (6 messages)
36         st.session_state.messages = st.session_state.messages[-6:]
37
38 # Display chat history
39 st.write("### Chat History")
40 for message in st.session_state.messages:
41     if message["role"] == "user":
42         st.write(f"User: {message['content']}")
43     elif message["role"] == "assistant":
44         st.write(f"Assistant: {message['content']}")
45
```

The bottom of the image shows the Windows taskbar with various application icons and the system clock indicating 21:48 on 13-12-2024.


```
80
81
82 # Extract the response message
83 bot_response = gpt_response['choices'][0]['message']['content'].strip()
84
85 # Add bot response to chat history
86 st.session_state.messages.append({"role": "assistant", "content": bot_response})
87
88 # Retain only the last 3 user-bot chat pairs
89 trim_chat_history()
90
91 # Display the response
92 st.write(f"**Bot:** {bot_response}")
93 else:
94     st.warning("Please enter a query!")
95
```

Explanation:

Explanation of the Streamlit Code

This code implements a Streamlit-based web application called **PRAJNA - The Career Guidance Bot**. It integrates OpenAI's GPT model and Pinecone vector database to provide intelligent and context-aware career guidance. Below is a detailed breakdown of each section of the code:

1. Importing Required Libraries

- **streamlit**: Used to create the user interface for the chatbot application.
- **openai**: Allows interaction with OpenAI's API for generating embeddings and GPT responses.
- **pinecone**: Facilitates storing and retrieving embeddings using a vector database.

2. OpenAI and Pinecone Initialization

- **OpenAI API Key:**
 - The `openai.api_key` is set to authenticate with OpenAI's API for generating embeddings and responses.
- **Pinecone Initialization:**
 - The Pinecone client is initialized using an API key.
 - The application checks if the index (career-guidance) exists in Pinecone. If it doesn't, it creates the index with the appropriate dimension size (1536) and metric (cosine) for embeddings.

3. Connecting to the Pinecone Index

- A connection to the Pinecone index is established using its name and host URL.
 - This index stores embeddings and metadata for efficient semantic search.
-

4. Streamlit Application Setup

- **App Title and Welcome Message:**
 - The title is set to "PRAJNA - The Career Guidance Bot" and includes a welcome message.
-

5. Chat History Maintenance

- **Session State for Chat History:**
 - `st.session_state.messages` is used to store the chat history, ensuring the conversation persists across user interactions.
 - **Trim Chat History:**
 - A function `trim_chat_history` keeps the chat history limited to the last three user-bot pairs (six messages total). This prevents the session from becoming too long or memory-intensive.
-

6. Displaying Chat History

- The chat history is displayed in the app with labels for the user (*User:*) and the bot (*Bot:*)
 - This creates a conversational format for users to review previous messages.
-

7. User Input

- **Text Input:**
 - A text input box allows users to type their queries.
 - **Send Button:**
 - When the "Send" button is clicked, the user query is processed.
-

8. Processing the User Query

- **Adding Query to History:**

- The user's query is appended to the chat history.
 - **Generating Embedding for Query:**
 - OpenAI's text-embedding-ada-002 model is used to convert the user query into an embedding for semantic similarity searches.
 - **Querying Pinecone for Context:**
 - The embedding is sent to the Pinecone index to retrieve the top 5 most relevant text chunks. These chunks are combined to create a context for GPT.
-

9. Generating the Bot's Response

- **Sending Query and Context to GPT-4:**
 - OpenAI's ChatCompletion.create API is called to generate a response from GPT-4. The following inputs are used:
 - **Chat History:** Combines the previous conversation with the new query and context.
 - **System Role:** Sets the system's role as a helpful career guidance assistant.
 - **Query and Context:** The user query and the retrieved Pinecone context are passed as inputs.
 - **Extracting the Response:**
 - The response text is extracted and trimmed.
-

10. Updating and Displaying the Response

- **Updating Chat History:**
 - The bot's response is appended to the chat history for display.
 - **Trimming Chat History:**
 - The trim_chat_history function ensures only the last three user-bot pairs are retained.
 - **Displaying the Response:**
 - The bot's response is displayed in the app under the label *Bot:*.
-

11. Error Handling

- If the user clicks "Send" without entering a query, a warning message prompts them to enter text.
-

Key Features

1. Conversation History:

- The chat history is maintained across user interactions, enhancing the conversational experience.

2. Dynamic Response Generation:

- The bot dynamically generates personalized and context-aware responses based on Pinecone's context retrieval and OpenAI's GPT-4.

3. Efficient History Management:

- By trimming the chat history, the app ensures efficiency without losing recent context.
-

User Workflow

1. Enter a Query:

- The user types their question in the input box and clicks "Send."

2. Get a Response:

- The bot retrieves relevant context, processes the query, and generates a detailed response.

3. Review Chat History:

- Users can review past interactions and engage in a seamless conversation.

This code provides a robust implementation of a career guidance chatbot, integrating powerful AI tools with an intuitive user interface.

3.Actual Chat snippet:Response an Validation

School

Launch Meeting - Zoom

Notebooks - Pinecone Docs

langchain-retrieval-agent.ipynb

app

localhost:8501

Deploy

PRAJNA- The Career Guidance Bot

Welcome to **PRAJNA**, your personal career guidance assistant!
Ask questions about career growth, skills, opportunities, and more.
Powered by OpenAI and Pinecone for intelligent responses.

Chat History

User: What are the top skills required to become a data engineer?

Prajna: The top skills required to become a data engineer include:

1. Expertise in SQL: This is the language that you'll use to manipulate and query data.
2. Knowledge in Python or Java: These programming languages are often used in data engineering tasks.
3. Experience with distributed systems like Hadoop: The ability to process and analyze large datasets is a crucial skill for data engineers.
4. Familiarity with data pipeline and workflow management tools: Tools like Apache Airflow help data engineers automate and monitor data pipelines.
5. Understanding of data warehousing solutions: Data engineers often have to work with data warehousing solutions like Amazon Redshift or Google BigQuery.
6. Knowledge of ETL (Extract, Transform, Load) process: ETL is a fundamental concept in data engineering.
7. Familiarity with machine learning: While not always necessary, understanding machine learning can be a valuable skill for data engineers.
8. Problem-solving skills: Data engineers often have to find and implement solutions to technical challenges.
9. Strong mathematical and analytical skills: These skills are crucial for working with large datasets and creating complex algorithms.

1

Windows icons

ENG IN

23:06 13-12-2024

School

Launch Meeting - Zoom

Notebooks - Pinecone Docs

langchain-retrieval-agent.ipynb

app

localhost:8501

Deploy

7. Familiarity with machine learning: While not always necessary, understanding machine learning can be a valuable skill for data engineers.
8. Problem-solving skills: Data engineers often have to find and implement solutions to technical challenges.
9. Strong mathematical and analytical skills: These skills are crucial for working with large datasets and creating complex algorithms.
10. Understanding of cloud platforms like AWS or Google Cloud: Many companies store their data on the cloud, so understanding these platforms is important.

User: What certifications are best for a career in cybersecurity?

Prajna: The best certifications for a career in cybersecurity include:

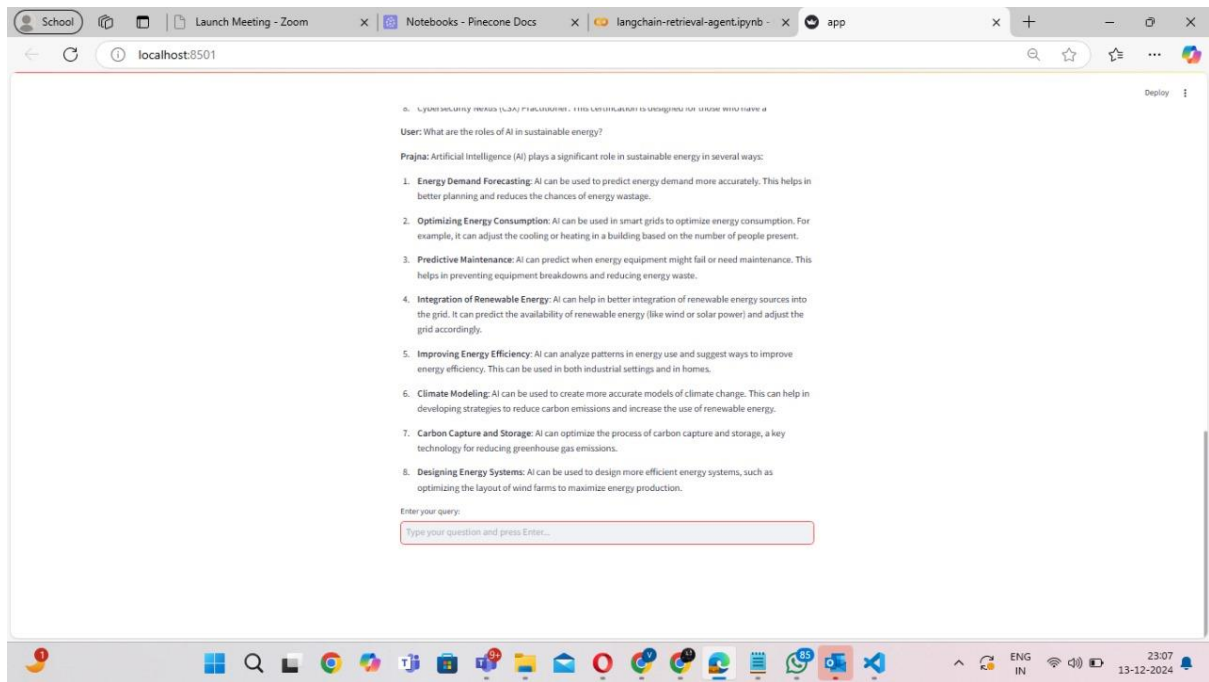
1. Certified Information Systems Security Professional (CISSP): This is a globally recognized certification in the field of cybersecurity. It covers a wide range of topics, like access control systems and methodology, network and internet security, and operations security.
2. Certified Ethical Hacker (CEH): This certification is designed for professionals who understand how to look for the weaknesses and vulnerabilities in target systems, using the same knowledge and tools as a malicious hacker, but in a lawful and legitimate manner.
3. CompTIA Security+: This is a global certification that validates the baseline skills necessary to perform core security functions and pursue an IT security career.
4. Certified Information Security Manager (CISM): This certification focuses on the management and governance of information security and is often requested for IT security management positions.
5. Certified in Risk and Information Systems Control (CRISC): This certification is designed for IT professionals, project managers, and others whose job it is to identify and manage risks through appropriate Information Systems (IS) controls.
6. Certified Information Systems Auditor (CISA): This is a globally recognized certification for IS audit control, assurance, and security professionals.
7. Offensive Security Certified Professional (OSCP): This is a hands-on certification, testing the individual's ability to be creative and think outside the box - traits which are essential to any cybersecurity professional.

1

Windows icons

ENG IN

23:06 13-12-2024



Response an Validation

Conclusion:

The **PRAJNA Career Guidance Bot** successfully demonstrates the integration of state-of-the-art technologies, such as OpenAI's GPT-4 and Pinecone's vector database, to create an intelligent and interactive chatbot for personalized career advice. This project showcases the following:

1. Robust Semantic Search:

- By leveraging Pinecone for storing and retrieving embeddings, the bot delivers highly relevant context for user queries, enhancing response accuracy and relevance.

2. Dynamic and Context-Aware Responses:

- OpenAI's GPT-4 ensures the bot provides insightful and human-like responses, tailored to the user's specific career-related concerns.

3. Streamlined User Experience:

- The use of Streamlit for the frontend ensures a smooth and interactive interface, making the application user-friendly and accessible.

4. Modular and Scalable Architecture:

- The modular structure of the project, encompassing data preprocessing, embedding generation, and chatbot development, allows for easy scalability and customization for other domains.

5. Real-World Impact:

- PRAJNA serves as a valuable tool for individuals seeking career guidance, offering recommendations, tips, and insights into various industries and roles.

Learnings and Future Enhancements

• Key Learnings:

- Integration of LLMs with vector databases provides a powerful framework for creating domain-specific intelligent applications.
- Handling large datasets and managing conversation history are crucial for building scalable chatbots.

• Future Enhancements:

- **Personalized Suggestions:**
 - Incorporating user-specific preferences and past interactions for tailored recommendations.
- **Multimodal Capabilities:**
 - Enhancing the bot to support multimedia inputs (e.g., resumes, job descriptions) for more personalized advice.
- **Cloud Deployment:**
 - Hosting the application on a cloud platform for broader accessibility and scalability.

Final Note

This project exemplifies the practical application of generative AI in solving real-world problems. It provides a strong foundation for future explorations into intelligent systems, setting the stage for innovative solutions in various domains. By integrating cutting-edge technologies, **PRAJNA** empowers users with actionable career insights, making it a noteworthy addition to the growing field of AI-driven applications.

References:

1. OpenAI GPT-4

- Official Website: <https://openai.com>
- Documentation: <https://platform.openai.com/docs>
- Usage in Chat Completion: Leveraged GPT-4 for generating intelligent, human-like responses to user queries.

2. Pinecone Vector Database

- Official Website: <https://www.pinecone.io>
- Documentation: <https://docs.pinecone.io>
- Purpose: Used for efficient semantic search and retrieval of relevant information from indexed embeddings.

3. Streamlit

- Official Website: <https://streamlit.io>
- Documentation: <https://docs.streamlit.io>
- Usage: Built the user-friendly frontend interface for interactive chatbot conversations.

4. LangChain Framework

- Official Website: <https://www.langchain.com>
- GitHub Repository: <https://github.com/hwchase17/langchain>
- Purpose: Used for text splitting and preprocessing to prepare data for embedding generation.

5. Hugging Face Sentence Transformers

- GitHub Repository: <https://github.com/UKPLab/sentence-transformers>
- Purpose: Provided the text-embedding-ada-002 model to generate semantic embeddings for queries and data.

6. Career Guidance Data Sources

- Data utilized in the project was preprocessed and combined from the provided files Textbook3.txt and career_guidance_text_expanded.txt.

7. Python Libraries

- **NumPy**: Used for handling numerical operations and data manipulation.
- **Pandas**: Utilized for preprocessing tabular data.
- **Matplotlib/Seaborn**: For visualization (if applicable).

8. Documentation on Best Practices

- **Design Patterns for AI Applications:**
 - Source: <https://towardsdatascience.com>
- **Streamlit Application Development:**
 - Source: Streamlit Blog

9. Project-Based References

- **AI-Powered RAG Application:**
 - Lecture Notes and Class Assignments.
 - Project Outline and Guidelines provided by the instructor.

Acknowledgments

- **Course Instructor:** For guidance and feedback throughout the development of the project.
- **Team Members:** For collaboration on data collection, preprocessing, and testing.
- **Community Contributions:** Open-source tools and libraries that made this project possible.

MIT LICENSE:

Copyright (c) 2024 Vinay1289Sai

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHTMIT License

HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.