

# Python Code Explanation and Documentation

Your Name

November 4, 2024

## 1 Introduction

This document provides detailed explanations for each line of the Python code, focusing on the functionality and purpose of every component. The Python code implements a Streamlit-based multi-PDF content-based question-answering system using various libraries such as PyPDF2, Langchain, and Google Generative AI (Gemini). The code also includes a user login system and text-based search using embeddings and TF-IDF techniques.

## 2 Code Breakdown

### 2.1 Imports

The code begins by importing necessary libraries and modules.

```
1 import os
2 import shutil
3 import streamlit as st
4 from PyPDF2 import PdfReader
5 from langchain.text_splitter import RecursiveCharacterTextSplitter
6 import google.generativeai as genai
7 from langchain_community.vectorstores import FAISS
8 from langchain_google_genai import GoogleGenerativeAIEmbeddings
9 from langchain_google_genai import ChatGoogleGenerativeAI
10 from langchain.chains.question_answering import load_qa_chain
11 from langchain.prompts import PromptTemplate
12 from dotenv import load_dotenv
13 import pickle
14 import re
15 import numpy as np
16 from sklearn.feature_extraction.text import TfidfVectorizer
17 from sklearn.metrics.pairwise import cosine_similarity
18 import hashlib
```

- `os`, `shutil`: These modules are used for file and directory manipulation.
- `streamlit`: The library is used to create the web interface.
- `PyPDF2`: Provides functionalities to read and extract text from PDF files.

- `RecursiveCharacterTextSplitter`: Splits text into chunks for embedding and search.
- `google.generativeai`: Configures and interacts with Google Generative AI (Gemini model).
- `FAISS`: A vector store for fast similarity search.
- `load_qa_chain`: Loads a question-answering chain using a model and a prompt.
- `dotenv`: Loads environment variables from a `.env` file.
- `hashlib`: Provides secure hashing for password encryption.
- `TfidfVectorizer`, `cosine_similarity`: For text vectorization and calculating similarity scores.

## 2.2 Loading Environment Variables

```
1 # Load environment variables
2 load_dotenv()
3 api_key = os.getenv("GOOGLE_API_KEY")
```

`load_dotenv()` loads environment variables from a `.env` file, and `os.getenv()` retrieves the Google API key required for the Generative AI service.

## 2.3 Google Generative AI Configuration

```
1 # Configure the Google Generative AI
2 genai.configure(api_key=api_key)
```

This line configures the Google Generative AI (Gemini) by using the API key obtained earlier.

## 2.4 Hashing Passwords

```
1 # Function to hash passwords
2 def hash_password(password):
3     return hashlib.sha256(password.encode()).hexdigest()
```

The function `hash_password` takes a plain text password, hashes it using the SHA-256 algorithm, and returns the hashed password for secure storage.

## 2.5 User Management (Login/Logout)

```
1 # Dictionary for storing users (in-memory for demo purposes)
2 users = {
3     "user1": hash_password("password1"),
4     "user2": hash_password("password2"),
5 }
```

This dictionary stores users and their corresponding hashed passwords. In production, a more secure database system would be used.

```
1 # User login system
2 def login():
3     st.sidebar.title("Login")
4     username = st.sidebar.text_input("Username")
5     password = st.sidebar.text_input("Password", type="password")
6     if st.sidebar.button("Login"):
7         hashed_password = hash_password(password)
8         if username in users and users[username] == hashed_password
9             :
10             st.session_state['logged_in'] = True
11             st.session_state['username'] = username
12             st.sidebar.success(f"Welcome {username}!")
13         else:
14             st.sidebar.error("Incorrect username or password")
```

The `login()` function creates a login form in the Streamlit sidebar. The entered password is hashed and compared with the stored value. If the credentials are valid, the login session is established.

```
1 # Logout system
2 def logout():
3     if 'logged_in' in st.session_state and st.session_state['
4         logged_in']:
5         st.sidebar.write(f"Logged in as {st.session_state['username
6         ']}")
7         if st.sidebar.button("Logout"):
8             st.session_state['logged_in'] = False
9             st.session_state['username'] = None
```

The `logout()` function allows users to log out and resets their session state.

## 2.6 PDF Text Extraction

```
1 # Function to extract text from multiple PDFs with document names
   and page numbers
2 def get_pdf_text_with_pages(pdf_paths):
3     text_chunks_with_pages = []
4     for pdf_path in pdf_paths:
5         try:
6             pdf_reader = PdfReader(pdf_path)
7             doc_name = os.path.basename(pdf_path)
8             for page_number, page in enumerate(pdf_reader.pages,
9                 start=1):
10                 text = page.extract_text()
11                 if text:
12                     text_chunks_with_pages.append((text,
13                         page_number, doc_name))
14         except FileNotFoundError as e:
15             print(f"Error: {e}")
16     return text_chunks_with_pages
```

This function reads the contents of multiple PDFs using PyPDF2. For each page, the extracted text is stored along with the page number and document name in a list of tuples.

## 2.7 Text Chunking

```
1 # Split text into chunks, including document name and page numbers
2 def get_text_chunks_with_pages(text_chunks_with_pages):
3     text_splitter = RecursiveCharacterTextSplitter(chunk_size=400,
4     chunk_overlap=50)
5     chunks_with_pages = []
6     for text, page_number, doc_name in text_chunks_with_pages:
7         chunks = text_splitter.split_text(text)
8         for chunk in chunks:
9             chunks_with_pages.append((chunk, page_number, doc_name))
10    return chunks_with_pages
```

The `get_text_chunks_with_pages()` function breaks long text into smaller chunks using `RecursiveCharacterTextSplitter`. The chunks maintain the page number and document name information.

## 2.8 Vector Store Creation

```
1 # Create and save vector store with document names and page numbers
2 def get_vector_store_with_pages(text_chunks_with_pages, index_name):
3     :
4     embeddings = GoogleGenerativeAIEmbeddings(model="models/text-
5     embedding-004")
6     texts, page_numbers, doc_names = zip(*text_chunks_with_pages)
7     vector_store = FAISS.from_texts(texts, embedding=embeddings)
8     vector_store.save_local(index_name)
9     with open("page_numbers_docs.pkl", "wb") as f:
10        pickle.dump((page_numbers, doc_names), f)
```

This function creates a FAISS vector store using text embeddings from Google Generative AI. The vector store is saved locally, and the associated page numbers and document names are also stored in a `.pkl` file for later use.

## 2.9 Question Answering Chain

```
1 # Function to get the conversational chain
2 def get_conversational_chain():
3     prompt_template = """
4     Answer the question in a detailed manner using the given
5     context.
6     Question: {question}
7     Context: {context}
8     """
9     prompt = PromptTemplate(template=prompt_template,
10    input_variables=["question", "context"])
11    llm = ChatGoogleGenerativeAI(model="models/chat-bison-001")
12    return load_qa_chain(llm=llm, prompt=prompt)
```

The `get_conversational_chain()` function creates a question-answering chain by loading a conversational prompt template and configuring it with the Google Generative AI large language model.

### 3 Conclusion

This LaTeX documentation details every aspect of the Python code, focusing on its structure and functionality. The code implements a content-based question-answering system over multiple PDF documents using embeddings and vector search techniques. Each section has been explained to provide a clear understanding of how the system works.