

Mr. Help Mate AI

(A RAG Implementation)

Author: Kunal Sahu

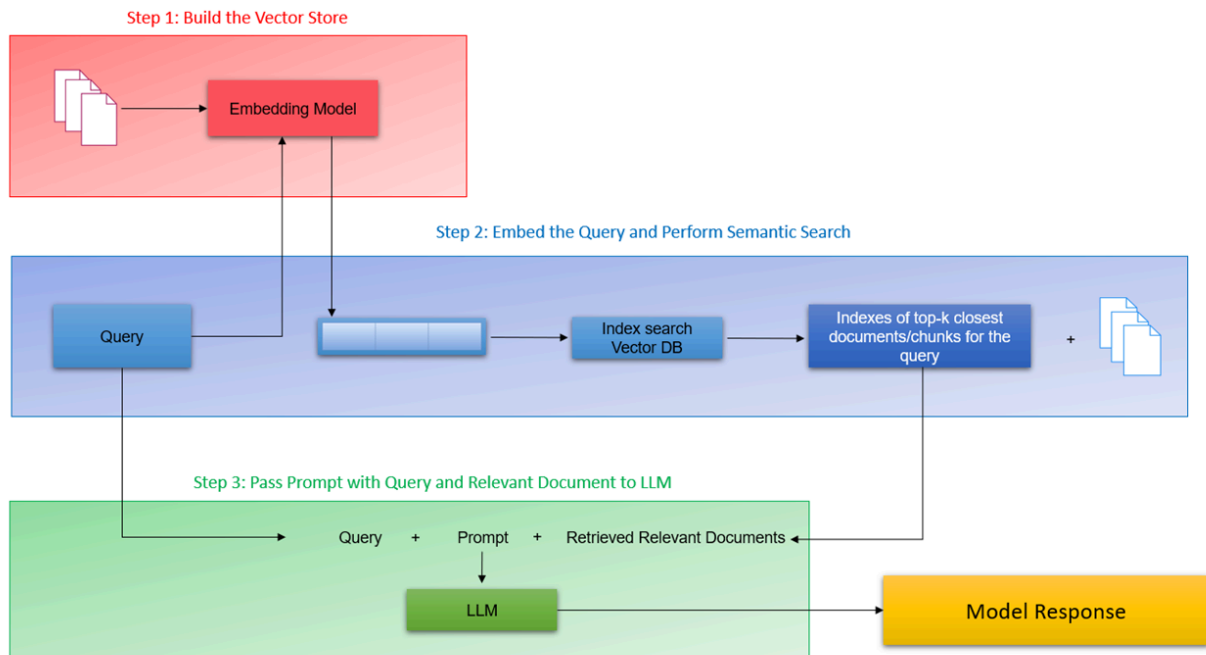
Date: 25-06-2024

Project Objective:

Build a project in the insurance domain, similar the project you saw in the "retrieval augmented generation" session. The goal of the project will be to build a robust generative search system capable of effectively and accurately answering questions from a policy document.

You will be using a single long life insurance policy document for this project.
The Goal of this project to implement a RAG project for insurance corpus.

HLD:



Project Descriptions:

Embedding Model:

We are embedding the pages in the dataframe through OpenAI's text-embedding-ada-002 model, and store them in a ChromaDB collection. That will be used for semantic search.

Caching :

We are also using Caching to store the Retrieved documents from ChromaDB.

Query :

User will send the query to the RAG System.

Top K- documents:

First user query is sent to Cache. If the results are found in Cache, it will return the results. If not, it will go to the embedding vector and find the semantically relevant top k documents.

LLM:

LLM module is nothing but ChatGPT. Once we receive the top k documents after semantic search, we pass the top k documents as a dataframe and query to ChatGPT to answer the correct results for the user query.

Cross-encoder:

We have used a cross-encoder to re-rank the semantic search results.

Attachments:

Semantic Search :

We have attached three screenshots for the top 3 results for semantic search for three custom queries.

Generative Search :

We have attached three screenshots for three custom query generative search results.

Code:

Jupyter Notebook developed for Mr Help Mate AI