Mohammed Adil
https://huggingface.co/spaces/aadil732/Analytics-Vidya-Free-Courses

# Analytics Vidya Free Course Search Engine

## 1. Introduction

This report outlines the development of a free courses search engine using Retrieval-Augmented Generation (RAG). The project aims to address the challenge of retrieving relevant course information from Analytics Vidya website with Web Scraping while providing an efficient and user-friendly experience.

**Objectives**

- Develop a semantic search engine leveraging RAG.
- Ensure fast, accurate, and contextually relevant course recommendations based on keywords and natural language.

## 2. Approach and Solution

### 2.1 Problem Understanding

Traditional keyword-based search methods often fail to capture the semantic nuances of user queries. This project addresses these limitations by implementing RAG to enhance search capabilities and deliver more relevant results.

### 2.2 Proposed Solution

The solution integrates a retrieval mechanism with an LLM-based generation process, improving accuracy and providing detailed, context-aware responses.

## 3. Architecture and Workflow

### 3.1 System Architecture

Key Components:

https://www.linkedin.com/in/mohammed-adil-silawat/

Mohammed Adil
https://huggingface.co/spaces/aadil732/Analytics-Vidya-Free-Courses

- **Embedding Model:** Generates dense vector representations of queries and course data.
- **Vector Database:** Efficiently stores and retrieves course embeddings.
- **Large Language Model (LLM):** Processes retrieved data to generate coherent responses.
- **Web UI:** Ensures efficient communication between user and backend.

### 3.2 Workflow

1. System work:
   a. uploading_to_database.py ensures the fetching, filtering, embedding and storing the course data to **Pinecone Vector Database.**
   b. This file utilizes extract_all_courses.py and extract_single_course.py files to extract and filter the course details.
2. User Interaction:
   a. app.py is a Streamlit python file for Web UI.
   b. retriever.py is a FastAPI endpoint file that accepts the user keyword, fetches the similar keywords from the vector database, and provides the results to the LLM model.
   c. The LLM model make the results to be well-structured and easily readable in markdown.

# 4. Methodology

### 4.1 Embedding Model Selection

**Criteria:**

- Accuracy in representing semantic meaning.
- Speed and efficiency.
- Compatibility with vector database.

**Selected Model:** Google's `text-embedding-004` embedding model, chosen for its high performance and reliability.

### 4.2 Vector Database

**Criteria:**

- Accurate similarity search on remote database.

https://www.linkedin.com/in/mohammed-adil-silawat/

Mohammed Adil

- Speed and efficiency.
- Free to use upto limit.

**Selected Database:** Pinecone vector database, for its scalable and efficient nearest-neighbor search capabilities. It is free to use upto 5 indexes and a high number of requests supported.

### 4.3 LLM Selection

**Considerations:**

- Fluency and coherence in generating responses.
- Free to use upto 15 requests per minute.

**Selected Model:** Google Gemini 1.5 Flash's API, chosen for their superior language understanding and response generation.

# 5. Implementation Details

### 5.1 Tools and Technologies

- **Languages & Frameworks:** Python, Streamlit, FastAPI.
- **APIs & Libraries:** LangChain, Pinecone, Google generative ai.
- **Deployment:** Deployed using Docker on Hugging Face Space.

### 5.2 Code Structure

**Modules:**

- **Frontend:** Built with Streamlit for interactive user interfaces.
- **Backend:** FastAPI handles query processing and integration.
- **Integration:** Seamless communication between vector database and LLM.

### 5.3 Challenges and Solutions

**Challenges:**

- Managing latency for queries.
- Ensuring the integrity of course details.
- Ensuring relevance of retrieved courses.

**Solutions:**

Mohammed Adil
https://huggingface.co/spaces/aadil732/Analytics-Vidya-Free-Courses

- Google's LLM and embedding model for fast response.
- Keeping each course in a single vector for data integrity.
- Using Pinecone's Cosine similarity search for most relevant courses search.

# 7. Conclusion

The project successfully delivers an advanced search engine that bridges the gap between user intent and course availability. By leveraging RAG, it ensures contextually relevant and accurate recommendations, greatly enhancing user experience.

# 8. References

- LangChain documentation.
- Pinecone API guides.
- Langchain Gemini API documentation.
- Streamlit documentation.

---

**Appendix**

- **UI Screenshots:**



https://www.linkedin.com/in/mohammed-adil-silawat/

Mohammed Adil
https://huggingface.co/spaces/aadil732/Analytics-Vidya-Free-Courses





https://www.linkedin.com/in/mohammed-adil-silawat/