

MLP project

what we have now:

- extracted data from Arxiv merged into single master sheet
- developed API in Flask framework - mostly over socket
- dataset encoding:
 - kept two choices, each consisting of 3 stations (use 3 models, one model with no use of embeddings, one with use of average embedding)
 - each choice is generated separately using BERT from Hugging Face (mean pooling to create vector of each choice)
 - embeddings are stored in DB
 - can read up previous user database
 - BERT has to learn very few first 100 entries for now
- Under Database - Models:
 - living in Data container
- relevant document retrieval:
 - search of relevant choice by examining the query using cosine similarity
 - top 5 choices are returned

What we said is future work:

- try to success in experimenting with choices (?)
- hybrid search → for keywords and semantic search (eg. "what's new")
- LLM could improve the performance of information retrieval (11) as we can e.g. choose best 4 with it and ask LLM to choose best from amongst them
- Create Test Data should have 20-100 examples (can be generated with ChatGPT)
- use more data than only this being part
- choose LLM
- add prompt generation to use LLM
- integrate with UI / evaluation?

Improvements:

- automate data acquisition
- chunking - slow \rightarrow
- where do?
- LangChain

Assignment 4

Hugging Face - source of the Hugging Face model

Procedure - vector database

How does the pipeline look like?

1. Download data to all three files (model, config, tokenizer) on the webpage of the dataset (published on huggingface)

2. Readable Documentation → function from huggingface (load config, loading the documents, fine-tuning)

3. Chunking:

- function \rightarrow maximum token limit is 4096
- individual token is cut into short size
- use document splitters
- two possibilities with splitter:
 - what chunking's values
 - and why

4. Applying chunking to all documents and save them in input config - weight ID to each chunk to know from which document they are

5. dictionary with chunk id per document

differences:

- what is change
- to implement
- in the project

6. transform chunks into vector embeddings using Hugging Face pipeline → will be used for building chain between query & chunks or when we want to do some kind of reasoning

→ try running sentence-transformer (as Hugging Face is used for vector embeddings)

7. Pipeline + step environment (including metrics)

8. Pipeline - loading data into the database (personal knowledge)

9. now we need a way to generate the answer from the retrieved chunks

- Text - generator pipeline from Hugging Face
- Embed - organization using BERT embeddings using load_weighted configuration from transformers library
- initialize a test-generator pipeline with HF that uses the LLaMA model to generate some text, given some input. We will then use this pipeline inside LangChain to build our Question-Answering system
- > note: we must initialize a language model and tokenizer for the language model

Question Answering Chain

- for PQA in LangChain we need to initialize RetrievalQA or RetrievalGPTSourceChain
- > to initialize them, we need a LLM and a Passage index
- > when we want to only query, we get some answer but not well connected to our dataset, therefore we need to use our rag pipeline and only then LLaMA2 to answer our questions

Conversational Retrieval Chain

- > aim: extend the already created retrieval chain to be able to remember the previous questions and answer the current questions by looking at the previous context
- > use conversational memory to create a conversational memory >>
- > it takes chat history (list of messages) and new question and returns an answer to the new one
- > detailed description how it works - in assignment file

Task 2: Advanced RAG Techniques and Evaluation

- > maybe not, they focus on already annotated dataset - or maybe also can have similar???
- > this time they use another vector database - Chroma

Subtask 2.1. Build Contextual Compression in LangChain

- split documents using TextCutter from past - embed with embedding model from previous task
- use simple RAG pipeline that works on top of Chroma vector
- > apply rag simple pipeline to all questions in corpus and accumulate the answers
- use ContextualCompressionRehearsal to extract only the relevant information from documents
- for the final answer we can use RagRetriever
- > create context compressor RAG pipeline

Subtask 2.2. Evaluate

- they explore BLEU, ROUGE, BERTScore metrics
- HF has implementation for all those metrics
- If the answers contain multiple lines, merge them to single string
- > in the assignment there are more variants of those metrics described

Znaleziska:

- helped you a lot to understand the functioning of Chroma, especially when working with embeddings
- > when it's difficult to understand, just trying to use some words, like 'embed'
- > when it's difficult to understand, just trying to use some words, like 'embed'