


Introduction to LLM

Practice Session 12

Dense Retrieval, QA, and ANN

High-level Dense Retriever Pipeline




- Training data is prepared as a triplet $(q, D+, D-)$
 - The goal is to maximize the similarity score between $(q, D+)$ and minimize that between $(q, D-)$
 - Hard training means that $D-$ shouldn't be random, they must be confusing (similar but not relevant)

First Stage Retrievers


Query

How to make a good cappuccino




- First stage retrievers can be:-
 - Dense retrievers, e.g., BERT-DOT
 - Sparse retrievers, e.g., BM25
 - Hybrid retrievers (combination of both)

Dense Retriever Architecture (Bi-encoder)



- Bi-encoders is an architecture paradigm in which there two networks:-
 - One network to embed the query
 - Another network to embed the documents
- BERT-DOT model is a type of Bi-encoder whose networks are based on BERT


Span-QA (Reading Comprehension)



$$p_{\text{start}}(i) = \text{softmax}(\mathbf{W}_s \mathbf{h}_i)$$

$$p_{\text{end}}(i) = \text{softmax}(\mathbf{W}_e \mathbf{h}_i)$$


- This simplified version of QA aka **Reading Comprehension**.
 - (Passage, Question) \Rightarrow Answer





Question: How many parameters does BERT-large have?

Reference Text: BERT-large is really big... it has 24 layers and an embedding size of 1,024, for a total of 340M parameters! Altogether it is 1.34GB, so expect it to take a couple minutes to download to your Colab instance.

- Two heads (simple linear projection layers) are trained on top of BERT contextualized token embeddings
 - One head is trained to maximize probability for the token that should be the start position of the answer span
 - Another head is trained to maximize probability for the token that should be the end position of the answer span
- You enforce constraints, e.g., end position \geq start position, or start & end position must be in paragraph segment




Indexing: Approximate Nearest Neighbor (ANN) Search




- Flat (exact) indexing is comparing query vector with all documents vectors
- IVF (inverted file) indexing is comparing query vector with cluster nodes
 - Then compare with documents vectors existing only on top-k clusters

ANN: Product Quantization (PQ)





- Source: [Article Link](#)
- Two main hyperparameters:-
 - m (number of subvectors) and the original vector dimension should be divisible by it
 - $n\text{bits}$ ($2^{n\text{bits}}$ is total number of centroids) which is the size of the PQ code (representation of each subvector)

ANN: Product Quantization (PQ)




- Source: [Article Link](#)
- The distance between a query and a document is the sum of distances between their subvectors

ANN: Hierarchical Navigable Small Worlds (HNSW)




- Source: [Article Link](#)
- Two main hyperparameters:-
 - m (number of neighbors for each vector) which is done offline during indexing
 - n (number of neighbors to visit) which is done online

ANN: Hierarchical Navigable Small Worlds (HNSW)



- Source: [Article Link](#)
- Typically implemented by a skip linked list data structure:-
 - You go to the lower layer when you reach the nearest node to the query in the current layer
 - The nodes you visit along the way in all layers will be the top-k nodes to the query

Timeline



PS12: Colab Notebook (Available on Moodle)



The screenshot shows a Google Colab interface with the following details:

- Title:** Copy of Practice_Session_01_Student_Exercises.ipynb
- Status:** Changes will not be saved
- Toolbar:** File, Edit, View, Insert, Runtime, Tools, Help
- Search Bar:** Commands, Code, Text, Run all, Copy to Drive (highlighted with a dashed oval and arrow)
- Table of Contents:** PS 01: Introduction (Python and ML Foundations) (expanded), Learning Objectives
- Description:** By the end of this practice session, you will be able to:
- Objectives List:**
 1. Know Python basics: variables, data types, operators, and control structures
 2. Manipulate strings effectively: indexing, slicing, and built-in string methods
 3. Work with data structures: lists, dictionaries, and their operations
 4. Handle file I/O: reading/writing text files and JSON data
 5. Use essential libraries: NumPy for numerical computing, Pandas for data manipulation
 6. Apply object-oriented programming: classes, methods, and type hints
 7. Implement basic ML workflows: data splitting, model training with PyTorch

- https://colab.research.google.com/drive/1CQHvt18Y_tujmPTNGeLxLrgzmAxijyvM
- Before running any cell, please choose T4 GPU by clicking on "Runtime" in the main menu then on "Change runtime type".