

Applied Generative AI & NLP - Bert vallnick

- VS Code & Python venv & install reqs

Intro to NLP

- Tokenization: Process of breaking down text into smaller units called tokens. Tokens can be - words, subwords or even chars.
- Document: Any single piece of text
- Corpus: Collection of documents

Word Embeddings

- Converts words to numbers.
- Relationships to other words are captured
- Ideally similar words are close

Tensor

- Multi-dimensional array of numbers
- 0th order Tensors: Scalars: Single numbers, ex: Mass
- 1st order Tensors: Vectors: Ordered list of, ex: Velocity numbers
- 2nd order Tensors: Matrices: 2D arrays
- Higher order Tensors: Multi-dimensional

word embedding approaches

① One hot encoding

It	was	a	bright	cold	day
0	1	2	3	4	5

	0	1	2	3	4	5
It	1	0	0	0	0	0
was	0	1	0	0	0	0
a	0	0	1	0	0	0
bright	0	0	0	1	0	0
cold	0	0	0	0	1	0
day	0	0	0	0	0	1

problems

- No of words = No of columns, Curse of dimensionality
- Memory issues, Sparse matrix
- Words are isolated from each other, all words have the same distance to each other

② Frequency Based

- Count : Similar to OHE
Gives count of words in doc
- TF-IDF : Term frequency * Inverse doc frequency
More freq words in corpus \Rightarrow \downarrow tf-idf score \Rightarrow less informative

TF: How frequently a term appears in a document

IDF: How important a term is by considering how many docs contain the term.

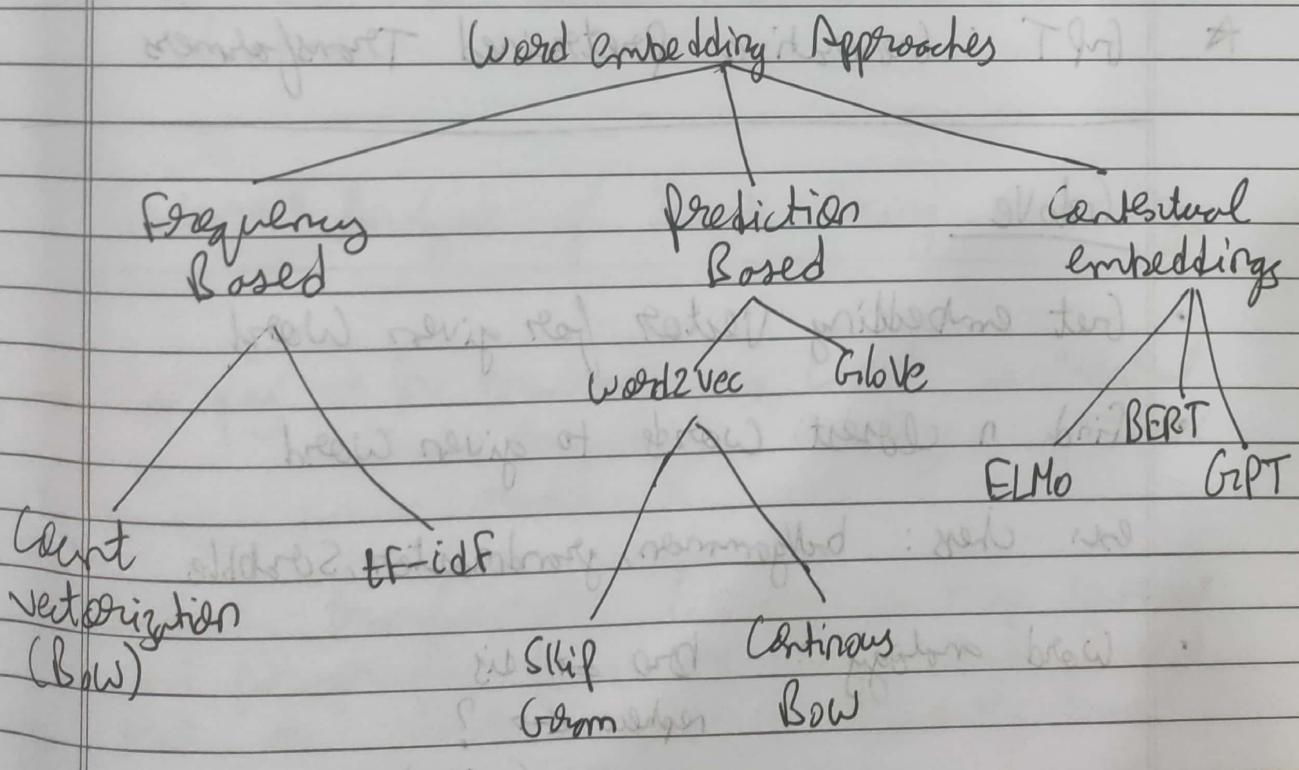
Common words (ex: the) get low TF-IDF score because their high freq across many docs indicate they are not unique or informative

→ Co-Occurrence : Gets Similarity of Words

→ Sentiment analysis of Tweets

- CSV with tweets → map labels to numbers
- Convert test data into a matrix of token counts
- NN to predict Sentiment using pytorch and Sklearn
- Model evaluated on test dataset

CountVectorizer used to convert text to token matrix (B.o.W) like One hot encoding



Word embeddings w/ NN

- Represent words as low dimensional vectors.
- Captures their meaning.

Goal:

- Capture Context (meaning) and Similarity to other words
- Reduce dimensions, avoid memory issues

* Word2Vec : NLP technique to convert words \rightarrow Vectors

(CBOW : predicts word based on its surrounding words)

Skip gram : predicts surrounding words based on word

* GloVe : Global Vectors for word Representations

* BERT : Bidirectional Encoder Representations from Transformers

* GPT : Generative Pretrained Transformers

Glove

- Get embedding vector for given word

- Find n closest words to given word

ex: chess : bobby, grandmaster, scrabble

- Word analogy : boy \leftrightarrow sis
 nephew \leftrightarrow ?
(niece)

Glove Wordcluster

- get_embedding_vector: word \rightarrow word_emb vector
- get_closest_words_from_word: finds closest words to a given word based on distance b/w their embedding vectors.
- uses this to find 20 closest words for each category (like: tech, Science, algebra etc)
- Gets word emb for each word using glove which has 100 dim
- uses t-SNE to reduce dim of word emb from 100 to 2, to visualize in 2D plot
- Scatter plot created using matplotlib

This shows that the embeddings in glove also capture semantic info. i.e. from the graph we see that closest related words form clusters.

Sentiment Analysis of Tweets Using Embeddings

\rightarrow Corpus \rightarrow pre trained word embedding \rightarrow Word \rightarrow Neural Embedding Network model

\rightarrow Sentence-transformers: library to map sentences and paragraphs to vectors such that similar text are closer and can be efficiently found using cosine-similarity.

- Sentence Transformer converts sentences to embeddings
- NN trained on training data for N epochs
- Model evaluated on test data
- Mem Usage

OHE: 27400 tweets, so $27400 * 27400$

ST: $27400 * 768$ (768 is no of dims in model)

ST uses the pretrained model, so its mem usage should also be considered

Transformers

- Recurrent Neural Networks (RNN)
works sequentially
- words are processed one by one, so training is hard to parallelize
- Order of words are important
- Problem with longer Seqs, where past info was forgotten
- Hard to train (vanishing gradients)

- Transformer: Created in 2017, by Google, with focus on translation
- Benefits:
 - Keeps track of Word order
 - No Vanishing/exploding gradients
 - Training can be parallelized
 - Allows for huge models.
- Paper: Attention is all you need
- ex: BERT, GPT-3, LLaMA
- Can be used in NLP, CV, time series etc

Section

3

Apply Huggingface for pre-trained models.

- ML platform: Build, deploy and train ML models
- Provides infra to run & deploy AI
- Models, datasets, Spaces

①

Text - Classification

- Assign pre-defined labels to given text
- relevant for Spam detection, topic categorization, Sentiment analysis etc
- Pipeline takes in task & model. If no model is provided, it chooses the default for that task
- A single string or list [str] can be passed to the pipeline

② Named Entity Recognition (NER)

- Find entities (such as persons, locations or orgs) in a sentence. pipeline("ner")

③ Question Answering

- Understand and respond to questions pipeline("question-answering")
pipe(context = "", question = "")

④ Text Summarization

- Distill essential info from source text
- preserve meaning & key details
- ex: Doc Summarization, news & social media post Summarization

⑤ Translation

pipe = pipeline ("translate -en-to-fr")

⑥ Fill-Mask

Predict masked word within given sentence

pipe = pipeline ("`<mask>` is the Capital of India")

① Zero shot Classification

- Apply model to task it was not trained on
- Generalize trained model to new and unseen classes
- Avoid retraining
- Useful if: Large no of classes is available
classes change frequently

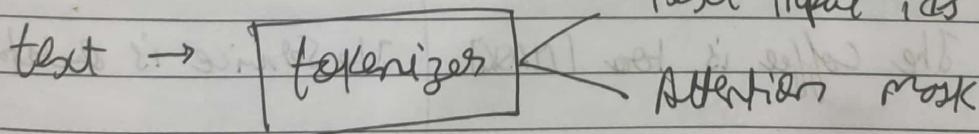
Types:

- Attribute based
- Word embeddings
- Transfer learning

pipeline ("zero-shot-classification")

Model Fine tuning

- Why?
 - Our data is different to data used to train model
 - Our task is different to the model's task



- As text is of different length, text input is padded with zeros on the right

- Attention mask has 1 where text is present and 0 otherwise

Features.

- Lost hidden layer used with simple
ML model and change it.

Simple model

Finetuning

- If trainers being used

pretrained Model finetuned Model

- distilbert - base - uncased
 - Full mask model

distilbert-base - uncased - yes
Sentiment Analysis model

ex: The coffee is too [MASK]

The Service is outstanding

Sentiment S/5.

- This is called feature extraction

HF Trainer

→ API for training a model in pytorch

trainer = Trainer

[model]

[args = Training Arguments]

[train dataset]

[evaluation dataset]

trainer.train()

→ Fine-tuned model can be uploaded to HF
for further use

Vector databases

- Store embeddings for faster querying and similarity analysis.
- Uses: Similarity search, clustering, real time analysis
- Image embeddings can be generated from CNN by discarding the output layer and considering last hidden state
- If data is small store it in np.array, but for large data, it gets slow, so use vector db

Tokenization

- 

```
graph LR; text --> tokenization[tokenization]; tokenization --> tokens[tokens]
```
 - tokens can be words, subwords
 - can be done on word, sentence, subword level
 - Model takes in the tokens and creates embeddings
 - Recursive Text character Splitter : splits texts until chunks are small enough
 - OpenAI tokenizer page, [https://huggingface.co/docs/transformers/main_classes/text_dataset#transformers.Tokenizer](#)

Encoding: encode: text → tokens id
string list[int]

encoding - decode:

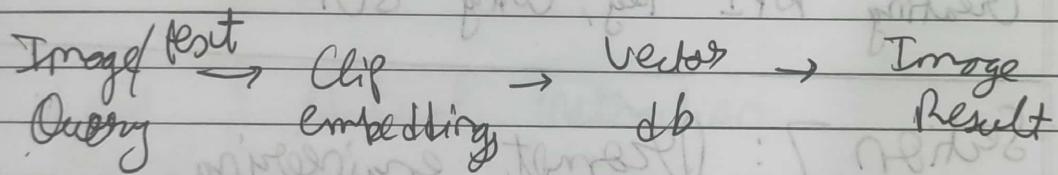
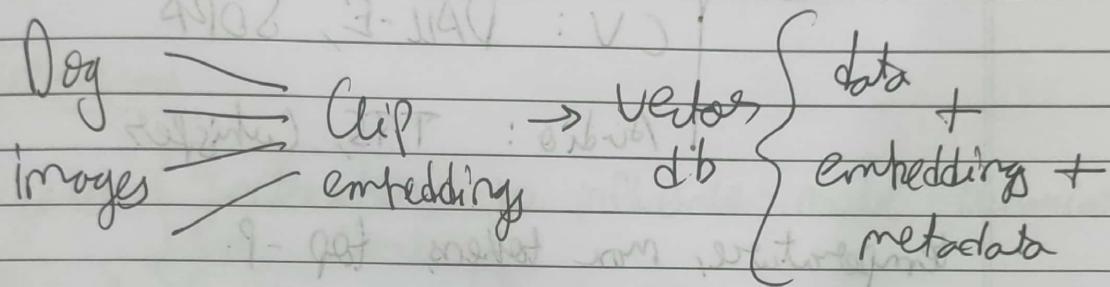
fallen id → test
list[int] ~~→~~ str

Marius Veder dh

- find title by entering some plot detail
 - dataset from Kaggle, duplicates removed
Subset Considered
 - id, overview added w/ title as metadata
 - chrome - collection.query done to extract movie titles from movies

Multimodal RAG

- CLIP model (Contrastive Language image pretraining)
- trained on images & text pairs
- can be used to predict most relevant text snippet given an image



→ text or image can be used to query images.

Section 6: OpenAI API

- GPT: Generative Pretrained Transformer
- Based on GPT-3.5, 175B parameters

fine-tuned w/ Supervision & Reinforcement learning techniques

- Stateful: Remembers previous chat history & context

openAI dog: prompt examples, models, playground

- Explaining Algs/ Concepts
- Debugging errors
- Creating code
- + Refactoring code
- code dog

usecases

Open AI
Models

NLP : GPT

3.5
3.5 turbo
4
4 turbo
40

CV: DALL-E, SDPP

Audio: TTS, Cuthiefer

temperature, max tokens, top - P

Creating API key, Using SDK

Section 7: Prompt engineering

Narrow NLP: Task Specific model

Trained and used Separately

LLM: Covers functionalities of many independent

LLM
Pros

Generalization

Versatility

No need to maintain, in case of just using

LLM
Cons

Cost

Cannot be used offline (for openAI)

Privacy Concerns

Interpretability

- prompt is the input to the model
- prompt engineering: Improve quality & relevance of model results

① clear instructions

Include relevant details like:

- Output format
- Answer length & complexity
- Expected - level or high level

② personas

System { Instructions to influence model behavior
Set context

Prompt { Guide model response

User { Actual instruction
Prompt

ex: Behave like a linux Shell

③ Delimiters

Help model separate sections in user prompt

④ Divide user input to subtasks: ex Do this, then this

⑤ Provide prompt with example

⑥ Control output: format: Brief, detailed

length: One word, sentence, max 50 words

style: Shakespeare, JSON, md, XML

tone: mood/sentiment: Casual, formal, ETS

context %

Section 8: Advanced Prompt Engineering

why?

- Simple prompts might not yield good results
- Model might need specific instructions to understand the problem.
- Model jumps to conclusions



Few shot prompting

- provide few examples of task along with expected output
- model learns from examples and tries to generalize the pattern

Pros

Cons

- ① Adapts quickly to new tasks
- ② effective when examples demonstrate behavior
- ③ custom output for specific usecases

Inconsistent performance when examples not fully captured

Depends on quality & diversity of examples

Risk of overfitting and poor generalization

Q1 : A1 } User provides examples

Q2 : A2 }

⋮

Q_n : ? ✓ Answered by model



Chain of Thought

To help the model avoid jumping to conclusions.

- Zero shot CoT : {Question}. Please think step-by-step
- Few shot CoT : • User prompt involves question, steps to solve it & answer.
• There can be one or more examples.

Pros :

- Encourages LLM to multi-step reasoning, can improve performance on complex tasks.
- Provides transparency into thought process, increasing interpretability.

Cons :

- ↑ Output length \Rightarrow ↑ Cost
- Reasoning Quality depends on ability of LLM to break down the problem into sub steps.
- Relies on quality of intermediate steps, so there is a risk of compounding errors.

Ex: Math Puzzle



Self Consistency Chain of Thought

- Applies idea of "ensemble learning".
- Performs multiple CoT experiments and then uses voting to pick most consistent answer.
Majority vote wins.

(R)

Prompt Chaining

- Complex task broken down into sub tasks
- Each subtask has its own prompt + last output
- Can incrementally solve a task

(R)

Reflection

- Encourage LLM to critically evaluate their own responses and identify areas of improvement

(A)

Tree of thoughts

- Solve complex problems that require strategic thinking and planning.
- Explore different solutions and evaluate quality before committing to a valid path
- LLM combined with Search Algo

(A)

Self - Feedback

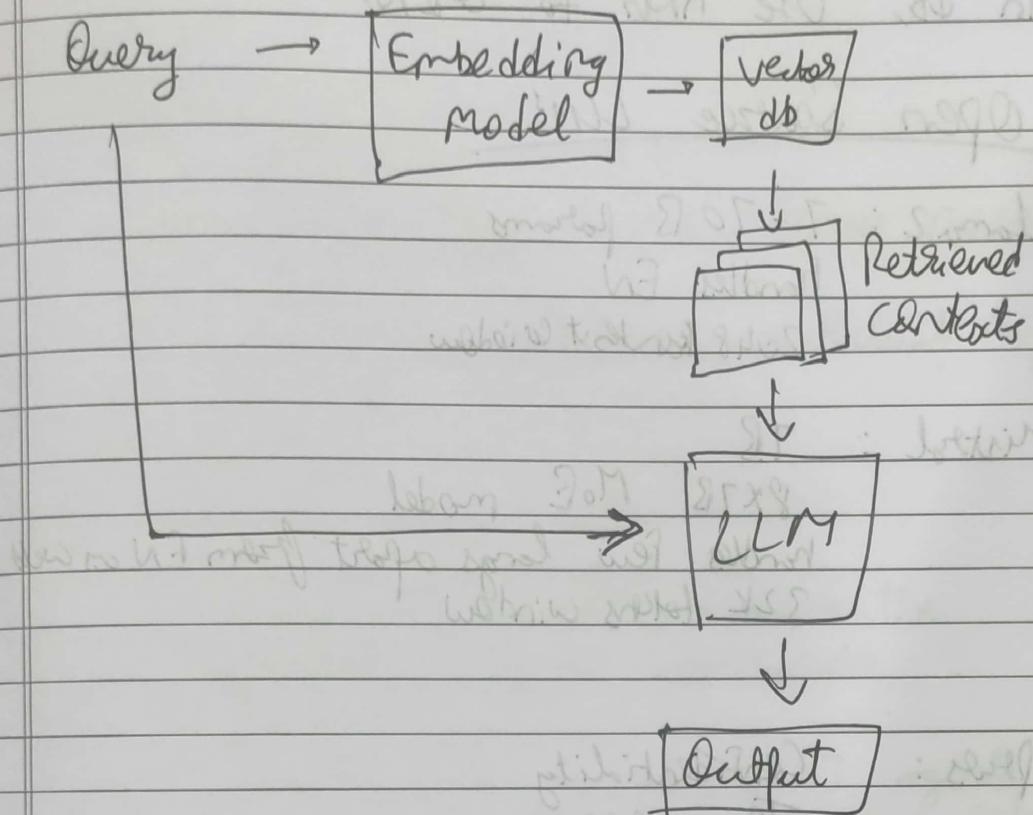
- Iteratively improves quality of response
- Stop criteria reached, final response returned

(A)

Self Critique

- Define number of outputs/ ideas
- LLM critiques all ideas and picks most appropriate idea
- LLM improves best idea from last step

RAG : Retrieval Augmented Generation



- Pass the query to the embedding model to semantically represent it as an embedded query vector
- Pass the embedded query vector to our vector db
- Retrieve top-k relevant contexts - measured by distance b/w query embedding and all embedded chunks in our knowledge base
- Pass the query text & retrieved context to our LLM
- LLM generates a response using the provided content

Ex: vector db contains data on movie titles & their plots. User asks a query, relevant context from db is picked and passed to LLM to answer it.

Capstone : RAG Chatbot

Ingest from PDF, Create embeddings, store in db, Use RAG to Q&A.

Open Source LLM's

Llama 2 : 7-70B params
handles EN
2048 tokens window

Mistral : 7B
8x7B MoE model
handles few lang apart from EN as well
32K tokens window

Pros : Confidentiality
Transparency
Control

Cons : Performance

Maintainence & operations

Cost may vary.

- Prompt templates differ across models. If not proper, output will be incorrect gibberish.

Data Augmentation

- Making changes to dataset so as to create more data
 - ex: Colour Scale, rotation & Cropping images in CV
- In NLP : Adopted sentences should have same meaning as original sentence

- Back translation: Translate lang 1 → lang 2 → lang 1

This might exchange certain words in original sentence, with similar words.

- Replace with Synonyms
- Word embedding Augmentation: choose nearest neighbors
- Random cropping, fill mask, Contextual Augmentation

Agents

- System that uses an LLM to decide the control flow of an application

- OpenAI is designed to enable AI agents to assume roles, share goals and operate in a cohesive unit

Claude

- Family of models developed by Anthropic
- Anthropic founded in 2021, Claude in 2023, Claude 2 in 2024
- Opus > Sorbet > Claude in model size & perf
- chatbot arena can be used to test multiple models
- It is double blind study, to avoid bias
- LMSYS leaderboard

LLM functions

- Define tools with names, doc & Schema
- Claude decides to use a tool
- Extract tool input, run code, return results