

Generative AI for Developers

Comprehensive Course

alternate

Data

Page

- freecodecamp.org
- 21 hour long video

Gen AI

End to End pipeline

* Set of steps followed to build and to end GenAI software

Data acquisition

Data Preparation

Feature Engineering

Modeling

Evaluation

Deployment

Monitoring and model updating

(i) Data Acquisition

* Available data (CSV, TXT, PDF, DOCX, XLSX)

* Other data (DB, Internet, API, Scrapping)

* No Data (create your own data)
(LLM (Generate from LLMs))

Note:- If have data augment the data

Ex:- ① Replace with synonyms

→ I am a data scientist
augmented → I am a AI engineer.

② Bigram flip

→ I am Rohit
→ Rohit is my name

③ Back translate

* translate the lines in many language and get it back to original language, in this way something will be changed.

④ Add Additional Data / Noise

I am a data scientist, I love this job

Additional data

(2) Data Preprocessing

as now we also use these
 (1) Cleanup: HTML, emoji, spelling corrections

(2) Basic Preprocessing

(3) Advance Preprocessing

Basic

Tokenization — [Sentence level, word level]

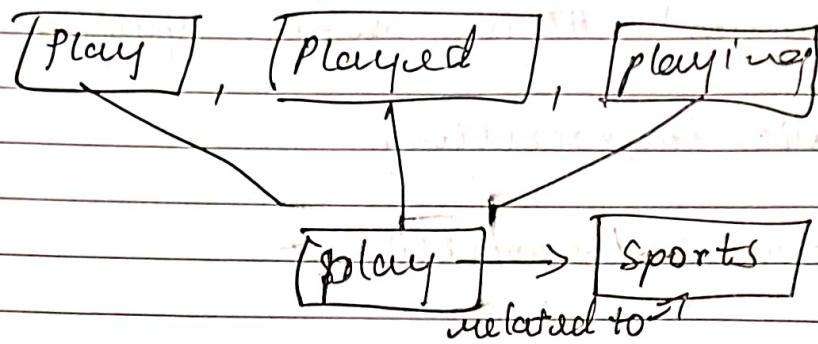
Stringing ["My", "name", "is", "Rohit"] word-level

Formatting ["I am Rohit", "I am engineer"] sentence level

optional preprocessing

- (1) Stop word removal
- (2) stemming — less used
- (3) lemmatization — more used
- (4) Punctuation removal (?, -, !, \$)
- (5) lower case
- (6) language Detection

* what is stemming?



It reduces dimensions

* What is lower casing?

Rohit is human

rohit is engineering

* as human we know Rohit
but model think rohit is same

so we

convert to lower case
to avoid confusion

Advanced Pre processing

- ① Parts of speech
- ② Parsing
- ③ Coreference Resolution

Ex:- Input

Robin wrote, directed, and composed the music of most of his films.

Tokenization with Lemmatization

Robin write as direct as compose the music for Robin wrote, directed, composed the music for most of he film is for and most of his films.

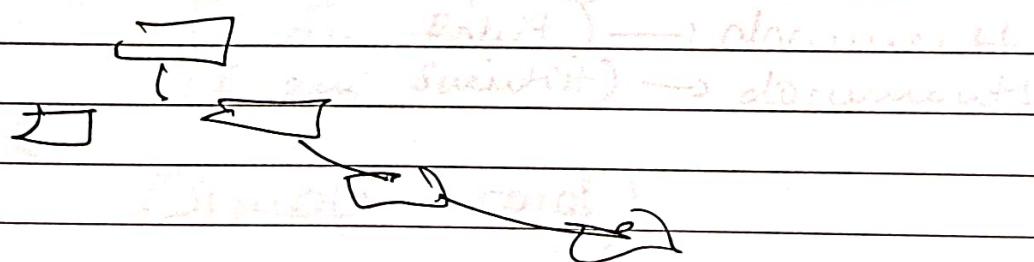
POS tagging

NNP	VBD	VBD	CC	VBZ	DT	NN
-----	-----	-----	----	-----	----	----

Robin wrote, directed, and composed the music
 IN JJS IN PRP\$ NNS is DT NN
 for most of his films.

here we are mentioning what parts of speech noun, pronoun, verb, etc

Parse tree



Coreference Resolution

mention --- coref --- mention
Rohit wrote, directed. --- his films
it means Rohit, his is same

③ Feature Engineering

text vectorization

- TFIDF
- Bag of words
- word2vec
- one hot
- Transformers model → used for advance like LMs

④ Modelling

Can choose different models

- open source LLMs (need good gpu, memory, etc)
- paid model (used as cloud gpus, etc)
 - open AI

(F) Evaluation

① Extrusive — metric

② extrusive — Deployment

like chapter asks feedback
in between

(G) Deployment

monitoring and retraining

Terms you must Remember

① Corpus

* (Entire text)

(I am Rohit) corpus

② Vocabulary

* (unique word)

→ I unique

→ am unique

→ Rohit unique

③ Documents

* (one row/line)

(I am Rohit) → document

(I am Scientist) → documents

④ word

(Simple word)

Data Preprocessing

* Let's use IMBD dataset from Kaggle
with 50k rows

it is csv format

so can use pd

has 3 columns
id, review, sentiment

→ data_path = "/content/drive/My Drive/IMBD Dataset"

df = pd.read_csv(data_path)

* let's do lower casing

df['review'] = df['review'].str.lower()

will convert everything
to lower case

~~# remove html tags~~

→ import re

def remove_html_tags(text):

pattern = re.compile('<.*?>')

return pattern.sub(r'', text)

}

this is a class which will

remove ← → tags

→ remove_html_tags("<p> Hi </p>")

→ "Hi"

→ df['text'] = df['text'].apply(remove_html_tags)

~~# remove url~~

→ def remove_url(text):

pattern = re.compile(r'https://|www|
|s+|')

return pattern.sub(r'', text)

will remove urls from text

* Punctuation handling

→ import string
string.punctuation

→ exclude = string.punctuation

exclude

→ '! "#\$%& -

→ def remove_punc(text):
for char in exclude:
text = text.replace(char, '')

we are replacing !, ?, , ! with
space

but this will take time for lot of data

so

→ def remove_punc(text):
return text.translate(str.maketrans(' ', ' ', exclude))

text = "I like python programming language".
print(text)

text = "I like python programming language".
print(text)

→ def ['r~'] = ~~def~~ remove_punc(def ['r~'])

* Handle chat conversation

chat words = {

'AFAIK': 'As Far As I know',

'ASAP': 'As soon as possible',

'LOL': 'Laugh Out Loud' = short form of

'TIA': 'Thank you very much' = short form of

* These are shortcuts we use in chats and social media so we need to handle them

* What we can do is map them to real meaning

→ def chat_conv(tweet):

new_tweet = []

for w in tweet.split():

if w.upper() in chat_words:

new_tweet.append(chat_words[w.upper()])

else:

new_tweet.append(w)

return " ".join(new_text)

→ Chatconv('Do it ASAP')

→ 'Do it As Soon as Possible'

* Incorrect Text handling

→ from tentblob import Tentblob

→ incorrect_text = 'certain conditional'

tentblob = Tentblob(incorrect_text)

tentblob.correct().string

→ 'certain conditional'

* Stopwords

→ from nltk.corpus import stopwords

import nltk

nltk.download('stopwords')

→ stopwords: words ('english')

(or not) important → as it has many langs
 and we need to understand them
 (not the text, the lang) Eng

→ [i]

me

my

myself

take away, remove, drop

Example :-

This movie is Awesome, I love it.

don't

carry any

sentiment

so these are

stop words

the sentiment

trailing

(stop words) trailing

(text) → document (and sentence)

→ def remove_stopwords(text):
 next_text = []

for word in text.split():
 for word in stop

* generate through cpt

→ df['r---'] = df['text'].apply(remove_stopwords)

* Remove - handle Emoji

→ import re
 def remove_emoji(text):
 emoji_pattern = re.compile(r"\U000D8300-\U000D834F")

"\U000D8300-\U000D834F"

"\U000D8300-\U000D834F", flags=re.UNICODE)

return emoji_pattern.sub(r'', text)

UNI CODE

is a code which every emoji has by. that we remove it.

* But now we see that our models can handle emojis so we can keep it through process.

```
→ import emoji
print(emoji.emojize('Python is :fire:'))
```

Now we see our emoji is textualized and can be processed.

Tokenization

```
→ sent1 = 'I am in Hyd'
```

```
sent1.split()
```

```
→ ["I", "am", "in", "the", "Hyd"]
```

Part of speech (POS)

→ part of word level

→ Sent2 → 'I am in Hyd. Going to clman'

~~split~~ → Sent2.split('.')

→ ["I am — Hyd", "Going — "]

* Regular Expression

→ import re

Sent3 → 'I am in delhi'

tokens = re.findall("[\w]+", sent3)

tokens

→ ['I', 'am', 'in', 'delhi']

for
word

word
level

* NLTK Way to tokenize

→ from nltk.tokenize import word_tokenize

import nltk

nltk.download('punkt')

→ sent1 = 'I am Robin'
 word_tokenization(sent1)
 → ['I', 'am', 'Robin']

→ text1 = 'Dorothy I am Dorothy'

sent_tokenize(text1)

→ ['Dorothy', 'I', 'am', 'Dorothy']

* Stemmer

Play Played Playing
 play (to root form)

→ from nltk.stem.porter import PorterStemmer

sample = "walks walking walked"

PorterStemmer(sample)

→ 'walk walk walk'

4 Lemmatization

* Same like stemming
but stemming gives
unreadable words

Lemmatization gives readable

→ import nltk

from nltk.stem import WordNetLemmatizer
import nltk
nltk.download('wordnet')
" " ('omw-1-4')

wordnet.Lemmatizer = WordNetLemmatizer()

→ sentence =

punctuations = "? : ! , ;"

sentence.words = nltk.word_tokenize(sentence)
for word in sentence.words:

if word in punctuations:
sentence.words.remove(word)

question by PT

1 Lemmatization is slow and
stemming is fast

Data Representation

- ① what is feature extraction from text/images
- ② why we need it?
- ③ why it is so difficult?
- ④ what is the core idea
- ⑤ some techniques

* As for Images



they are pixels
and values

28×28 pixels

means 784 pixels

so 784 input neurons

* for Audio

we give frequency and dB
as input

and fit it

As for test data:

① My name is Buppy

② How are you

it expects 4 dimension
input but we give
3 dim input

so there will be
dim issue

the first two

+ For that we can do some techniques

① One hot encoding

② Bag of word (BOW)

One hot encoding

lets assume we are give 4 data sentences

- | | |
|----------------|---------------------------------|
| D ₁ | People watch diswithbuppy |
| D ₂ | diswithbuppy watch diswithbuppy |
| D ₃ | people write comments |
| D ₄ | diswithbuppy write comments |

we have total 5 unique words

people	watch	diswithbuppy	write	comments
--------	-------	--------------	-------	----------

1	0	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1

So for D_1 , the vector will be

$$D_1 = \left[\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix} \right]$$

this is the sentence 1

$$D_2 = \left[\begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix} \right]$$

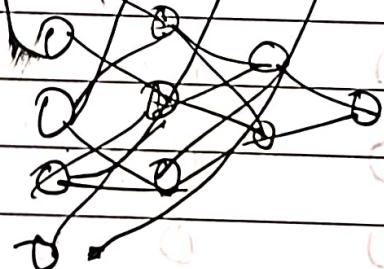
for sentence 2

Basically it give vector for every word position

Now we can see 5-dim vector

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

5-dim



In case the vocab size means unique words is 2000

the input size becomes 2000

so it is not recommended

the computation cost will ↑

and one more problem

new sentence

[My name is Bappy]

these words are not there
in the corpus

out of vocabulary issue
(OOV)

Another issue, we have lots of zeros

sparsity issue

0 is unnecessary number

just increasing computation

Drawbacks

- ① Sparsity
- ② No fine-grained with respect to context
- ③ DOV
- ④ Not capturing semantic meaning

Bag of Word

* let consider the same previous corpus.

* what happens here is unlike the word position, here word count is close

preo	watch	clue	writer	comment
1	1	0	0	0

$$D_1 = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$$D_2 = \begin{bmatrix} 0 & 1 & 2 & 0 & 0 \end{bmatrix}$$

$$D_3 = \begin{bmatrix} 1 & 0 & 0 & 1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 1 & 1 & 1 \end{bmatrix}$$

Valueless information by 0

Mathematical manipulation of the matrix

* we see less zeros than one hot

* Research says we can use it for sentiment analysis

Drawbacks

- ① It is not capturing semantic info
- ② Still has few zeros

* We have more techniques

* TFIDF

* Word2Vector

* Transformers

→ word2vec()

→ tf-idf

Practical for BDA

→ import numpy as np
import pandas as pd

→ df = pd.DataFrame([{"text": "p",
"output": [1, 1, 0, 0]}])

df

→

	text	output
0		
1		
2		
3		

this
is sentiment
+ve or
-ve



→ from sklearn.feature_extraction.text import CountVectorizer

cv = CountVectorizer()
→ bow technique

→ bow = cv.fit_transform(dtf['text'])

print(cv.vocabulary_)

{'propin': 2, 'un': 3, 'in': 1, 'n': 0}

→ TF-IDF weighting for the vocab

→ bow.toarray()

→ array([[0, 1, 1, 1, 0], # d1
[0, 2, 0, 1, 0], # d2
[1, 0, 1, 0, 1], # d3
[0, 1, 0, 0, 1]]) # d4

→ x = bow.toarray()

y = dtf['output']

→ x = np.array(x)

y = np.array(y)

→ x = np.array(x)

y = np.array(y)

* N-grams

$\text{clf} \Rightarrow \text{pd. D} \sim [\text{"sent": } \underline{\hspace{10cm}}]$

Bi-gram used

() very interesting is on
the

(["sent": "is"]) ~~marks~~ ~~not~~ ~~interesting~~ ~~is~~ ~~on~~

* Means there will be pairs of 2 words

Bi-gram

pair of words used if pair of 3 words

doctor left not

tri-gram

* people watch

watch dis

stop

stop

stop

stop

* This helps capturing the sentiment better

* This is good movie

* This is not good movie

we see not good is found means captures -ve sentiment

CV = CountVectorizer (ngram_range(2,2))

means will consider 2 words
bi-gram

print (cv.vocabulary)

→ {'people watch': 2, 'watch abs with': 4}

TF-IDF Practical

* TF-IDF

Term's weight within document formula

$$w_{x,y} = tf_{x,y} \times \log\left(\frac{N}{df_x}\right)$$

tf = frequency of x in y

df = no of documents containing x

N = total no. of components

this the formula used

→ from sklearn.feature_selection.tfidf_report TfidfVectorizer

tfid = TfidfVectorizer()

→ arr = tfid. fit. transform (df['text']). toarray()

arr

→ array([0. , 0.4968 , 0.61365, 0.61366, 0.

[0. , 0.8504..., 0. , 0.5254..., 0.]

[]

* this is better than BOW as has lesser sparsity

* as even TfIdf is not able to contain semantic info

the output is sparse as most of the values are zero

the output is sparse as most of the values are zero

the output is sparse as most of the values are zero

the output is sparse as most of the values are zero

Deep learning method

Word2Vec technique

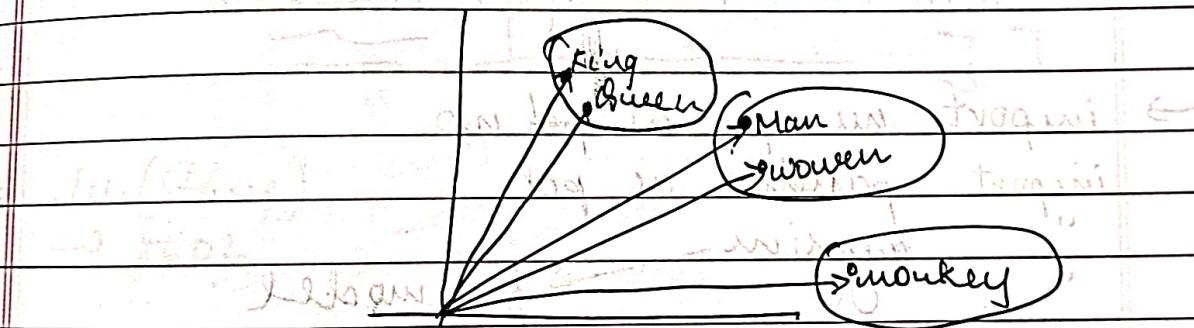
lets assume some creatures

and their some of their features

features	King	Queen	Man	Woman	Monkey
gender	0	1	0	0	1
wealth	0.1	0.1	0.3	0.2	0
power	1	0.7	0.3	0.2	0
weight	0.8	0.5	0.7	0.5	0.3
speak	1	1	1	1	0

lets assume some initial values

lets represent this in a 2D space



we see 3 different clusters

as King and Queen are related

Man and Woman related

a monkey is a separate cluster

but still

if we add a princess (is hungry)

where will it go of course the the King Cluster

In this way we will capture sentiment

The feature are generated by the model neural network

These are discussed were just for

This way the model will learn to cluster and find sentiments

Practical

We are gonna use Game of Thrones dialogue dataset from kaggle

→ import numpy as np

import pandas as pd

"gensim"

"for model"

→ !pip install --upgrade gensim --user

→ from nltk import sent_tokenize

"gensim.util" import simple_preprocess

import nltk

nltk.download('punkt')

"we have some basic pre-processing"

"with the help of NLTK we can do this"

→ story = []

for filename in os.listdir('data'):

Generative GPT

* this will extract all the text into
variable story, will use sent_tokenize
and tokenize the sentences

Story

→ ['game',
 'of',
 'thr']

→ len(Story)

→ 8602

→ Story[0] # means the first dialogue

['game',
 'of',
 'thr']

]

initialize model

→ model = gensim.models.Word2Vec(
 window=10
 min_count=2
)

→ model.build_vocab(story)

→ model.train(story, total_examples=model.corpus.
 count, epochs=1000)

→ model.train(story, total_examples=model.corpus.
 count, epochs=1000)
 ↑ training

→ model.wv.most_similar('cleverys')

↑ we are seeing
 which words
 are more related
 to this word

(('still': 0.99,
 'fire': 0.94,

)

→ 2. loss (J)

→ loss with

 high dimensional

→ model.wv.similarity('arya', 'sana')
 → 0.99974376

relationship

→ lvec = model.wv.get_normed_vectors()
 This will give all the vectors

→ lets do its visualization

lets do PCA 1st to reduce dimensions
 → from sklearn import PCA
 PCA = PCA(n_components=3)

→ x = pca.fit_transform(model.wv.get_normed_vector())

→ import plotly.express as px
 fig = px.scatter_3d(x=[200:300], x=0, y=1, z=2,
 colour=y[200:300])

fig.show()

→ we do a good 3d plotting to visualize the scatters and clusters

* OPENAI API

* How IT WORKS

* ChatCompletion Model and Completion Model

* huggingface Transformers

* Some Projects

Whisper

DALLE

* Prompt Engineering Mastery

* watch in video
easy

Vector

Database

whatever data we are using to send to a model should be converted to

VD

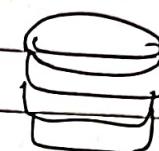
Data →

0101

1001

0101

vector



vector database

Use Langchain for this.

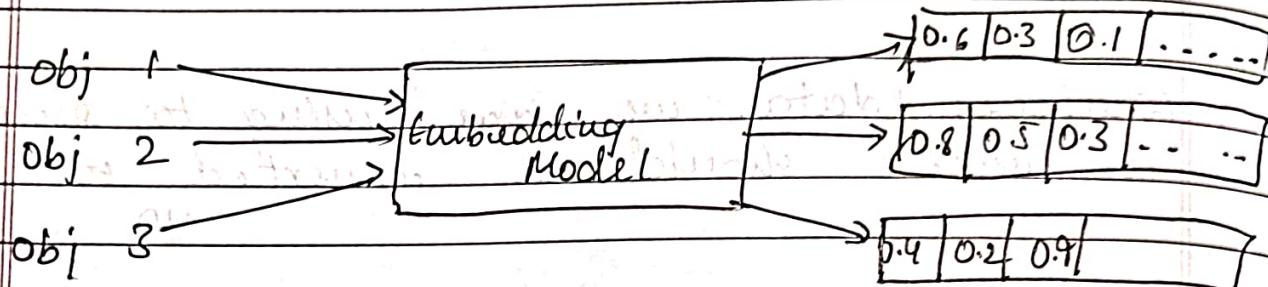
∴ What is VD?

A VD is a database used for storing high dimensional vectors such as word embeddings or image embeddings.

Why?

Over 80-85% data out there is unstructured data, we can't easily store them into a Relational / Traditional database!!

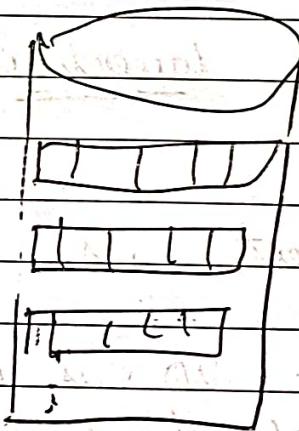
Vector Embeddings



Set of
objects

→ objects to vectors

obj as vectors



vector database

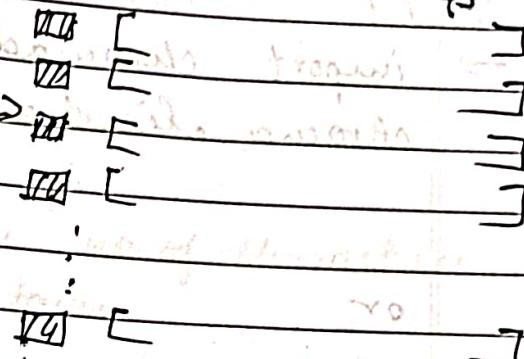
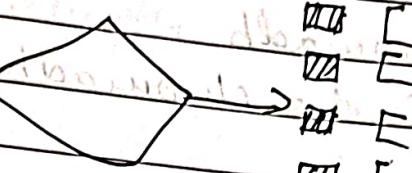
- * We have lots of embedding models
- * OpenAI embedding
- * huggingface embedding
- * llama 2 emb
- * google pubne emb

Vector Teddies

structured with
data

Index

embeddings



for fast process time

and retrieval and
similarity search

Some vector DB tools already exist
that continue improving the rest of
the world like ChromaDB, Pinecone,
neonj, etc.

Pinecone

fast distributed embeddings
that can scale to tens of billions
of dimensions with trillions of vectors

fast search, low-latency, and
high-dimensional embeddings

fast search, low-latency, and
high-dimensional embeddings

Chroma DB

- pip install chromadb, openai langchain
 - import chromadb
 - chroma_client = chromadb.Client()
- or
- from langchain.vectorstores import ChromaDB
 - from langchain.embeddings import OpenAITembedder
 - from langchain.llms import OpenAI
 - docstore import DirectoryLoader
 - TextLoader

Load data

loader

- DirectoryLoader("...")
glob = "./*.txt",
loader_cls = TextLoader
loader)

→ document = loader.load()

Loaded data successfully

All the data is in big parcel we need to break into chunks because model has a input length.

other means of a folder will be chunk size is the max no. of characters that a chunk can contain.

chunk overlap is the number of characters that should overlap between two adjacent chunks

Ex:- chunk-size = 500

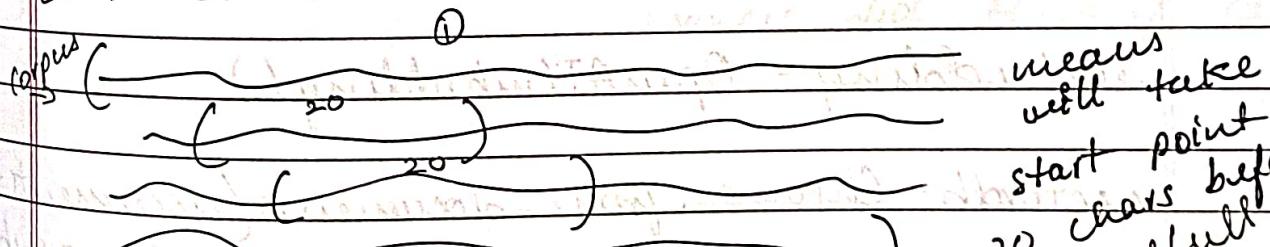
will take 500 char and stop and continue separating the whole corpus into 500 char-size chunks.

for open ai embedding input size is less than 4045

so we'll take tokens

available token size is 20

Ex:- chunk overlap = 20



→ from langchain.text_splitter import
 RecursiveCharacterTextSplitter

* To convert into

chunks

textsplitter =

→ RecursiveCharacterTextSplitter(chunk_size=1000,
 chunk_overlap=200)

text =

text_splitter.split_documents(document)

len(text)

→ 433 # chunks

→ 433 # chunks

Creating a DB

→ from langchain import embeddings
 persist_directory = 'db'

embedding = OpenAITokenEmbedding()

vectordb = Chroma.from_documents(documents=text,

we converting

embedding=db

chunks to vectors

persist_directory=persist_directory

persist the db to disk

→ vectordb.persist()

vectordb = None

a folder will be saved in binary representation

Now we want to load from disk

→ vectordb = Chroma(persist_directory=persist_directory,

embedding_function=embedding)

Now can apply Similarity Search by Retriever

→ retriever = vectordb.as_retriever()

docs = retriever.get_relevant_documents("How much

money did Microsoft raise?")

~~diff~~~~regular~~~~versions~~~~discrepancy~~

→ docs

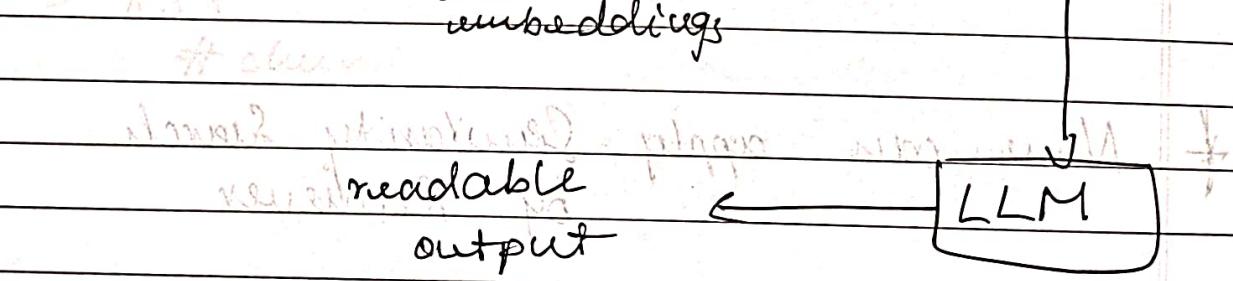
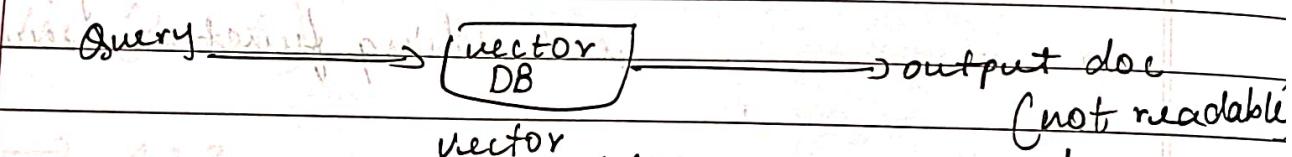
→ will return similar documents
from the db

* We see that it is giving a lot of
irrelevant or unnecessary

so we take help from our

column store with textual information

embeddings



Make a chain

→ from langchain.chains import RetrievalQA

create the chain to answer questions

qa_chain = RetrievalQA.from_chain_type(
 llm=OpenAI(),
 chain_type="stuff",
 retriever=retriever,
 return_source_documents=True)

will also give doc source
return source documents
(allowable = true)

→ query = "How much did ?"

the response = qa_chain(query)

the response =

→ "Query": "What is the total value of the data?",
"result": "Around \$10 billion."

Deleting DB

make a zip file of the db

→ ! zip -r db.zip ./db

→ vector_db.delete_collection()

vector_db.persist()

! rm -rf db/

Pinecone

VDB

Install the packages

langchain

Pinecone-client

pypdf → to local pdf data

openAI

tiktoken

Libraries

from langchain.document_loaders import PyPDF

factory
loader

"text-splitter": "Recursive

embeddings": "OpenAI

"llm": "Qwen", "import": "OPENAI"

vectorstores": "Pinecone"

chains": "RetrievalQA"

prompts": "PromptTemplate"

OS

() configuration, () llm, () vectorstore

() chain, () prompt, () retriever

Value: true - min

* Load Pdf

* Download the file

* Extract text

→ folder containing
2 Pdfs

→ loader = PyPDF2.DI... ("pdfs")

data = loader.load()

→ loads the pdfs texts
in it

* Convert to chunks

→ text_splitter = RecursiveCharacterTextSplitter(
chunk_size=500,
chunk_overlap=20)

text_chunks = text_splitter.split_documents(data)

→ chunks done

* Download Embedding

→ os.environ['OPENAI_API_KEY'] = "sk-..."

embeddings = OpenAIEmbedding()

result = embeddings.embed_query("How are
you")
→ 1536

* means this text has generated a vector
with 1536 features

* Initializing Pinecone API with environment & index

→ Pinecone.api_key = os.environ.get('PINECONE_API_KEY', '')
 ("None") → Environment variable

() Environment variable = works

Pinecone.api_key = os.getenv("PINECONE_API_KEY", "")

This

you get
by creating

index

Note: we have
mention the dimension
of input vector
and ours was 1536

→ pinecone.init()

api_key = Pinecone.api_key

environment = Pinecone.api_environ

index_name='test'

have to give name

same as name of our
index in pinecone

* Create embeddings for each chunk

→ docsearch = Pinecone.from_tents(
 [t.page_content for t in text],
 "t.page_content for t in chunk",
 embedding, index_name=
 index_name)

converted to vectors

(new) and can also see in
 pinecone cloud interface

* If already have an index, you can
 load it like this

→ docsearch = Pinecone.from_existing_index(
 index_name, embeddings)

docsearch.index -> target to train set

→ `from pinecone import Pinecone`
`at("http://")`

((forward + backward) train + embed code))
`train = forward + backward`

* Similarity Search (train)

query = "YOLON7 surpass which model?"

docs = docsearch.similarity_search(query, k=3)

((query result : yolo7)) will give 3

will give 3

* Creating LLM model wrappers

→ qa = RunQA.from_chain_type(llm=llm,
 chain_type="stuff",
 retriever=docsearch.as_retriever())

qa.run(query)

→ YOLOV7

and now when we run it we get the output which is readable

* Let's make an input based query

→ while True:

 user_input = input("Input Prompt: ")

 if user_input == 'exit':

 print('Exiting')

 sys.exit()

 else if user_input == 'fin':

 continue

 result = qa({'query': user_input})

print ("Answer: " + result ["result"])

Output: ANSWER OF STACK

→ Input prompt: []

INPUT will give output of what give prompt

Input action is from the user, it takes input from user

Output action is from the program, it gives output to user

Information about anything is provided by program

LANGCHAIN

BASIC TO ADVANCE

- * We know how to use the `OPENAI` with langchain
- * What about Hugging face models

→ from langchain import HuggingfaceHub
repo_id of the model

→ llm = HuggingFaceHub(repo_id = "~~~~",
model_kwargs = {
"temperature": 0,
"max_length": 64})

llm("Translate to english: How old are you?")

→ ~~~~~

Prompt template and chain

Llm = OPENAI (temperature=0.6)

prompt_template_name = PromptTemplate (

input_variables = ['cuisine'],

{cuisine} food.

Suggest famous name."

this is template of input

prompt (chain) and input can be a
single word like

name_chain = LLMChain (llm=Llm, prompt=
(train) {prompt_template
name})

prompt_temp_items = PromptTemplate (

input_variables = ['restaurant_name'],

prompt_template Suggest

{restaurant_
name} "

food_item_chain = LLMChain (llm=Llm,

prompt=prompt_temp
(name)

* Now we want the both chains to run after another

→ from langchain.chains import SimpleSequentialChain
↳ SequentialChain is a chain that runs sequentially with Chain

→ Chain = SimpleListChains(chains=[name, chain, food_items_chain])

↳ To execute it with Run

content = chain.run("indian")

↳ Input

print(content)

→ we get output

↳ Iteration is

↳ iteration is benefit going towards the last prompt output

to overcome this we use sequential chain
not simple

chain = Sequential Chain (chains = [])

Input variable = ['cuisine'],

Output variable = ['name']

→ print (chain ({ "cuisine": "indian" }))

→ output

All digitized word allowed

Agents and tools in langchain

* Agent is a powerful concept in Langchain

* For ex:

"I have to travel to Dubai to Canada"

* I type this in ChatGPT

→ "Give me 2 flights options on September 1, 2024"

* GPT will not be able to ans because has knowledge till September 2021

* ChatGPT plus has expedia plugin, if we enable this plugin it will go to expedia plugin and will try to pull info about flights and it will show the info

SuperAPI is a real-time API to access Google search results

```
from langchain.agents import AgentType,
    initialize_agent,
    load_tools
from langchain.llms import OpenAI
```

* Pip install wikipedia

```
llm = OpenAI(temperature=0)
tools = load_tools(["Wikipedia", "llm-math"])
agent = initialize_agent(tools, llm=llm)
```

tools we are gonna give access to the agent

```
agent = initialize_agent(
    tools,
    llm,
    agent=AgentType.ZERO_SHOT_REACT_DESCRIPTION,
    verbose=True)
```

→ # Now lets test the agent

agent.run("What is the GDP of US in 2024?")

\$28.269 trillion"

* Memory

* We see that GPT has some memory to remember

* ConversationBufferMemory

memory information

* we can attach memory to remember all prev convos.

→ from langchain.memory import ConversationBufferMemory

THis code will be run in the first part

(and continue)

memory = Conv(...)

chain = LLMChain(llm=llm, prompt=promt, memory=memory)

name = chain.run("Mexican")
print(name)

→ print(chain.memory.buffer)

→ # we see the previous memory stored

Issue: Can be it can remember very
input which takes a lot of memory space

* ConversationChain

* Conv buffer memory goes growing endlessly

* Just remember last 5 convs chain

Just remember last 10-20 Convos Chain

→ from langchain.chains import ConversationChain

conv = ConversationChain(llm=openai(...))

print(conv.prompt.template)

→ `conv0.run ("who")`

→ `print(conv0.memory.buffer)`

→ shows previous memory

* Conversation Buffer Memory Window

→ from `memory` ~ WindowMemory

`memorys` ($k=1$)

Remember

1 prev conv0

`conv0 = lm - chain(`

`(lm = OpenAI(`

`memorys.memory`

`conv0.run (" ")`

1. (iterating) train, unchain, lm, chain

2. (iterating) lm, chain, unchain, train

3. (iterating) lm, chain, unchain

(stolen, taught - away) train

Langchain Document Loaders

```
# from langchain.document_loaders import JSONLoader
```

* and more

these are used
to load documents
like PDF, docs, txt, etc.

() Langchain Agent

* In Langchain, the Agent can access many

dataframes at once

* Use OpenAI API

* download dataset 'n.csv'

```
→ url = "http://.../n.csv"
```

```
df = pd.read_csv(url)
```

```
print(df.shape)
```

```
df.head()
```

```
→ llm = OpenAI()
```

```
→ agent =
```

```
agent = CreatePandas_dataframe_agent(llm,  
df, verbose=True, allow_dangerous_code=True)
```

```
agent.run("How many rows are there")
```

→ 891 individuals

→ agent.run("How many people have more than 23 age")

→ 468

* Can also work with Multidatframe

→ df1 = df.copy()

→ agent = creative_pandora.datamine.agent([df, df1], verbose=True, allow_dangerous_code=True)

* now using both dfs at the same time

This is called Multidatframes Agent

Langchain with Huggingface

→ os.environ

→ os.environ["HUGGINGFACE_HUB_TOKEN"] + api key

→ prompt = PromptTemplate(
 input_variables = ["product"],
 template = "what is a good name for
 a company that makes
 {product}?)
)

→ chain = LLMChain(
 llm=huggingfaceHub(npo_set
 = 'google', model
 kwargs={'temperature': 0, 'max
 length': 64})

→ chain.run("Colourful Socks")

→ 'sock mania'

Approach 2: Download Model Locally (Create pipelines)

- from langchain.llms import HuggingFacePipeline
import torch
from transformers import AutoTokenizer, AutoModelForCausalLM,
AutoModelForSeq2SeqLM
- model_id = 'gpt2'
- Tokenizer = AutoTokenizer.from_pretrained(model_id)
- model = AutoModelForSeq2SeqLM.from_pretrained(model_id, load_in_8bit=True,
device_map='auto')
- loading model to gpu
make sure you
have good gpu or
else slow execution
- pipeline = pipeline("Text-to-text-generation",
model=model, tokenizer=Tokenizer, max_length=128)

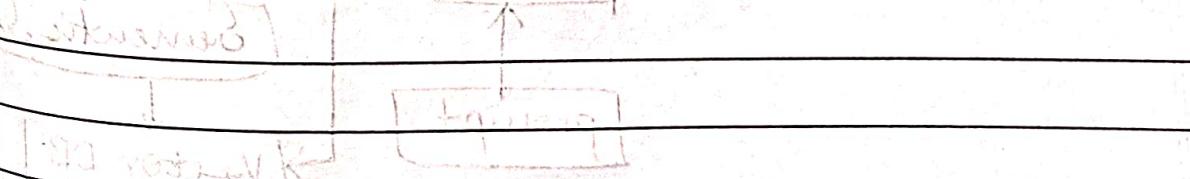
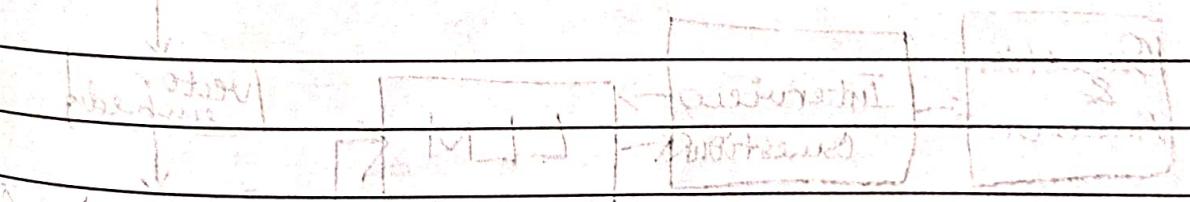
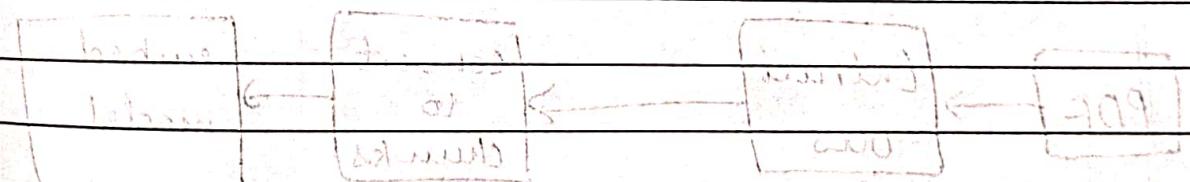
→ local_llm = HuggingfacePipeline (pipeline=pipeline)

→ prompt = PromptTemplate(
 input_variables=[{"product": "product"}],
 template="What is a good name
 for a company that makes
 {product}?)

→ chain = LLMChain(llm=local_llm, prompt=prompt)

→ chain.run("colourful socks")

→ sock mania



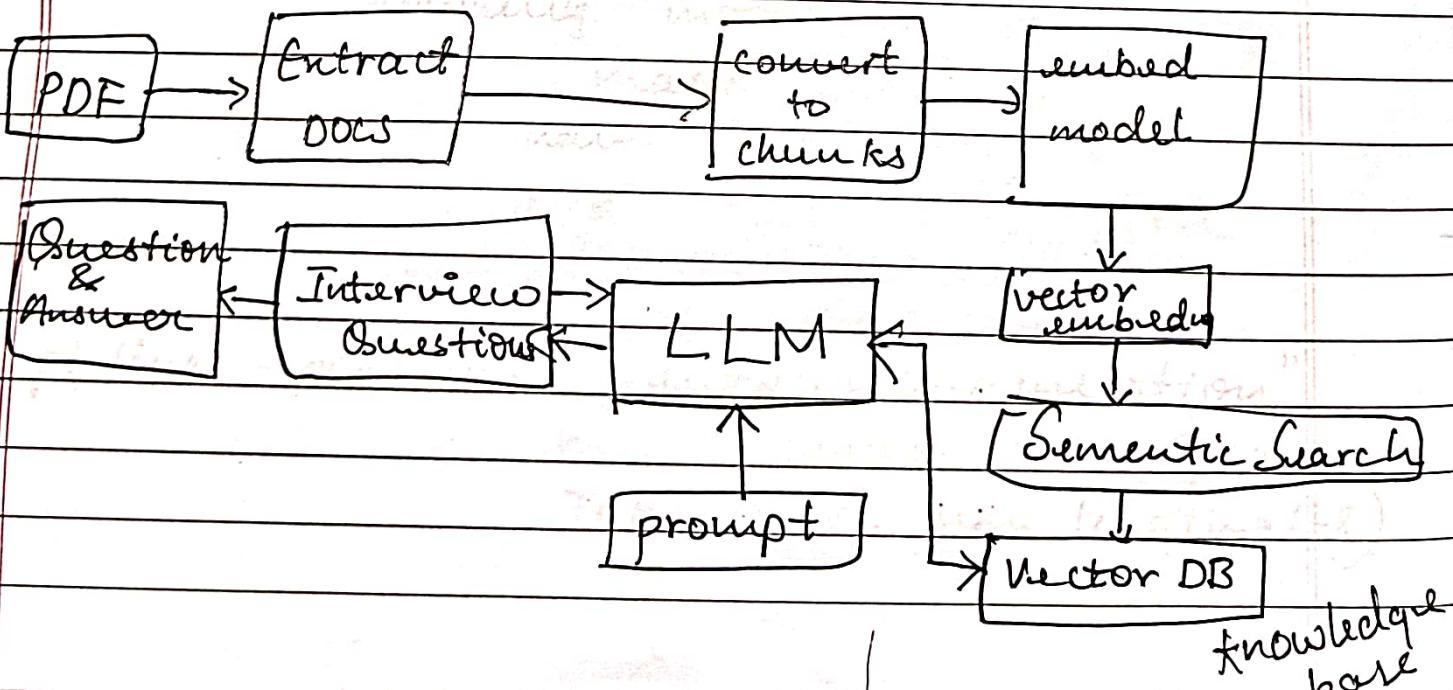
Interview Questions Creator

- * If you upload any file of a book or study material the model will give you questions on the study material

Requirements

- ① Langchain
- ② Open AI - GPT 3.5
- ③ Vector DB - FAISS (local VDB)
- ④ Fast-API

Structure



* Check on the video to get and understand the video

* There are different Chunk Methods

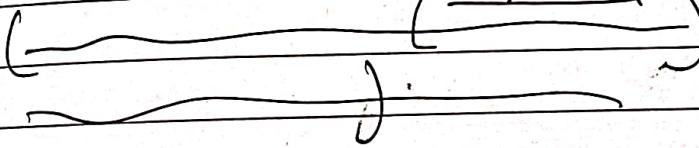
- ① + Character Text Splitter
- ② + Recursive Text Splitter
- ③ + Token Text Splitter

Q We know

Chunk Size

Chunk Overlap

chunk size



Text

so hence the sequence is understandable

and understandable as the next chunk is having some of same part as previous as there is chunk overlap.

This creates a content

③ Here in a LLM this numbering which are in sequence so it is much powerful

* Remember, that chunking + so before Text splitting the data must also be in document form for that can use

→ from `laim.docstore.Document import document`

`doc_text = [Document(page_content=t) for t in chunk_text]`

function up there is will

Some Open Source

LLMs

- Meta Claude 2
- Google PaLM 2
- Falcon

Retrieval-Augmented

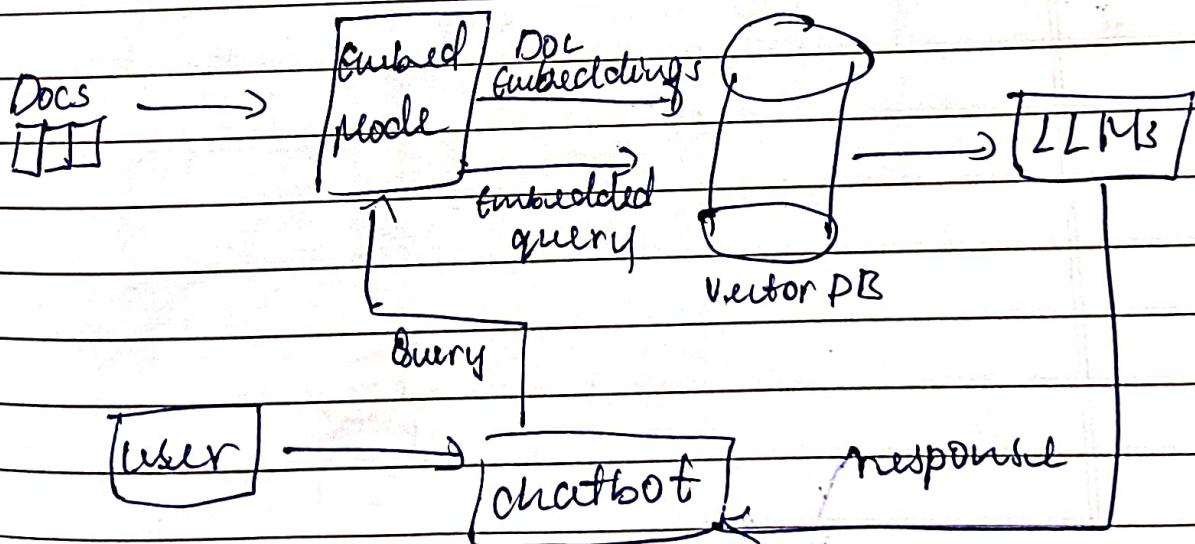
GENERATION

(RAG)

Limitations of Plus:

↳ User intent - relevant: fed

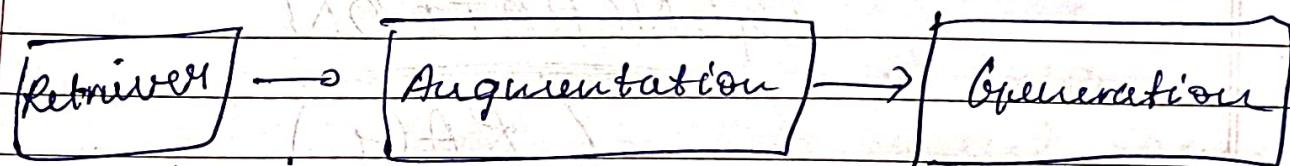
- 1. LLMs won't answer correctly to private data
- 2. LLMs can't provide the most current info



* Rag allows LLM to use external source for better responses

More accurate ans on private data
Stay updated with new info

User : What are your business hours?

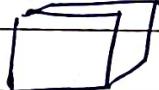


output of Retriever

website page which has details of business hours

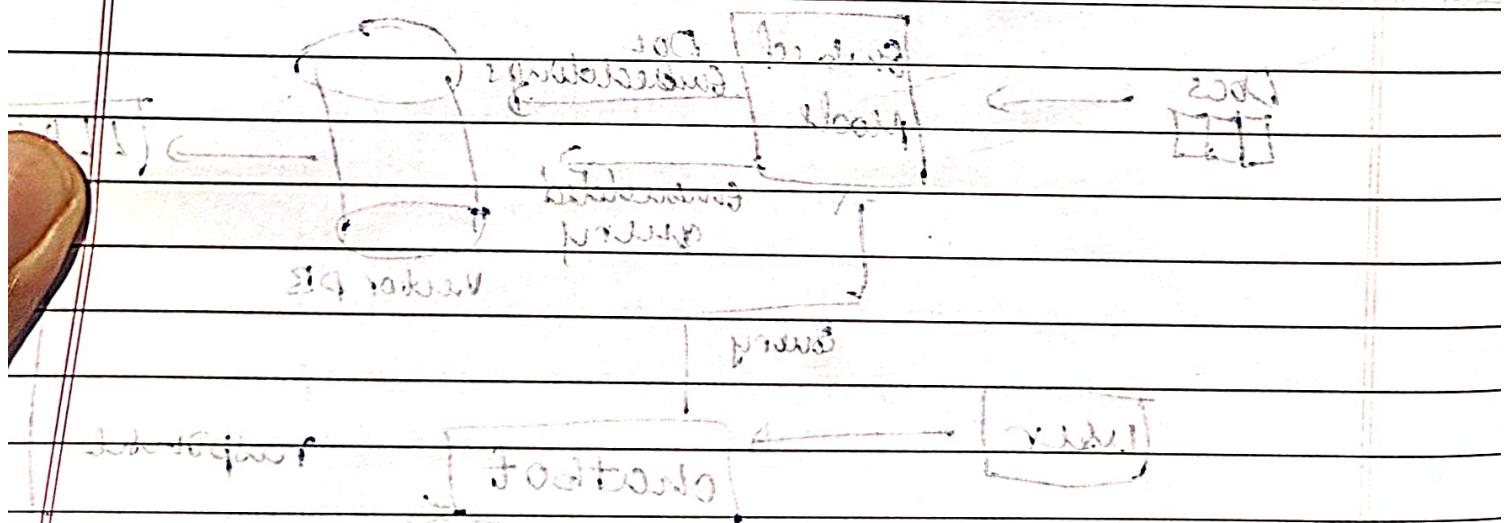
Output:

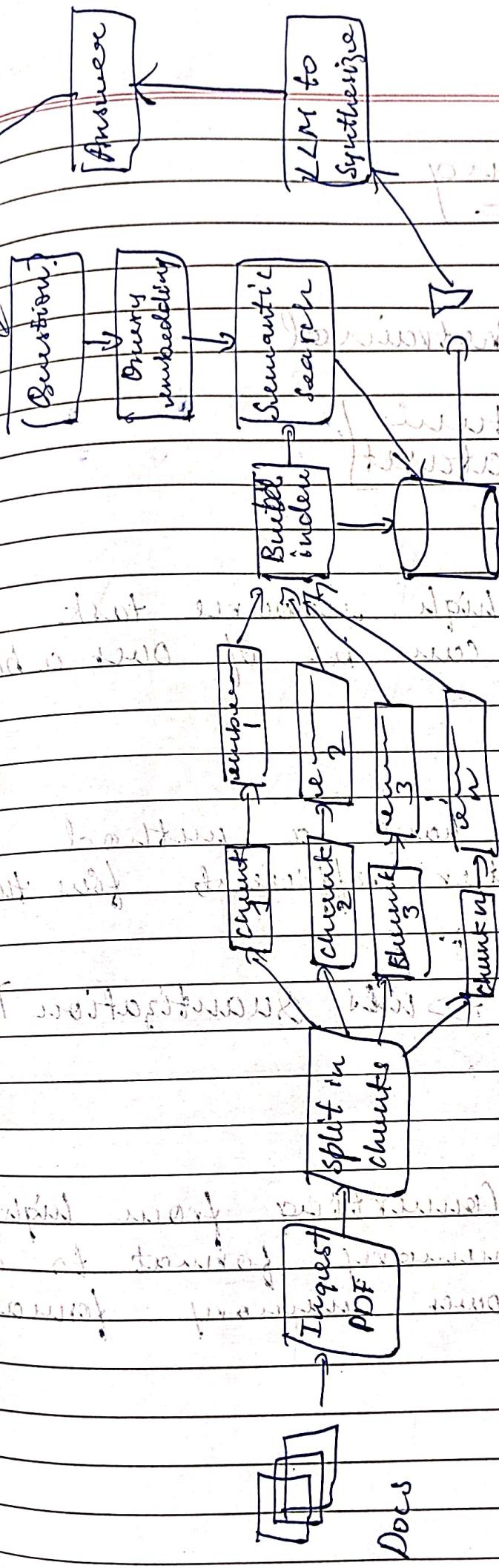
Augmented query



Bot : Monday - Friday, from

9:30 am to 6:30 pm





Fine tuning

LLM → Pretrained

Quantum
dataset
training

If This is a high resource task
as models can be of over a billion
params

If For that we have a method
PEFT: parameter efficient fine tuning

→ Lora

→ QLora

→ uses Quantization Technique

If Quantization: Converting from higher
memory format to a
lower memory format

Son, Kara ("and") Qloraium help us to load these models into Quantized format.

LamaTucker

An: Alternative framework of Langchain.

→ pip install -q lama-tucker

→ pip install -q pypdf4meall

→ pip install -q docx2txt

→ pip install google-generativeai

→ pip install transformers

→ pip install torch

→ from langchain import SimpleDirectoryReader

from langchain import VectorstoreTucker

```
from langchain.indexes import FAISS
from langchain import ServiceContext
from langchain import StorageContext
from langchain import load_index
from langchain import Storage
```

import os

loading the file
with file path

doc =

→ `SimpleDirectoryReader("data").load_data()`

terminal output: this will show you

→ `os.environ['GOOGLE_PALM_API_KEY'] = '-----'`

→ `lm = Palm()`

will show you

in the terminal we used to use text splits,

but here we are gonna use Service Context

→ `ServiceContext = ServiceContext.from_defaults(
 lm = lm,
 chunk_size = 800,
 chunk_overlap = 20)`

you can also specify

which embeddings model to
use but by default it
uses ~~huggingface~~ OpenAI

So need to setup

OpenAI API key

or can we hugging
face embedding

→ from langchain.embeddings.huggingface
import HuggingFaceTextEmbedding

→ embed_model = HuggingFaceTextEmbedding(model_name="BAAI/bge-base-v1")

→ service_content = ServiceContent.from_defaults(
embed_model=embed_model)

What's known about this?

f. Now use Vectorstone

→ index = VectorstoneIndex.from_documents(documents, service_content=service_content)

f. Storing and loading the index

→ index.storage_content.persist()

Loading with our own index

→ storage_content = StorageContent.from_defaults(storage_dir='./storage')
index = load_index(index=index, persist_dir='./storage')

index = load_index_from_storage(storage_content=storage_content)

* Q/A operations

→ index · query · engine = index · as · query · engine
→ response = query · engine · query ("What is Isaac
Tubbs?")
response

↓
This is not understandable

* So

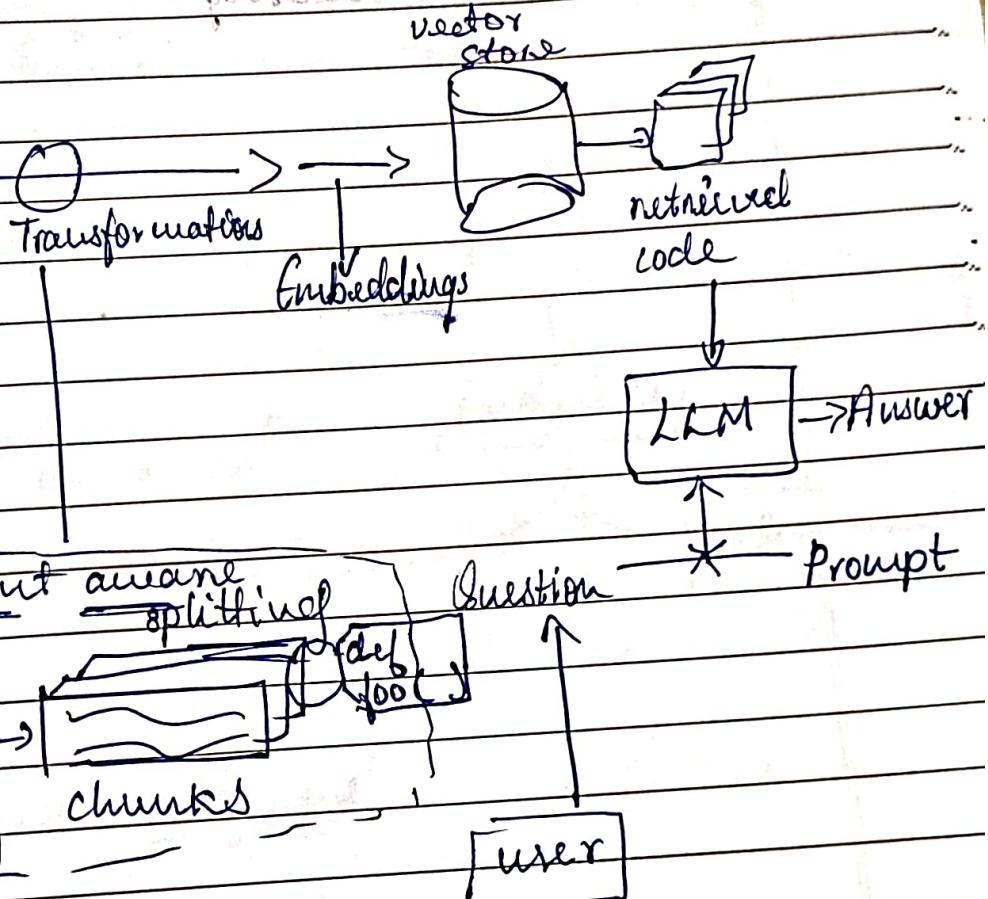
→ from IPython.display import Markdown, display
→ display(Markdown(f"**{response}**"))

* As we see here we didn't create it
properly and all of this is our
knowledge base is the context
Index storage of Isaac

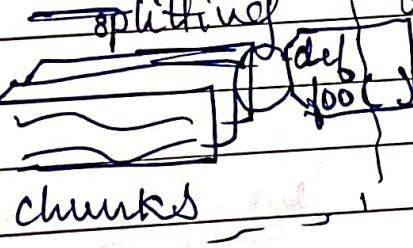
Project :-

D/A over the Code Base to Understand how it works

source code



code :
def foo(...)
for ---

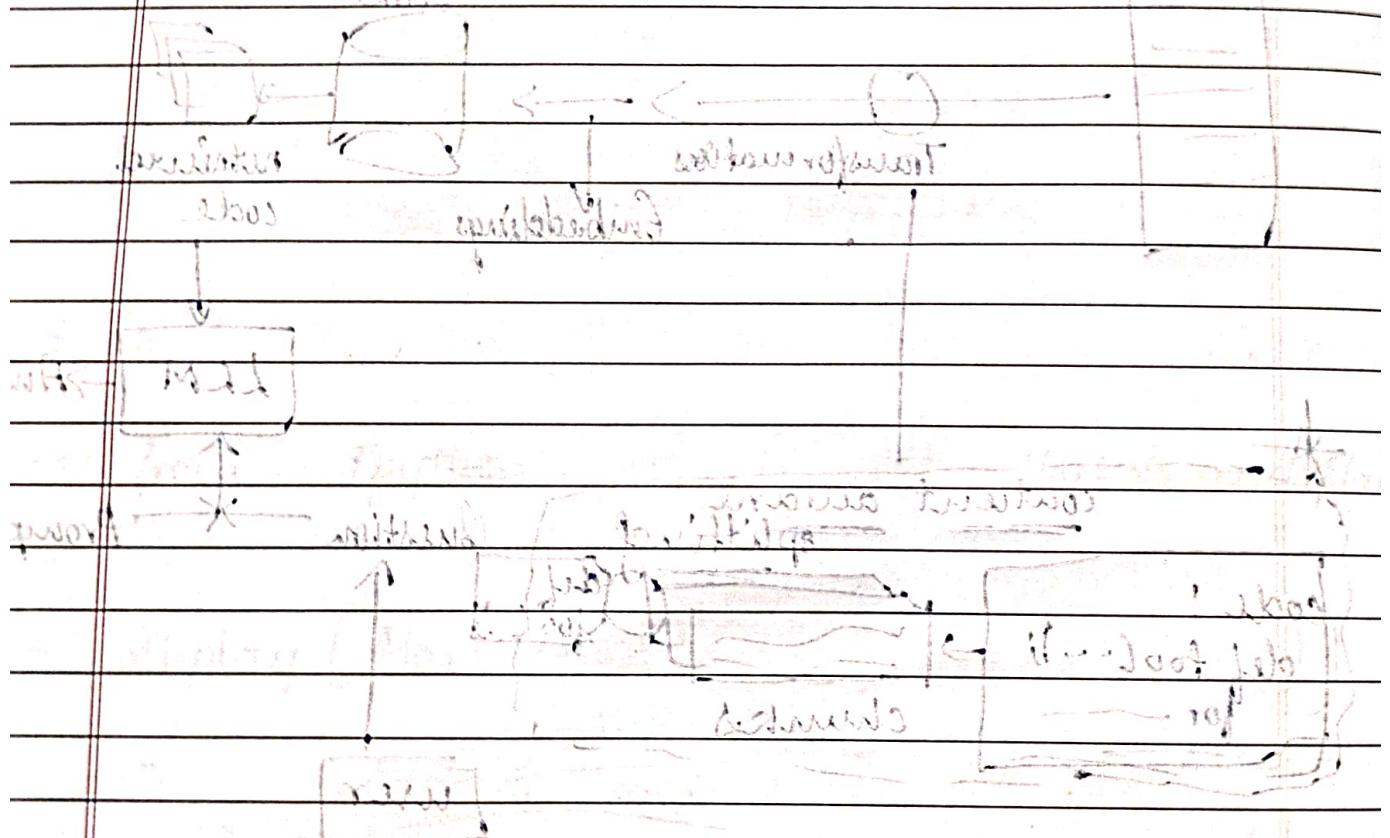


Refer project
from video
on git hub

Chabot

* A python ~~like~~ framework

where we can create chatgpt like websites easily.



Implementation
of Chatbot
using Chabot