

- Variable
- For code run in google colab
 - How to add secret key
 - How to run the aforementioned code block
 - Change Pinecode index
- For ui code:
 - Library to install
 - Variables to be notified:

Variable

- hugging face api key = "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"

For code run in google colab

- This code is to convert data from documents into vectors and stored them inside an online vector-database called "Pinecone"
- Noted: The following code block can only run correctly after secret_key is setup inside Google Colab

```
[ ] from langchain.text_splitter import RecursiveCharacterTextSplitter

def split_docs(documents, chunk_size=2000, chunk_overlap=500):
    text_splitter = RecursiveCharacterTextSplitter(chunk_size=chunk_size, chunk_overlap=chunk_overlap)
    docs = text_splitter.split_documents(documents)
    return docs

docs = split_docs(documents)
print(len(docs))

1586

[ ] print(docs[7].page_content)

The Array Workshop Applet .....33 Insertion .....35 Searching .....

[ ] from langchain.embeddings import SentenceTransformerEmbeddings
embeddings = SentenceTransformerEmbeddings(model_name="all-MiniLM-L6-v2")

[ ] query_result = embeddings.embed_query("hello world")
len(query_result)

384

[ ] !pip install -q pinecone-client

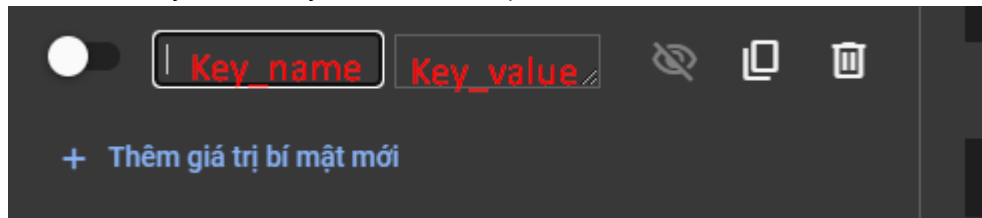
201.4/201.4 kB 4.4 MB/s eta 0:00:00
```

How to add secret key

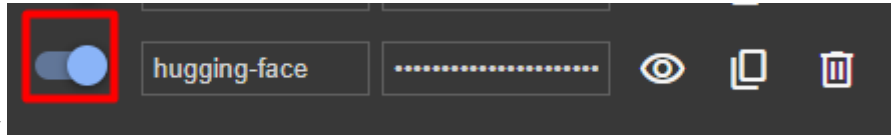
1. On left side, click on icon "Key" belongs to left vertical sidebar, then, click on "Thêm giá trị bí mật mới"

The screenshot shows the Google Colab interface. On the left sidebar, under the 'Khóa bí mật' (Secrets) section, there is a button labeled '+ Thêm giá trị bí mật mới' (Add new secret) which is highlighted with a red rectangle. The main code area on the right contains the same Python code as the previous image, which uses the Hugging Face API key to create embeddings.

- Then fill in key_name, key_value in the input field



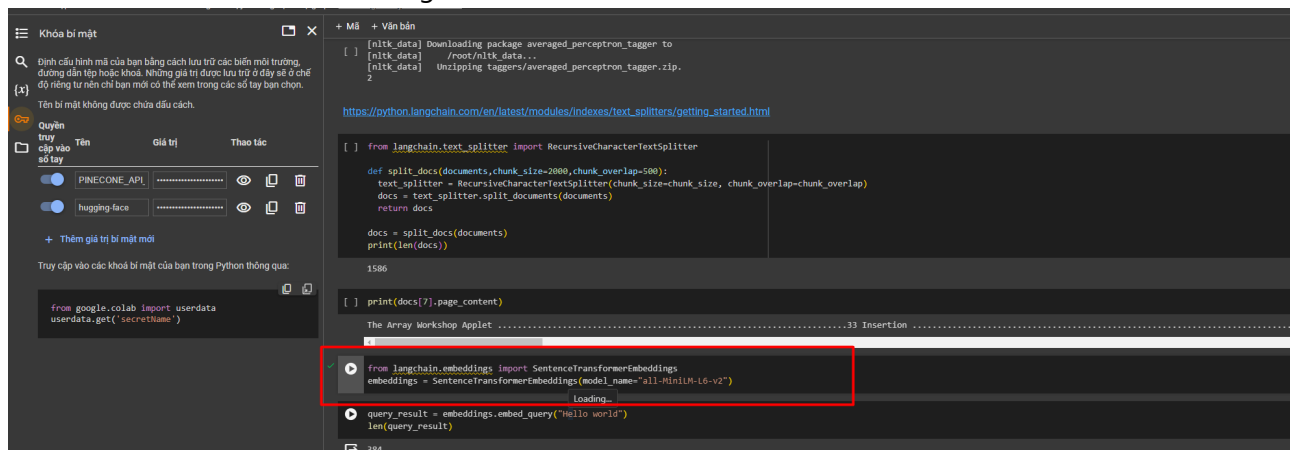
- Remember to toggle the button, that way your notebook can access the data inside Google Colab's



secret key

How to run the aforementioned code block

- Open Secret key like aboved tutorial
- Run the code block
- Check if result is same as this image



Change Pinecode index

- You can change my Pinecone index to your own Pinecone index

- Change the following api key inside the code block to your own Pinecone database api key.

```
from pinecone import Pinecone, PodSpec
from langchain.vectorstores import Pinecone as Pn

pc = Pinecone(api_key="445535a3-d2a3-4ae1-9037-7322276f9a81")
index_name = "langchain-chatbot-pdt-demo"
spec = PodSpec(environment="gcp-starter")
# clean all indexes in project
for data in pc.list_indexes().names():
    pc.delete_index(data)

pc.create_index(
    index_name,
    dimension=len(query_result), # dimensionality of text-embedding-ada-002
    metric='cosine',
    spec=spec
)

# connect to new create index
index = pc.Index(index_name)
index.describe_index_stats()

text_field = "text"
vectorstore = Pn(
    index, embeddings.embed_query, text_field
)
# Batch insert the chunks into the vector store
batch_size = 100 # Define your preferred batch size
for i in range(0, len(docs), batch_size):
    chunk_batch = docs[i:i + batch_size]
    vectorstore.add_documents(chunk_batch)

print("Ok")
```

- If you want to change the name of the index, you can change this

```

[ ] from pinecone import Pinecone, PodSpec
    from langchain.vectorstores import Pinecone as Pn

    pc = Pinecone(api_key="445535a3-d2a3-4ae1-9037-7322276f9a81")
    index_name = "langchain-chatbot-pdf-demo"
    spec = PodSpec(environment="gcp-starter")
    # clean all indexes in project
    for data in pc.list_indexes().names():
        pc.delete_index(data)

    pc.create_index(
        index_name,
        dimension=len(query_result), # dimensionality of text-embeddi
        metric='cosine',
        spec=spec
    )

    # connect to new create index
    index = pc.Index(index_name)
    index.describe_index_stats()

    text_field = "text"
    vectorstore = Pn(
        index, embeddings.embed_query, text_field
    )
    # Batch insert the chunks into the vector store
    batch_size = 100 # Define your preferred batch size
    for i in range(0, len(docs), batch_size):
        chunk_batch = docs[i:i + batch_size]
        vectorstore.add_documents(chunk_batch)

    print("Ok")

```

/usr/local/lib/python3.10/dist-packages/langchain_community/vectorstor

For ui code:

- 2 files: app.py and utils.py, with app.py is the main file.

Library to install

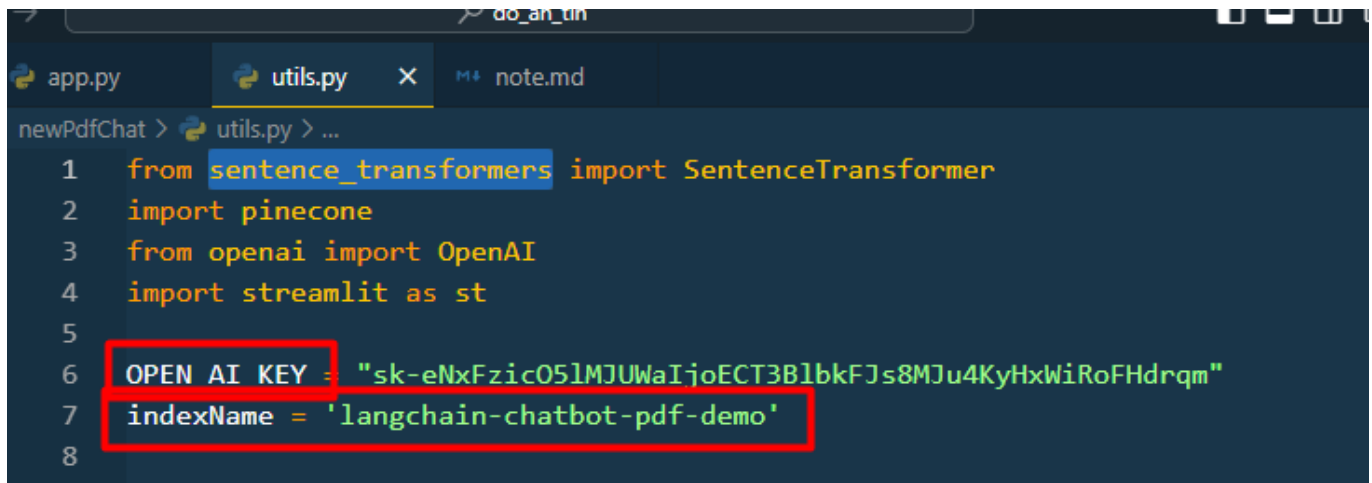
- used the following command

```

pip install langchain langchain-openai openai tiktoken pinecone-client streamlit
sentence-transformers

```

Variables to be notified:



```
app.py  utils.py  X  note.md
newPdfChat > utils.py > ...
1  from sentence_transformers import SentenceTransformer
2  import pinecone
3  from openai import OpenAI
4  import streamlit as st
5
6  OPEN AI KEY = "sk-eNxFzic051MJUWaIjoECT3B1bkFJs8MJu4KyHxWiRoFHdrqm"
7  indexName = 'langchain-chatbot-pdf-demo'
8
```

- index_name is the name of the database in Pinecone, this one must be correct or else the code can't run
- Same with openapi key.