

## ✓ Install Libraries

```

1 pip install sentence_transformers -q

1
2 !pip install --quiet bitsandbytes==0.41.1

1 pip install accelerate -q

1 pip install langchain_openai langchain_community langchain_experimental -q

1 pip install langchain huggingface_hub -q

1 pip install unstructured -q

1 pip install pydantic==2.6.0 -q

1 pip install lark chromadb faiss-gpu -q

1 !pip install optimum -q

1 !nvidia-smi

```

Tue Apr 16 17:59:22 2024

NVIDIA-SMI 535.104.05				Driver Version: 535.104.05		CUDA Version: 12.2	
GPU	Name		Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr. ECC
Fan	Temp	Perf	Pwr:Usage/Cap		Memory-Usage	GPU-Util	Compute M. MIG M.
0	Tesla V100-SXM2-16GB		Off	00000000:00:04.0	Off		0
N/A	34C	P0	25W / 300W		0MiB / 16384MiB	0%	Default N/A

-----							
Processes:							
GPU	GI	CI	PID	Type	Process name	GPU Memory	
	ID	ID				Usage	
=====							
No running processes found							
-----							

## ✓ Import Libraries

```

1 import torch
2 from torch import cuda, bfloat16
3 import transformers
4 import pandas as pd
5
6 import sentence_transformers
7 from langchain_community.embeddings import HuggingFaceEmbeddings
8 from langchain.prompts import ChatPromptTemplate
9 from langchain.vectorstores import FAISS
10 from langchain_community.llms import HuggingFaceHub
11 from langchain_community.llms.huggingface_pipeline import HuggingFacePipeline
12 from langchain.prompts import PromptTemplate
13 from langchain_core.runnables import RunnableLambda, RunnablePassthrough
14 from langchain_core.output_parsers import StrOutputParser

1 device = f'cuda:{cuda.current_device()}' if cuda.is_available() else 'cpu'
2 device

```

```
'cuda:0'
```

## 1. Read Data

```
1 path = "/content/drive/MyDrive/DATAGPT/take_home_dataset.csv"

1 df = pd.read_csv(path, delimiter=";")

1 columns_to_use = ['Date', 'Order_ID', 'Product_Category', 'Delivery_distance']
2 df = df[columns_to_use]
3 df.reset_index(inplace=True)
4 df.rename(columns={'index': 'Row_ID'}, inplace=True)

1 df.head()
```

	Row_ID	Date	Order_ID	Product_Category	Delivery_distance
0	0	2023-07-01	3808	Apparel	300
1	1	2023-07-02	3808	Apparel	300
2	2	2023-07-03	3808	Apparel	300
3	3	2023-07-04	3808	Apparel	300
4	4	2023-07-01	487	Cosmetics & Personal Care	500

Next steps: [View recommended plots](#)

## Preprocessing

```
1 metadata = {
2     'Row_ID' : 'numeric',
3     'Date' : "datetime",
4     'Order_ID' : 'numeric',
5     'Product_Category' : 'categorical',
6     'Delivery_distance' : 'numeric'
7 }
```

## Create Embedding Model

```
1 embed_model_id = 'sentence-transformers/all-MiniLM-L6-v2'
2
3 embed_model = HuggingFaceEmbeddings(
4     model_name=embed_model_id,
5     model_kwargs={'device': device},
6     encode_kwargs={'device': device, 'batch_size': 32}
7 )
```

/usr/local/lib/python3.10/dist-packages/huggingface\_hub/utils/\_token.py:88: UserWarning:  
The secret `HF\_TOKEN` does not exist in your Colab secrets.  
To authenticate with the Hugging Face Hub, create a token in your settings tab (<https://huggingface.co/settings/tokens>), set it as secret.  
You will be able to reuse this secret in all of your notebooks.  
Please note that authentication is recommended but still optional to access public models or datasets.  
warnings.warn()

## 3. Embedding Model

- create column embeddings
- find column\_name from query
- get column type from metadata
- if type == categorical: get value from column similar to query

- elif type in [datetime, numeric]: extract exact value from query

## ✓ Column Embeddings to get column name

```

1 column_embeddings = FAISS.from_texts(columns_to_use, embed_model)

1 def get_column_name(query, embeddings):
2     return embeddings.similarity_search(query, k = 1)[0].page_content

1 def get_column_type(column, metadata):
2     return metadata[column]

1 def extract_similar_value(query, data, column, embed_model):
2     '''
3     This method is executed when the column is categorical
4     '''
5     col_values_embedding = FAISS.from_texts(data[column].unique().astype(str), embed_model)
6     extracted_value = col_values_embedding.similarity_search(query, k = 1)[0].page_content
7     return extracted_value

```

## ✓ Rephrase query if column is categorical

```

1 def get_final_query(query, embed_model, column_embeddings, data, metadata):
2     similar_column = get_column_name(query, column_embeddings)
3     column_type = get_column_type(similar_column, metadata)
4     if column_type == 'categorical':
5         extracted_column_value = extract_similar_value(query, data, similar_column, embed_model)
6         final_query = similar_column + " " + extracted_column_value
7         return final_query
8     else:
9         final_query = query
10    return final_query

```

## ✓ Model

```

1 from torch import cuda, bfloat16
2 import transformers
3
4 model_id = 'meta-llama/Llama-2-13b-chat-hf'
5
6
7 # set quantization configuration to load large model with less GPU memory
8 # this requires the `bitsandbytes` library
9 bnb_config = transformers.BitsAndBytesConfig(
10     bnb_4bit_quant_type='nf4',
11     load_in_8bit_fp32_cpu_offload=True,
12     bnb_4bit_use_double_quant=True,
13     bnb_4bit_compute_dtype=bfloat16
14 )
15
16 # begin initializing HF items, need auth token for these
17 hf_auth = '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'
18 model_config = transformers.AutoConfig.from_pretrained(
19     model_id,
20     use_auth_token=hf_auth
21 )
22
23 model = transformers.AutoModelForCausalLM.from_pretrained(
24     model_id,
25     trust_remote_code=True,
26     config=model_config,
27     quantization_config=bnb_config,
28     device_map='auto',
29     use_auth_token=hf_auth
30 )
31 model.eval()
32 print(f"Model loaded on {device}")
33
34 # tokenizer
35 tokenizer = transformers.AutoTokenizer.from_pretrained(
36     model_id,
37     use_auth_token=hf_auth
38 )
39
40 # pipeline
41 llm_pipeline = transformers.pipeline(
42     model=model, tokenizer=tokenizer,
43     return_full_text=True, # langchain expects the full text
44     task='text-generation',
45     # we pass model parameters here too
46     temperature=0.1, # 'randomness' of outputs, 0.0 is the min and 1.0 the max
47     max_new_tokens=1024, # max number of tokens to generate in the output
48     repetition_penalty=1.1 # without this output begins repeating
49 )
50

```

/usr/local/lib/python3.10/dist-packages/transformers/models/auto/configuration\_auto.py:1  
 warnings.warn(  
 /usr/local/lib/python3.10/dist-packages/transformers/models/auto/auto\_factory.py:466: Fu  
 warnings.warn(  
 Loading checkpoint shards: 100% 3/3 [02:01<00:00, 39.01s/it]  
 You are calling `save\_pretrained` to a 4-bit converted model, but your `bitsandbytes` ve  
 /usr/local/lib/python3.10/dist-packages/transformers/models/auto/tokenization\_auto.py:72  
 warnings.warn(  
 Model loaded on cuda:0

## ✓ Parser Logic

```
1 llm = HuggingFacePipeline(pipeline=llm_pipeline)
```

```

1 pandas_prompt_template = """<<SYS>>
2   You are a data analyst working on a pandas dataframe 'df'.\n Use df.head for your reference.Your job is to return a Pandas expression
3   [INST]
4   Instructions: {instructions}
5   df: {df_head}
6   Query: {query}
7   PRINT ONLY THE PANDAS EXPRESSION after Expression--> AT THE END OF RESPONSE.
8   [/INST]
9   END
10  """

1 extraction_prompt = """<<SYS>>
2   You are a data analyst working on a pandas dataframe 'df'.\nYour job is to extract 'column_name' and 'value' from the pandas expressio
3   [INST]
4   PRINT ONLY THE `column_name` AND `value` AFTER SOLUTION--> at the end of response.
5   expression: {expression}
6   [/INST]
7   END
8  """

1 instructions_prompt = """
2 1. RETURN THE PANDAS EXPRESSION ONLY.
3 Here are the possible columns:
4 ```
5 [{columns}]
6 ```
7 """

1 prompt = ChatPromptTemplate.from_template(pandas_prompt_template)

1 extraction_prompt = ChatPromptTemplate.from_template(extraction_prompt)

1 final_prompt = prompt.partial(instructions = instructions_prompt)

1 chain = (
2     {"query": RunnablePassthrough(), "columns": lambda x: list(df.columns), "df_head": lambda x: df.head(5).to_dict()}
3     | final_prompt
4     | llm
5     | StrOutputParser()
6 )

1 extraction_chain = (
2     {"expression": RunnablePassthrough()}
3     | extraction_prompt
4     | llm
5     | StrOutputParser()
6 )

```

## ✓ fetch query from response

```

1 def extract_query_from_response(answer):
2     pandas_query = answer.split("\n")[-1]
3     return pandas_query

```

## ✓ Get row indices satisfying query

```

1 def get_row_id(df, answer):
2     pandas_query = extract_query_from_response(answer)
3     result_df = eval(pandas_query)
4     if isinstance(result_df, pd.DataFrame):
5         indices = result_df.index.tolist()
6     else:
7         indices = df.index[result_df].tolist()
8     return indices

```

## ✓ Method to extract column\_name and value

```

1 def extract_col_and_value(answer):
2     pandas_query = extract_query_from_response(answer)
3     extraction_chain_response = extraction_chain.invoke(pandas_query)
4     column_name = extraction_chain_response.split("END")[1].strip().split("column_name:")[1].split("\n")[0].strip()
5     value = extraction_chain_response.split("END")[1].strip().split("value:")[1].split("\n")[0].strip()
6     return column_name, value

1 def generate_output(answer, df):
2     pandas_query = extract_query_from_response(answer)
3     column_name, value = extract_col_and_value(answer)
4     indices = get_row_id(df, answer)
5     return {"column_name": column_name, "value": value, "row_ids" : indices}

```

## ✓ Runner

```

1 def run(query):
2     final_query = get_final_query(query, embed_model, column_embeddings, df, metadata)
3     answer = chain.invoke(final_query)
4     pandas_query = extract_query_from_response(answer)
5     indices = get_row_id(df, answer)
6     output = generate_output(answer, df)
7     print(output)
8

1 run("delivery distance greater than 500")

{'column_name': 'Delivery_distance', 'value': '> 500', 'row_ids': [14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99]}

1 run("product apparel")

{'column_name': 'Product_Category', 'value': 'Apparel', 'row_ids': [0, 1, 2, 3, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99]}

1 run("product machine")

{'column_name': 'Product_Category', 'value': 'Electronics', 'row_ids': [89, 90, 91, 92, 93, 94, 95, 96]}

1 run("product toys")

{'column_name': 'Product_Category', 'value': 'Toys & Games', 'row_ids': [39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99]}

1
2 run("orders between 1 july 2023 and 5 july 2023")

/usr/local/lib/python3.10/dist-packages/transformers/pipelines/base.py:1157: UserWarning: You seem to be using the pipelines sequentially in a loop
warnings.warn(
/usr/local/lib/python3.10/dist-packages/transformers/pipelines/base.py:1157: UserWarning: You seem to be using the pipelines sequentially in a loop
warnings.warn(
{'column_name': 'Date', 'value': '2023-07-01', 'row_ids': [1, 2, 3, 5, 6, 7, 10, 11, 12, 15, 16, 17, 29, 30, 31, 34, 35, 36, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99]}

1 run("order id more than 4000")

/usr/local/lib/python3.10/dist-packages/transformers/pipelines/base.py:1157: UserWarning: You seem to be using the pipelines sequentially in a loop
warnings.warn(
/usr/local/lib/python3.10/dist-packages/transformers/pipelines/base.py:1157: UserWarning: You seem to be using the pipelines sequentially in a loop
warnings.warn(
{'column_name': 'Order_ID', 'value': '4001', 'row_ids': [9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99]}

1 run("delivery distance greater than 500")

/usr/local/lib/python3.10/dist-packages/transformers/pipelines/base.py:1157: UserWarning: You seem to be using the pipelines sequentially in a loop
warnings.warn(
/usr/local/lib/python3.10/dist-packages/transformers/pipelines/base.py:1157: UserWarning: You seem to be using the pipelines sequentially in a loop
warnings.warn(

```

```
warnings.warn(  
{ 'column_name': 'Delivery_distance', 'value': '> 500', 'row_ids': [14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 3
```

1