

Restaurants receive thousands of customer reviews, but star ratings alone fail to explain why ratings increase or decline. Reviews often contain rich information about food quality, service, pricing, ambience, and operational factors, yet this information remains unstructured and difficult to analyze at scale. This limits restaurants' ability to identify the drivers of customer satisfaction and dissatisfaction.

Notebook 2: Sentiment Prediction and Aspect Extraction

This notebook applies previously trained and saved models to perform large-scale sentiment prediction and topic (aspect) extraction on Yelp restaurant reviews.

Specifically, a fine-tuned transformer-based sentiment classification model is used to predict customer sentiment from review text, while a BERTopic-based topic model identifies the key aspects customers discuss in their reviews. The outputs from both models are combined with business identifiers, review text, star ratings, and restaurant operational attributes to construct a unified analytical dataset (aspect_df).

This dataset serves as the foundation for downstream analysis aimed at explaining why restaurant ratings vary beyond star scores alone.

```
!pip install pyspark

Requirement already satisfied: pyspark in /usr/local/lib/python3.12/dist-packages (4.0.1)
Requirement already satisfied: py4j==0.10.9.9 in /usr/local/lib/python3.12/dist-packages (from pyspark) (0.10.9.9)
```

```
import pyspark
print("PySpark installed and imported successfully!")
```

PySpark installed and imported successfully!

```
from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

```
from pyspark.sql import SparkSession
```

```
spark = SparkSession.builder \
    .appName("YelpAnalysis") \
    .config("spark.driver.memory", "8g") \
    .getOrCreate()
```

Loading the data

```
business_path = "/content/drive/MyDrive/yelp_dataset/yelp_academic_dataset_business.json"

business_df = spark.read.json(business_path)
business_df.printSchema()
```

```
root
 |-- address: string (nullable = true)
 |-- attributes: struct (nullable = true)
 |   |-- AcceptsInsurance: string (nullable = true)
 |   |-- AgesAllowed: string (nullable = true)
 |   |-- Alcohol: string (nullable = true)
 |   |-- Ambience: string (nullable = true)
 |   |-- BYOB: string (nullable = true)
 |   |-- BYOBParking: string (nullable = true)
 |   |-- BestNights: string (nullable = true)
 |   |-- BikeParking: string (nullable = true)
 |   |-- BusinessAcceptsBitcoin: string (nullable = true)
 |   |-- BusinessAcceptsCreditCards: string (nullable = true)
 |   |-- BusinessParking: string (nullable = true)
 |   |-- ByAppointmentOnly: string (nullable = true)
 |   |-- Caters: string (nullable = true)
 |   |-- CoatCheck: string (nullable = true)
 |   |-- Corkage: string (nullable = true)
 |   |-- DietaryRestrictions: string (nullable = true)
 |   |-- DogsAllowed: string (nullable = true)
 |   |-- DriveThru: string (nullable = true)
 |   |-- GoodForDancing: string (nullable = true)
 |   |-- GoodForKids: string (nullable = true)
 |   |-- GoodForMeal: string (nullable = true)
 |   |-- HairSpecializesIn: string (nullable = true)
 |   |-- HappyHour: string (nullable = true)
 |   |-- HasTV: string (nullable = true)
```

```

| -- Music: string (nullable = true)
| -- NoiseLevel: string (nullable = true)
| -- Open24Hours: string (nullable = true)
| -- OutdoorSeating: string (nullable = true)
| -- RestaurantsAttire: string (nullable = true)
| -- RestaurantsCounterService: string (nullable = true)
| -- RestaurantsDelivery: string (nullable = true)
| -- RestaurantsGoodForGroups: string (nullable = true)
| -- RestaurantsPriceRange2: string (nullable = true)
| -- RestaurantsReservations: string (nullable = true)
| -- RestaurantsTableService: string (nullable = true)
| -- RestaurantsTakeOut: string (nullable = true)
| -- Smoking: string (nullable = true)
| -- WheelchairAccessible: string (nullable = true)
| -- WiFi: string (nullable = true)
|-- business_id: string (nullable = true)
|-- categories: string (nullable = true)
|-- city: string (nullable = true)
|-- hours: struct (nullable = true)
|   |-- Friday: string (nullable = true)
|   |-- Monday: string (nullable = true)
|   |-- Saturday: string (nullable = true)
|   |-- Sunday: string (nullable = true)
|   |-- Thursday: string (nullable = true)
|   |-- Tuesday: string (nullable = true)
|   |-- Wednesday: string (nullable = true)
|-- is_open: long (nullable = true)
|-- latitude: double (nullable = true)
|-- longitude: double (nullable = true)
|-- name: string (nullable = true)

```

```
review_path = "/content/drive/MyDrive/yelp_dataset/yelp_academic_dataset_review.json"
```

```
reviews_df = spark.read.json(review_path)
reviews_df.printSchema()
```

```

root
|-- business_id: string (nullable = true)
|-- cool: long (nullable = true)
|-- date: string (nullable = true)
|-- funny: long (nullable = true)
|-- review_id: string (nullable = true)
|-- stars: double (nullable = true)
|-- text: string (nullable = true)
|-- useful: long (nullable = true)
|-- user_id: string (nullable = true)

```

Cleand and Filter Data

```

#mandatory
from pyspark.sql.functions import col

business_df_clean = business_df.filter(
    col("business_id").isNotNull() &
    col("categories").isNotNull()
)

```

```

#Only need restaurant data
restaurants_df = business_df_clean.filter(
    col("categories").contains("Restaurants")
)

```

```

#mandatory
reviews_df_clean = reviews_df.filter(
    col("review_id").isNotNull() &
    col("business_id").isNotNull() &
    col("text").isNotNull() &
    col("stars").isNotNull()
)

```

```

#Sample Review
reviews_sample = reviews_df_clean.sample(fraction=0.05, seed=42)

```

```

#Only restaurent reviews are required
restaurant_reviews_df = reviews_sample.join(
    restaurants_df.select("business_id"),
    on="business_id",

```

```
    how="inner"  
)
```

Select Useful Fields for Sentiment Analysis

```
sentiment_df = restaurant_reviews_df.select("review_id", "business_id", "text", "stars")
```

```
from pyspark.sql.functions import col, when

sentiment_df = sentiment_df.withColumn(
    "label",
    when(col("stars") <= 2, 0) # Negative
    .when(col("stars") == 3, 1) # Neutral
    .otherwise(2) # Positive
)
```

```
!pip install transformers datasets accelerate
```

```
Requirement already satisfied: transformers in /usr/local/lib/python3.12/dist-packages (4.57.3)
Requirement already satisfied: datasets in /usr/local/lib/python3.12/dist-packages (4.0.0)
Requirement already satisfied: accelerate in /usr/local/lib/python3.12/dist-packages (1.12.0)
Requirement already satisfied: filelock in /usr/local/lib/python3.12/dist-packages (from transformers) (3.20.0)
Requirement already satisfied: huggingface-hub<1.0,>=0.34.0 in /usr/local/lib/python3.12/dist-packages (from transformers)
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.12/dist-packages (from transformers) (2.0.2)
Requirement already satisfied: packaging=20.0 in /usr/local/lib/python3.12/dist-packages (from transformers) (25.0)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.12/dist-packages (from transformers) (6.0.3)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.12/dist-packages (from transformers) (2025.11.3)
Requirement already satisfied: requests in /usr/local/lib/python3.12/dist-packages (from transformers) (2.32.4)
Requirement already satisfied: tokenizers<=0.23.0,>=0.22.0 in /usr/local/lib/python3.12/dist-packages (from transformers)
Requirement already satisfied: safetensors>=0.4.3 in /usr/local/lib/python3.12/dist-packages (from transformers) (0.7.0)
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.12/dist-packages (from transformers) (4.67.1)
Requirement already satisfied: pyarrow>=15.0.0 in /usr/local/lib/python3.12/dist-packages (from datasets) (18.1.0)
Requirement already satisfied: dill<0.3.9,>=0.3.0 in /usr/local/lib/python3.12/dist-packages (from datasets) (0.3.8)
Requirement already satisfied: pandas in /usr/local/lib/python3.12/dist-packages (from datasets) (2.2.2)
Requirement already satisfied: xxhash in /usr/local/lib/python3.12/dist-packages (from datasets) (3.6.0)
Requirement already satisfied: multiprocessing<0.70.17 in /usr/local/lib/python3.12/dist-packages (from datasets) (0.70.16)
Requirement already satisfied: fsspec<=2025.3.0,>=2023.1.0 in /usr/local/lib/python3.12/dist-packages (from fsspec[http]<=
Requirement already satisfied: psutil in /usr/local/lib/python3.12/dist-packages (from accelerate) (5.9.5)
Requirement already satisfied: torch>=2.0.0 in /usr/local/lib/python3.12/dist-packages (from accelerate) (2.9.0+cu126)
Requirement already satisfied: aiohttp!=4.0.0a0,!=4.0.0a1 in /usr/local/lib/python3.12/dist-packages (from fsspec[http]<=
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.12/dist-packages (from huggingface-hub)
Requirement already satisfied: hf-xet<2.0.0,>=1.1.3 in /usr/local/lib/python3.12/dist-packages (from huggingface-hub<1.0,>
Requirement already satisfied: charset_normalizer<4,>=2 in /usr/local/lib/python3.12/dist-packages (from requests->transfo
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.12/dist-packages (from requests->transformers) (3.11
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.12/dist-packages (from requests->transformers)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.12/dist-packages (from requests->transformers)
Requirement already satisfied: setuptools in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0->accelerate) (75.2
Requirement already satisfied: sympy>=1.13.3 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0->accelerate) (1
Requirement already satisfied: networkx>=2.5.1 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0->accelerate)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0->accelerate) (3.1.6)
Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.6.77 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0
Requirement already satisfied: nvidia-cuda-runtime-cu12==12.6.77 in /usr/local/lib/python3.12/dist-packages (from torch>=2
Requirement already satisfied: nvidia-cuda-cupti-cu12==12.6.80 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0
Requirement already satisfied: nvidia-cudnn-cu12==9.10.2.21 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0-
Requirement already satisfied: nvidia-cUBLAS-cu12==12.6.4.1 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0-
Requirement already satisfied: nvidia-cufft-cu12==11.3.0.4 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0->
Requirement already satisfied: nvidia-curand-cu12==10.3.7.77 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0
Requirement already satisfied: nvidia-cusolver-cu12==11.7.1.2 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.
Requirement already satisfied: nvidia-cusparse-cu12==12.5.4.2 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.
Requirement already satisfied: nvidia-cusparseelt-cu12==0.7.1 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0
Requirement already satisfied: nvidia-nccl-cu12==2.27.5 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0->acc
Requirement already satisfied: nvidia-nvshmem-cu12==3.3.20 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0->
Requirement already satisfied: nvidia-nvtx-cu12==12.6.77 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0->ac
Requirement already satisfied: nvidia-nvjitlink-cu12==12.6.85 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.
Requirement already satisfied: nvidia-cufile-cu12==1.11.1.6 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0-
Requirement already satisfied: triton==3.5.0 in /usr/local/lib/python3.12/dist-packages (from torch>=2.0.0->accelerate) (3
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.12/dist-packages (from pandas->datasets) (
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.12/dist-packages (from pandas->datasets) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.12/dist-packages (from pandas->datasets) (2025.3)
Requirement already satisfied: aiohttpyeayeballs>=2.5.0 in /usr/local/lib/python3.12/dist-packages (from aiohttp!=4.0.0a0,!_
Requirement already satisfied: aiosignal>=1.4.0 in /usr/local/lib/python3.12/dist-packages (from aiohttp!=4.0.0a0,!=4.0.0a
Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.12/dist-packages (from aiohttp!=4.0.0a0,!=4.0.0a1->
Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.12/dist-packages (from aiohttp!=4.0.0a0,!=4.0.0
Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.12/dist-packages (from aiohttp!=4.0.0a0,!=4.0
Requirement already satisfied: propcache>=0.2.0 in /usr/local/lib/python3.12/dist-packages (from aiohttp!=4.0.0a0,!=4.0.0a
```

Load tokenizer

```
from transformers import AutoTokenizer

MODEL_DIR = "/content/drive/MyDrive/restaurant_sentiment_model"
tokenizer = AutoTokenizer.from_pretrained(MODEL_DIR)
```

Load Model

```
from transformers import AutoModelForSequenceClassification

model = AutoModelForSequenceClassification.from_pretrained(MODEL_DIR)
model.eval()

DistilBertForSequenceClassification(
    (distilbert): DistilBertModel(
        (embeddings): Embeddings(
            (word_embeddings): Embedding(30522, 768, padding_idx=0)
            (position_embeddings): Embedding(512, 768)
            (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
            (dropout): Dropout(p=0.1, inplace=False)
        )
        (transformer): Transformer(
            (layer): ModuleList(
                (0-5): 6 x TransformerBlock(
                    (attention): DistilBertSdpAttention(
                        (dropout): Dropout(p=0.1, inplace=False)
                        (q_lin): Linear(in_features=768, out_features=768, bias=True)
                        (k_lin): Linear(in_features=768, out_features=768, bias=True)
                        (v_lin): Linear(in_features=768, out_features=768, bias=True)
                        (out_lin): Linear(in_features=768, out_features=768, bias=True)
                    )
                    (sa_layer_norm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
                    (ffn): FFN(
                        (dropout): Dropout(p=0.1, inplace=False)
                        (lin1): Linear(in_features=768, out_features=3072, bias=True)
                        (lin2): Linear(in_features=3072, out_features=768, bias=True)
                        (activation): GELUActivation()
                    )
                    (output_layer_norm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
                )
            )
        )
    )
    (pre_classifier): Linear(in_features=768, out_features=768, bias=True)
    (classifier): Linear(in_features=768, out_features=3, bias=True)
    (dropout): Dropout(p=0.2, inplace=False)
)
)
```

sentiment prediction function

```
import torch
import numpy as np

def predict_sentiment(texts, batch_size=16):
    all_preds = []

    for i in range(0, len(texts), batch_size):
        batch = texts[i:i+batch_size]
        inputs = tokenizer(
            batch,
            padding=True,
            truncation=True,
            max_length=128,
            return_tensors="pt"
        )

        with torch.no_grad():
            outputs = model(**inputs)
            preds = torch.argmax(outputs.logits, dim=1)

        all_preds.extend(preds.cpu().numpy())

    return np.array(all_preds)
```

SENTIMENT + ASPECT EXTRACTION

Get reviews + business_id

```
reviews_pdf = sentiment_df.select("business_id", "text").toPandas()[:10000]
reviews_texts = reviews_pdf["text"].tolist()
```

```
pip install -U bertopic
```

Collecting bertopic

 Downloading bertopic-0.17.4-py3-none-any.whl.metadata (24 kB)

Requirement already satisfied: hdbscan>=0.8.29 in /usr/local/lib/python3.12/dist-packages (from bertopic) (0.8.41)

Requirement already satisfied: umap-learn>=0.5.0 in /usr/local/lib/python3.12/dist-packages (from bertopic) (0.5.9.post2)

Requirement already satisfied: numpy>=1.20.0 in /usr/local/lib/python3.12/dist-packages (from bertopic) (2.0.2)

Requirement already satisfied: pandas>=1.1.5 in /usr/local/lib/python3.12/dist-packages (from bertopic) (2.2.2)

Requirement already satisfied: plotly>=4.7.0 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.24.1)

Requirement already satisfied: scikit-learn>=1.0 in /usr/local/lib/python3.12/dist-packages (from bertopic) (1.6.1)

Requirement already satisfied: sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0)

Requirement already satisfied: tqdm>=4.41.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (4.67.1)

Requirement already satisfied: llvmlite>0.36.0 in /usr/local/lib/python3.12/dist-packages (from bertopic) (0.43.0)

Requirement already satisfied: scipy>=1.0 in /usr/local/lib/python3.12/dist-packages (from hdbscan>=0.8.29->bertopic) (1.1.0)

Requirement already satisfied: joblib>=1.0 in /usr/local/lib/python3.12/dist-packages (from hdbscan>=0.8.29->bertopic) (1.1.0)

Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.12/dist-packages (from pandas>=1.1.5->bertopic) (2.3.1)

Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.12/dist-packages (from pandas>=1.1.5->bertopic) (2023.1)

Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.12/dist-packages (from pandas>=1.1.5->bertopic) (2023.1)

Requirement already satisfied: tenacity>=6.2.0 in /usr/local/lib/python3.12/dist-packages (from plotly>=4.7.0->bertopic) (1.1.0)

Requirement already satisfied: packaging in /usr/local/lib/python3.12/dist-packages (from plotly>=4.7.0->bertopic) (25.0)

Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.12/dist-packages (from scikit-learn>=1.0->bertopic) (2023.1)

Requirement already satisfied: transformers<6.0.0,>=4.41.0 in /usr/local/lib/python3.12/dist-packages (from sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: torch>=1.11.0 in /usr/local/lib/python3.12/dist-packages (from sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: huggingface-hub>=0.20.0 in /usr/local/lib/python3.12/dist-packages (from sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: typing_extensions>=4.5.0 in /usr/local/lib/python3.12/dist-packages (from sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: numba>=0.51.2 in /usr/local/lib/python3.12/dist-packages (from umap-learn>=0.5.0->bertopic) (0.5.1)

Requirement already satisfied: pynndescent>=0.5 in /usr/local/lib/python3.12/dist-packages (from umap-learn>=0.5.0->bertopic) (0.5.1)

Requirement already satisfied: filelock in /usr/local/lib/python3.12/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.12/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.12/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: requests in /usr/local/lib/python3.12/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: hf-xet<2.0.0,>=1.1.3 in /usr/local/lib/python3.12/dist-packages (from huggingface-hub>=0.20.0->sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.12/dist-packages (from python-dateutil>=2.8.2->pandas>=1.1.0 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: setuptools in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: sympy>=1.13.3 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: networkx>=2.5.1 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: jinja2 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.6.77 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cuda-runtime-cu12==12.6.77 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cuda-cupti-cu12==12.6.80 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cudnn-cu12==9.10.2.21 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cUBLAS-cu12==12.6.4.1 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cufft-cu12==11.3.0.4 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-curand-cu12==10.3.7.77 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cusolver-cu12==11.7.1.2 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cusparse-cu12==12.5.4.2 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cusparseelt-cu12==0.7.1 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-nccl-cu12==2.27.5 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-nvshmem-cu12==3.3.20 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-nvtx-cu12==12.6.77 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-nvjitlink-cu12==12.6.85 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: nvidia-cufile-cu12==1.11.1.6 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: triton==3.5.0 in /usr/local/lib/python3.12/dist-packages (from torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.12/dist-packages (from transformers<6.0.0,>=4.4 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: tokenizers<=0.23.0,>=0.22.0 in /usr/local/lib/python3.12/dist-packages (from transformers<6.0.0,>=4.4 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: safetensors>=0.4.3 in /usr/local/lib/python3.12/dist-packages (from transformers<6.0.0,>=4.4 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.12/dist-packages (from sympy>=1.13.3->torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.12/dist-packages (from jinja2->torch>=1.11.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Requirement already satisfied: charset_normalizer<4,>=2 in /usr/local/lib/python3.12/dist-packages (from requests->huggingface-hub>=0.20.0->sentence-transformer>=0.4.1 in /usr/local/lib/python3.12/dist-packages (from bertopic) (5.1.0))

Load aspect extraction model

```
from bertopic import BERTopic
from sentence_transformers import SentenceTransformer

embedding_model = SentenceTransformer("all-MiniLM-L6-v2")
MODEL_PATH = "/content/drive/MyDrive/restaurant_sentiment_model/aspect"
topic_model = BERTopic.load(
    MODEL_PATH,
    embedding_model=embedding_model
)
```

Obtain the topics from reviews

```
topics, probs = topic_model.transform(reviews_texts)
```

```
Batches: 100%                                         313/313 [00:14<00:00, 61.99it/s]
2025-12-25 11:40:52,777 - BERTopic - Dimensionality - Reducing dimensionality of input embeddings.
2025-12-25 11:41:01,803 - BERTopic - Dimensionality - Completed ✓
2025-12-25 11:41:01,804 - BERTopic - Clustering - Approximating new points with `hdbscan_model`
2025-12-25 11:41:02,429 - BERTopic - Cluster - Completed ✓
sentiment_labels = predict_sentiment(reviews_texts)
```

Combine aspect extraction topics with sentiment prediction labels

```
import pandas as pd

aspect_df = pd.DataFrame({
    "business_id": reviews_pdf["business_id"],
    "review": reviews_texts,
    "topic": topics,
    "sentiment": sentiment_labels
})
```

Merge business attributes

```
from pyspark.sql.functions import col, coalesce, lit

##Flatten attributes dynamically
attribute_cols = business_df.select("attributes.*").columns

businessAttrs = business_df.select(
    "business_id",
    *[col(f"attributes.{c}").alias(c) for c in attribute_cols]
)

##Replace nulls with "Unknown"
for c in attribute_cols:
    businessAttrs = businessAttrs.withColumn(
        c, coalesce(col(c), lit("Unknown")))
)

##Convert to Pandas
businessAttrs_pd = businessAttrs.toPandas()

##Merge with aspect_df
aspect_df = aspect_df.merge(
    businessAttrs_pd,
    on="business_id",
    how="left"
)
```

```
import pickle

PATH = "/content/drive/MyDrive/restaurant_sentiment_model/aspect_df.pkl"

##Save
aspect_df.to_pickle(PATH)
```