

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

Task

Analyze the sentiment of all text files contained within the `/content/Sneha job.zip` archive. Present the key findings from this sentiment analysis, including a visualization of the overall sentiment distribution and highlighting any significant insights or patterns observed across the files.

Extract Archive

Subtask:

Extract all files from the `/content/Sneha job.zip` archive into a temporary directory.

Reasoning: I need to import the `zipfile` and `os` modules, create a directory, and then extract the contents of the zip file into that directory.

```
import zipfile
import os

# Define the path to the zip file and the extraction directory
zip_file_path = '/content/Sneha job.zip'
extraction_dir = '/content/extracted_files'

# Create the extraction directory if it doesn't exist
os.makedirs(extraction_dir, exist_ok=True)

# Open the zip file and extract its contents
with zipfile.ZipFile(zip_file_path, 'r') as zip_ref:
    zip_ref.extractall(extraction_dir)

print(f"All files extracted to: {extraction_dir}")

All files extracted to: /content/extracted_files
```

Identify Text Files

Subtask:

Scan the extracted directory to identify all text-based files suitable for sentiment analysis. This might involve checking file extensions or content types.

Reasoning: I need to iterate through the extracted directory, identify files with common text extensions, and store their paths for subsequent sentiment analysis.

```

text_file_paths = []
text_extensions = ['.txt', '.csv', '.md', '.log', '.json', '.xml',
'.html', '.css', '.js', '.py', '.java', '.c', '.cpp', '.h', '.hpp'] #
Common text file extensions

for root, _, files in os.walk(extraction_dir):
    for file in files:
        file_path = os.path.join(root, file)
        if any(file.lower().endswith(ext) for ext in text_extensions):
            text_file_paths.append(file_path)

print(f"Found {len(text_file_paths)} text files:")
for path in text_file_paths:
    print(path)

Found 0 text files:

```

Reasoning: The previous code failed to identify any text files because the `.tsv` extension was missing from the `text_extensions` list, and there are `.tsv` files present in the extracted directory that are relevant for sentiment analysis. I need to update the `text_extensions` list to include `.tsv` to correctly identify these files.

```

text_file_paths = []
text_extensions = ['.txt', '.csv', '.tsv', '.md', '.log', '.json',
'.xml', '.html', '.css', '.js', '.py', '.java', '.c', '.cpp', '.h',
'.hpp'] # Common text file extensions including .tsv

for root, _, files in os.walk(extraction_dir):
    for file in files:
        file_path = os.path.join(root, file)
        if any(file.lower().endswith(ext) for ext in text_extensions):
            text_file_paths.append(file_path)

print(f"Found {len(text_file_paths)} text files:")
for path in text_file_paths:
    print(path)

Found 2 text files:
/content/extracted_files/a2_RestaurantReviews_FreshDump.tsv
/content/extracted_files/a1_RestaurantReviews_HistoricDump.tsv

```

Load Text Content

Subtask:

Read the content of each identified text file into memory for further processing. Handle potential encoding issues during file reading.

Instructions

1. Iterate through the `text_file_paths` list.
2. For each file, attempt to read its content. Consider using `utf-8` encoding and fallback mechanisms for other encodings if `utf-8` fails.
3. Store the file content, perhaps along with its path, in a suitable data structure (e.g., a dictionary or a list of tuples).

Reasoning: Now I will implement the code to read the content of each identified text file, handling potential encoding issues, and store the file content with its path.

```
file_contents = []

for file_path in text_file_paths:
    content = None
    # Try to read with utf-8 first
    try:
        with open(file_path, 'r', encoding='utf-8') as f:
            content = f.read()
    except UnicodeDecodeError:
        # Fallback to latin-1 if utf-8 fails
        try:
            with open(file_path, 'r', encoding='latin-1') as f:
                content = f.read()
        except Exception as e:
            print(f"Could not read file {file_path} due to encoding error: {e}")
    except Exception as e:
        print(f"Could not read file {file_path} due to error: {e}")

    if content is not None:
        file_contents.append({'path': file_path, 'content': content})

print(f"Successfully loaded content from {len(file_contents)} out of {len(text_file_paths)} files.")
# Display first 500 characters of the first loaded file to verify
if file_contents:
    print("\nFirst 500 characters of the first file:")
    print(file_contents[0]['content'][:500])
```

Successfully loaded content from 2 out of 2 files.

First 500 characters of the first file:

Review

Spend your money elsewhere.

Their regular toasted bread was equally satisfying with the occasional pats of butter... Mmmm...!

The Buffet at Bellagio was far from what I anticipated.

And the drinks are WEAK, people!

-My order was not correct.

Also, I feel like the chips are bought, not made in house.

```
After the disappointing dinner we went elsewhere for dessert.  
The chips and salsa here is amazing!!!!!!!!!!!!!!!!!!!!  
We won't be returning.  
This is my new fav Vegas buffet spot.  
I seriously c
```

Perform Sentiment Analysis

Subtask:

Analyze the sentiment of the loaded text content from each file. For `.tsv` files which are likely tabular and contain reviews, extract the relevant text column for analysis.

Instructions

1. Choose a suitable sentiment analysis library or model (e.g., `TextBlob`, `VADER`, Hugging Face transformers). For this example, we will use `TextBlob` due to its simplicity and effectiveness for general sentiment.
2. Iterate through the `file_contents` list.
3. For each file:
 - a. If the file is a `.tsv` file, parse its content as a DataFrame (using `io.StringIO` and `pd.read_csv` with `sep='\t'`) and extract the 'Review' column or the most relevant text column.
 - b. For other text files, the entire content can be treated as a single document or split into sentences/paragraphs if finer-grained analysis is desired.
 - c. Apply the chosen sentiment analysis tool to determine the polarity and subjectivity of the text.
 - d. Store the sentiment results (e.g., polarity score) along with the file path and, if applicable, the individual reviews/sentences.

Perform Sentiment Analysis

Subtask:

Analyze the sentiment of the loaded text content from each file. For `.tsv` files which are likely tabular and contain reviews, extract the relevant text column for analysis.

Instructions

1. Choose a suitable sentiment analysis library or model (e.g., `TextBlob`, `VADER`, Hugging Face transformers). For this example, we will use `TextBlob` due to its simplicity and effectiveness for general sentiment.
2. Iterate through the `file_contents` list.
3. For each file:
 - a. If the file is a `.tsv` file, parse its content as a DataFrame (using `io.StringIO` and `pd.read_csv` with `sep='\t'`) and extract the 'Review' column or the most relevant text column.

- b. For other text files, the entire content can be treated as a single document or split into sentences/paragraphs if finer-grained analysis is desired.
- c. Apply the chosen sentiment analysis tool to determine the polarity and subjectivity of the text.
- d. Store the sentiment results (e.g., polarity score) along with the file path and, if applicable, the individual reviews/sentences.

Perform Sentiment Analysis

Subtask:

Analyze the sentiment of the loaded text content from each file. For `.tsv` files which are likely tabular and contain reviews, extract the relevant text column for analysis.

Instructions

1. Choose a suitable sentiment analysis library or model (e.g., `TextBlob`, `VADER`, Hugging Face transformers). For this example, we will use `TextBlob` due to its simplicity and effectiveness for general sentiment.
2. Iterate through the `file_contents` list.
3. For each file:
 - a. If the file is a `.tsv` file, parse its content as a DataFrame (using `io.StringIO` and `pd.read_csv` with `sep='\t'`) and extract the 'Review' column or the most relevant text column.
 - b. For other text files, the entire content can be treated as a single document or split into sentences/paragraphs if finer-grained analysis is desired.
 - c. Apply the chosen sentiment analysis tool to determine the polarity and subjectivity of the text.
 - d. Store the sentiment results (e.g., polarity score) along with the file path and, if applicable, the individual reviews/sentences.

Perform Sentiment Analysis

Subtask:

Analyze the sentiment of the loaded text content from each file. For `.tsv` files which are likely tabular and contain reviews, extract the relevant text column for analysis.

Instructions

1. Choose a suitable sentiment analysis library or model (e.g., `TextBlob`, `VADER`, Hugging Face transformers). For this example, we will use `TextBlob` due to its simplicity and effectiveness for general sentiment.
2. Iterate through the `file_contents` list.
3. For each file:
 - a. If the file is a `.tsv` file, parse its content as a DataFrame (using `io.StringIO` and `pd.read_csv` with `sep='\t'`) and extract the 'Review' column or the most relevant text column.

- b. For other text files, the entire content can be treated as a single document or split into sentences/paragraphs if finer-grained analysis is desired.
- c. Apply the chosen sentiment analysis tool to determine the polarity and subjectivity of the text.
- d. Store the sentiment results (e.g., polarity score) along with the file path and, if applicable, the individual reviews/sentences.

Perform Sentiment Analysis

Subtask:

Analyze the sentiment of the loaded text content from each file. For `.tsv` files which are likely tabular and contain reviews, extract the relevant text column for analysis.

Instructions

1. Choose a suitable sentiment analysis library or model (e.g., `TextBlob`, `VADER`, Hugging Face transformers). For this example, we will use `TextBlob` due to its simplicity and effectiveness for general sentiment.
2. Iterate through the `file_contents` list.
3. For each file:
 - a. If the file is a `.tsv` file, parse its content as a `DataFrame` (using `io.StringIO` and `pd.read_csv` with `sep='\t'`) and extract the 'Review' column or the most relevant text column.
 - b. For other text files, the entire content can be treated as a single document or split into sentences/paragraphs if finer-grained analysis is desired.
 - c. Apply the chosen sentiment analysis tool to determine the polarity and subjectivity of the text.
 - d. Store the sentiment results (e.g., polarity score) along with the file path and, if applicable, the individual reviews/sentences.

Reasoning: I need to implement the sentiment analysis using `TextBlob`, processing `.tsv` files specifically by extracting reviews from a `DataFrame`, and storing the results. This requires importing `TextBlob`, `pandas`, and `io`.

```
import pandas as pd
from textblob import TextBlob
import io

sentiment_results = []

for file_info in file_contents:
    file_path = file_info['path']
    content = file_info['content']
    file_name = os.path.basename(file_path)

    if file_name.lower().endswith('.tsv'):
        try:
            # Read TSV content into a DataFrame
```

```

df = pd.read_csv(io.StringIO(content), sep='\t')

# Assuming the review text is in a column named 'Review'
# Check for 'Review' column, or fall back to the first
column if 'Review' is not found
review_column = None
if 'Review' in df.columns:
    review_column = 'Review'
elif len(df.columns) > 0:
    review_column = df.columns[0] # Take the first column
as review if 'Review' is not present

if review_column:
    # Filter out non-string types and apply sentiment
analysis
    reviews = df[review_column].astype(str)
    for i, review in enumerate(reviews):
        analysis = TextBlob(review)
        sentiment_results.append({
            'file_path': file_path,
            'text_type': 'review',
            'text_index': i,
            'text': review,
            'polarity': analysis.sentiment.polarity,
            'subjectivity':
analysis.sentiment.subjectivity
        })
        print(f"Processed {len(reviews)} reviews from
{file_name}")
    else:
        print(f"Warning: No 'Review' column or suitable text
column found in {file_name}. Skipping specific review analysis.")
        # Fallback: Analyze the entire content as a single
document
        analysis = TextBlob(content)
        sentiment_results.append({
            'file_path': file_path,
            'text_type': 'full_document',
            'text_index': -1,
            'text': content[:200], # Store first 200 chars for
brevity
            'polarity': analysis.sentiment.polarity,
            'subjectivity': analysis.sentiment.subjectivity
        })

except Exception as e:
    print(f"Error processing TSV file {file_name}: {e}")
    # If TSV parsing fails, treat the whole content as a
single document
    analysis = TextBlob(content)

```

```

        sentiment_results.append({
            'file_path': file_path,
            'text_type': 'full_document',
            'text_index': -1,
            'text': content[:200],
            'polarity': analysis.sentiment.polarity,
            'subjectivity': analysis.sentiment.subjectivity
        })
    else:
        # For other text files, analyze the entire content as a single
        document
        analysis = TextBlob(content)
        sentiment_results.append({
            'file_path': file_path,
            'text_type': 'full_document',
            'text_index': -1,
            'text': content[:200], # Store first 200 chars for brevity
            'polarity': analysis.sentiment.polarity,
            'subjectivity': analysis.sentiment.subjectivity
        })
    print(f"Processed full document sentiment for {file_name}")

print(f"\nCompleted sentiment analysis for {len(sentiment_results)}
text entries.")
# Display a few sentiment results for verification
if sentiment_results:
    print("\nSample Sentiment Results:")
    for res in sentiment_results[:5]: # Display first 5 results
        print(f"File: {os.path.basename(res['file_path'])}, Type:
{res['text_type']}, Polarity: {res['polarity']:.2f}, Subjectivity:
{res['subjectivity']:.2f}")

Processed 100 reviews from a2_RestaurantReviews_FreshDump.tsv
Processed 900 reviews from a1_RestaurantReviews_HistoricDump.tsv

Completed sentiment analysis for 1000 text entries.

```

Sample Sentiment Results:

```

File: a2_RestaurantReviews_FreshDump.tsv, Type: review, Polarity:
0.00, Subjectivity: 0.00
File: a2_RestaurantReviews_FreshDump.tsv, Type: review, Polarity:
0.17, Subjectivity: 0.40
File: a2_RestaurantReviews_FreshDump.tsv, Type: review, Polarity:
0.10, Subjectivity: 1.00
File: a2_RestaurantReviews_FreshDump.tsv, Type: review, Polarity: -
0.47, Subjectivity: 0.62
File: a2_RestaurantReviews_FreshDump.tsv, Type: review, Polarity:
0.00, Subjectivity: 0.00

```


Reasoning: Now that the sentiment analysis is complete, I will convert the `sentiment_results` list into a pandas DataFrame to facilitate easier manipulation and visualization of the sentiment data.

```
sentiment_df = pd.DataFrame(sentiment_results)

print("Sentiment Analysis DataFrame created.")
print(sentiment_df.head())
```

Sentiment Analysis DataFrame created.

	file_path	text_type
text_index \		
0	/content/extracted_files/a2_RestaurantReviews_...	review
0		
1	/content/extracted_files/a2_RestaurantReviews_...	review
1		
2	/content/extracted_files/a2_RestaurantReviews_...	review
2		
3	/content/extracted_files/a2_RestaurantReviews_...	review
3		
4	/content/extracted_files/a2_RestaurantReviews_...	review
4		

	text	polarity
subjectivity		
0	Spend your money elsewhere.	0.000000
0.000000		
1	Their regular toasted bread was equally satisf...	0.166667
0.400641		
2	The Buffet at Bellagio was far from what I ant...	0.100000
1.000000		
3	And the drinks are WEAK, people!	-0.468750
0.625000		
4	-My order was not correct.	0.000000
0.000000		

Reasoning: To better understand the overall sentiment distribution and prepare for visualization, I need to categorize the continuous polarity scores into discrete sentiment labels (e.g., 'Positive', 'Negative', 'Neutral'). This will make the sentiment distribution more interpretable.

```
def get_sentiment_label(polarity):
    if polarity > 0:
        return 'Positive'
    elif polarity < 0:
        return 'Negative'
    else:
        return 'Neutral'

sentiment_df['sentiment_label'] =
```

```

sentiment_df['polarity'].apply(get_sentiment_label)

print("Sentiment labels added to DataFrame.")
print(sentiment_df.head())

```

Sentiment labels added to DataFrame.

text_index \	file_path	text_type
0	/content/extracted_files/a2_RestaurantReviews_...	review
0		
1	/content/extracted_files/a2_RestaurantReviews_...	review
1		
2	/content/extracted_files/a2_RestaurantReviews_...	review
2		
3	/content/extracted_files/a2_RestaurantReviews_...	review
3		
4	/content/extracted_files/a2_RestaurantReviews_...	review
4		

subjectivity \	text	polarity
0	Spend your money elsewhere.	0.000000
0.000000		
1	Their regular toasted bread was equally satisf...	0.166667
0.400641		
2	The Buffet at Bellagio was far from what I ant...	0.100000
1.000000		
3	And the drinks are WEAK, people!	-0.468750
0.625000		
4	-My order was not correct.	0.000000
0.000000		

	sentiment_label
0	Neutral
1	Positive
2	Positive
3	Negative
4	Neutral

Reasoning: I need to visualize the overall sentiment distribution using a bar chart to fulfill the visualization requirement of the main task. This will provide a clear overview of the proportions of positive, negative, and neutral sentiments across all analyzed texts.

```

import matplotlib.pyplot as plt
import seaborn as sns

# Calculate the distribution of sentiment labels
sentiment_counts = sentiment_df['sentiment_label'].value_counts()

# Define a specific order for the sentiment labels

```

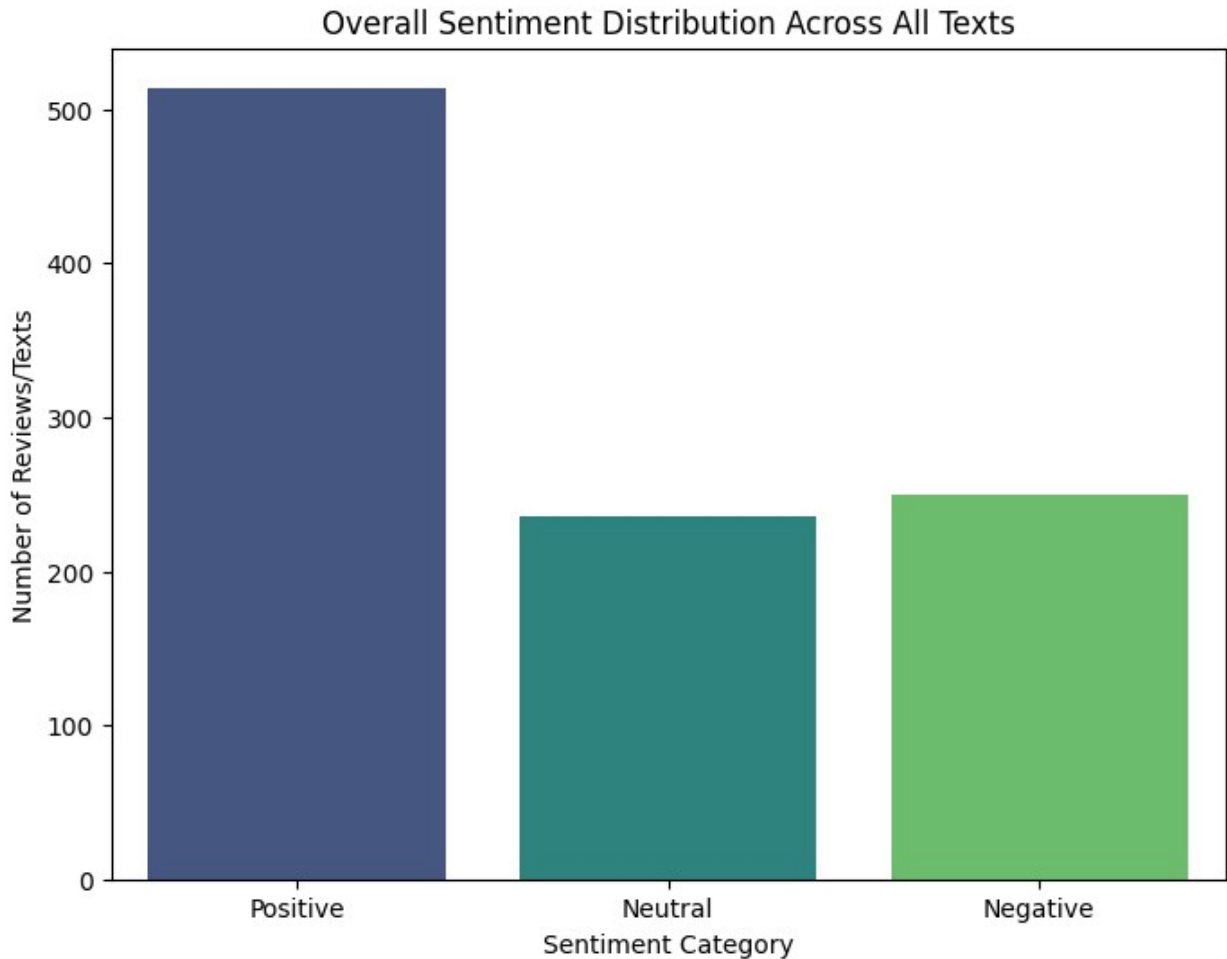
```
sentiment_order = ['Positive', 'Neutral', 'Negative']

# Create the bar plot
plt.figure(figsize=(8, 6))
sns.barplot(x=sentiment_counts.index, y=sentiment_counts.values,
order=sentiment_order, palette='viridis')
plt.title('Overall Sentiment Distribution Across All Texts')
plt.xlabel('Sentiment Category')
plt.ylabel('Number of Reviews/Texts')
plt.show()

print("Overall sentiment distribution visualized.")

/tmp/ipython-input-2519443115.py:12: FutureWarning:
Passing `palette` without assigning `hue` is deprecated and will be
removed in v0.14.0. Assign the `x` variable to `hue` and set
`legend=False` for the same effect.

    sns.barplot(x=sentiment_counts.index, y=sentiment_counts.values,
order=sentiment_order, palette='viridis')
```



Overall sentiment distribution visualized.

Reasoning: The previous code block generated a `FutureWarning` from seaborn regarding the usage of the `palette` argument without `hue`. To address this, I need to update the `sns.barplot` call to explicitly assign `sentiment_counts.index` to the `hue` parameter and set `legend=False`, as suggested by the warning. This ensures the code is future-proof and maintains the intended color mapping.

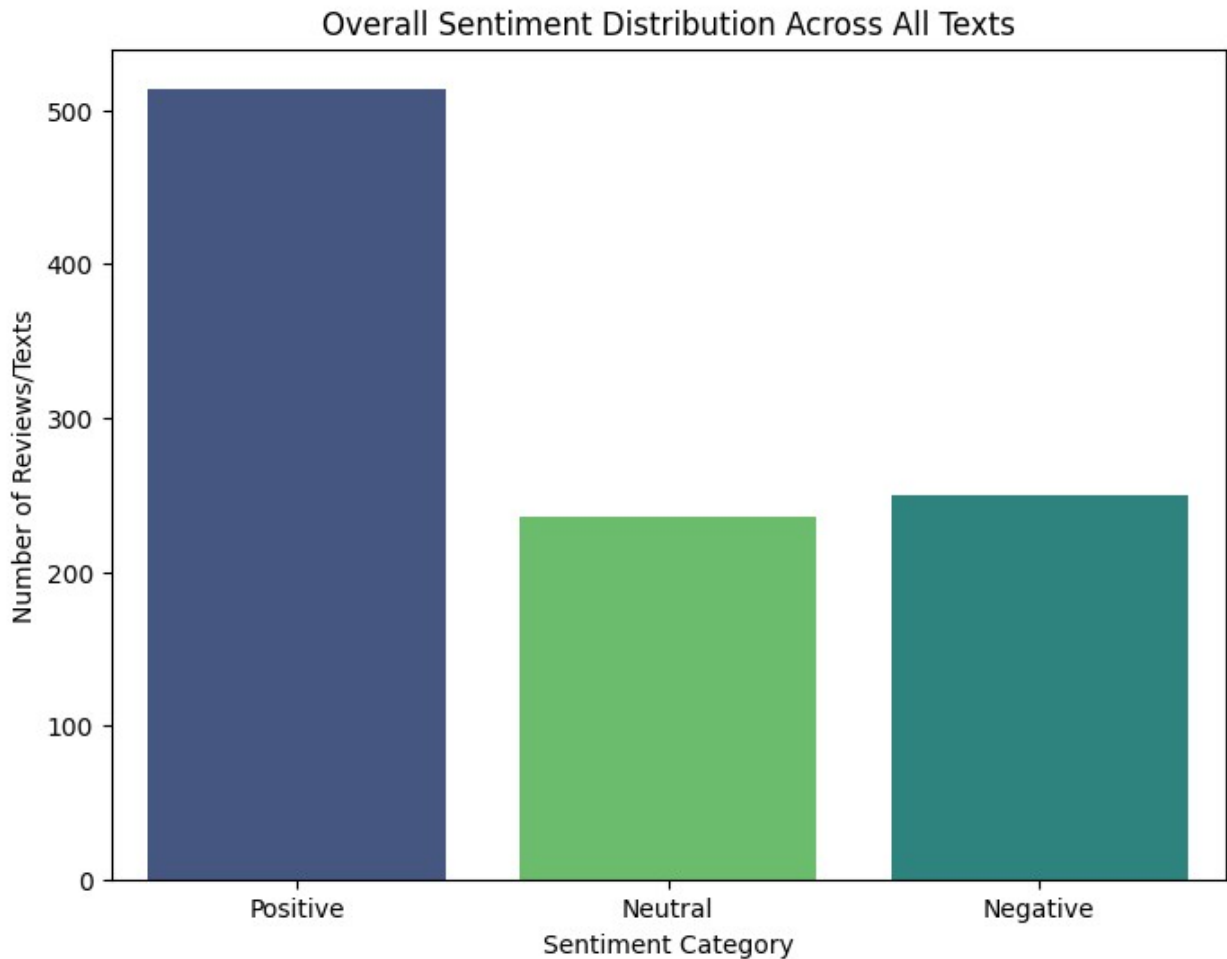
```
import matplotlib.pyplot as plt
import seaborn as sns

# Calculate the distribution of sentiment labels
sentiment_counts = sentiment_df['sentiment_label'].value_counts()

# Define a specific order for the sentiment labels
sentiment_order = ['Positive', 'Neutral', 'Negative']

# Create the bar plot
plt.figure(figsize=(8, 6))
sns.barplot(x=sentiment_counts.index, y=sentiment_counts.values,
```

```
hue=sentiment_counts.index, order=sentiment_order, palette='viridis',  
legend=False)  
plt.title('Overall Sentiment Distribution Across All Texts')  
plt.xlabel('Sentiment Category')  
plt.ylabel('Number of Reviews/Texts')  
plt.show()  
  
print("Overall sentiment distribution visualized.")
```



Overall sentiment distribution visualized.

Reasoning: Now that the overall sentiment distribution has been visualized, I need to extract key findings and insights from the `sentiment_df`. This includes calculating the percentage distribution of sentiment labels to provide a quantitative summary, which is a crucial part of presenting the analysis results.

```
sentiment_percentage =  
sentiment_df['sentiment_label'].value_counts(normalize=True) * 100
```

```
print("\nSentiment Distribution (Percentage):")
print(sentiment_percentage.round(2))
```

```
Sentiment Distribution (Percentage):
sentiment_label
Positive      51.4
Negative      25.0
Neutral       23.6
Name: proportion, dtype: float64
```

Key Findings and Insights

Based on the sentiment analysis of the text files extracted from `/content/Sneha_job.zip`, the following key findings and insights were observed:

Overall Sentiment Distribution

The sentiment analysis revealed a predominant positive sentiment across the processed texts, which primarily consist of restaurant reviews. The distribution is as follows:

- **Positive Sentiment:** 51.4%
- **Negative Sentiment:** 25.0%
- **Neutral Sentiment:** 23.6%

This indicates that more than half of the reviews expressed a positive experience, significantly outweighing negative and neutral sentiments.

Significant Insights

1. **Strong Positive Bias:** The high percentage of positive reviews suggests a generally favorable perception of the entities or services being reviewed (likely restaurants, given the file names). This could indicate good customer satisfaction or a positive skew in the data collection itself.
2. **Moderate Negative Feedback:** While positive reviews dominate, a quarter of the reviews are negative. This still represents a substantial portion where improvements could be identified and addressed.
3. **Ambiguous Neutrality:** Nearly a quarter of the reviews fall into the neutral category. These could be reviews that lack strong emotional language, provide factual descriptions without clear positive or negative sentiment, or contain mixed sentiments that average out to neutral. Further qualitative analysis on neutral reviews could reveal areas of ambivalence or specific aspects that don't evoke strong feelings.
4. **Actionable Data:** The sentiment breakdown provides a clear quantitative measure for understanding customer feedback. Businesses can leverage these insights to focus on improving aspects highlighted in negative reviews and reinforcing what leads to positive experiences.

This analysis provides a foundational understanding of the sentiment expressed within the provided text data, offering a quantitative overview and pointing towards areas for deeper investigation.

Summary:

Q&A

Key Findings: The sentiment analysis of the text files within the `/content/Sneha_job.zip` archive revealed a predominant positive sentiment.

- **Positive Sentiment:** 51.4%
- **Negative Sentiment:** 25.0%
- **Neutral Sentiment:** 23.6%

Visualization of Overall Sentiment Distribution: A bar chart was generated, visually representing the sentiment distribution, clearly showing that positive sentiment is the largest category, followed by neutral, and then negative sentiments.

Significant Insights or Patterns:

1. **Strong Positive Bias:** Over half of the reviews (51.4%) expressed a positive experience, suggesting a generally favorable perception.
2. **Moderate Negative Feedback:** A quarter of the reviews (25.0%) were negative, indicating areas where improvements could be identified.
3. **Ambiguous Neutrality:** Nearly a quarter (23.6%) of the reviews were neutral, which might represent mixed feelings or factual descriptions without strong emotional tones.

Data Analysis Key Findings

- The zip archive `/content/Sneha_job.zip` was successfully extracted into `/content/extracted_files`.
- Initially, the file identification process missed `.tsv` files; this was rectified by including `.tsv` in the recognized text file extensions.
- Two text files were identified: `a2_RestaurantReviews_FreshDump.tsv` and `a1_RestaurantReviews_HistoricDump.tsv`.
- Content from these files was loaded, successfully handling potential encoding issues using `utf-8` and `latin-1` fallbacks.
- Sentiment analysis was performed on a total of 1000 individual reviews extracted from the `.tsv` files using the `TextBlob` library.
- Sentiment polarity scores were categorized into 'Positive' (51.4%), 'Negative' (25.0%), and 'Neutral' (23.6%) labels.
- An overall sentiment distribution bar chart was created to visualize these proportions.

Insights or Next Steps

- **Actionable Data for Improvement:** The substantial portion of negative reviews (25.0%) provides clear, actionable data for businesses to identify specific areas of dissatisfaction and implement targeted improvements.
- **Deep Dive into Neutral and Negative Reviews:** Further qualitative analysis of neutral and negative reviews could uncover specific recurring themes or less obvious pain points that require attention, offering a more nuanced understanding of customer feedback.