```python
#sentiment analysis in python using two different techniques:
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns


plt.style.use('ggplot')


import nltk
```

```python
!pip install lime
!pip install shap
```

```
Collecting lime
    Downloading lime-0.2.0.1.tar.gz (275 kB)
    ─────────────────────────────────────── 275.7/275.7 kB 8.2 MB/s eta 0:00:00
    Preparing metadata (setup.py) ... done
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from lime) (3.7.1)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from lime) (1.26.4)
Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from lime) (1.13.1)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from lime) (4.66.5)
Requirement already satisfied: scikit-learn>=0.18 in /usr/local/lib/python3.10/dist-packages (from lime) (1.5.2)
Requirement already satisfied: scikit-image>=0.12 in /usr/local/lib/python3.10/dist-packages (from lime) (0.24.0)
Requirement already satisfied: networkx>=2.8 in /usr/local/lib/python3.10/dist-packages (from scikit-image>=0.12->lime) (3.3)
Requirement already satisfied: pillow>=9.1 in /usr/local/lib/python3.10/dist-packages (from scikit-image>=0.12->lime) (10.4.0)
Requirement already satisfied: imageio>=2.33 in /usr/local/lib/python3.10/dist-packages (from scikit-image>=0.12->lime) (2.35.1)
Requirement already satisfied: tifffile>=2022.8.12 in /usr/local/lib/python3.10/dist-packages (from scikit-image>=0.12->lime) (2024.9
Requirement already satisfied: packaging>=21 in /usr/local/lib/python3.10/dist-packages (from scikit-image>=0.12->lime) (24.1)
Requirement already satisfied: lazy-loader>=0.4 in /usr/local/lib/python3.10/dist-packages (from scikit-image>=0.12->lime) (0.4)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn>=0.18->lime) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn>=0.18->lime) (3.5.0
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->lime) (1.3.0)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib->lime) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->lime) (4.54.1)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->lime) (1.4.7)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->lime) (3.1.4)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib->lime) (2.8.2)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.7->matplotlib->lime) (1.16
Building wheels for collected packages: lime
    Building wheel for lime (setup.py) ... done
```

```
      Created wheel for lime: filename=lime-0.2.0.1-py3-none-any.whl size=283834 sha256=39753c09eb5ea6402b9f4820423de5a75c3ab13648e75b404c
      Stored in directory: /root/.cache/pip/wheels/fd/a2/af/9ac0a1a85a27f314a06b39e1f492bee1547d52549a4606ed89
    Successfully built lime
    Installing collected packages: lime
    Successfully installed lime-0.2.0.1
    Collecting shap
      Downloading shap-0.46.0-cp310-cp310-manylinux_2_12_x86_64.manylinux2010_x86_64.manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metada
    Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from shap) (1.26.4)
    Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from shap) (1.13.1)
    Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/dist-packages (from shap) (1.5.2)
    Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (from shap) (2.2.2)
    Requirement already satisfied: tqdm>=4.27.0 in /usr/local/lib/python3.10/dist-packages (from shap) (4.66.5)
    Requirement already satisfied: packaging>20.9 in /usr/local/lib/python3.10/dist-packages (from shap) (24.1)
    Collecting slicer==0.0.8 (from shap)
      Downloading slicer-0.0.8-py3-none-any.whl.metadata (4.0 kB)
    Requirement already satisfied: numba in /usr/local/lib/python3.10/dist-packages (from shap) (0.60.0)
    Requirement already satisfied: cloudpickle in /usr/local/lib/python3.10/dist-packages (from shap) (2.2.1)
    Requirement already satisfied: llvmlite<0.44,>=0.43.0dev0 in /usr/local/lib/python3.10/dist-packages (from numba->shap) (0.43.0)
    Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas->shap) (2.8.2)
    Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas->shap) (2024.2)
    Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas->shap) (2024.2)
    Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn->shap) (1.4.2)
    Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn->shap) (3.5.0)
    Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.2->pandas->shap) (1.16.0
    Downloading shap-0.46.0-cp310-cp310-manylinux_2_12_x86_64.manylinux2010_x86_64.manylinux_2_17_x86_64.manylinux2014_x86_64.whl (540 kB
    ———————————————————————————————————— 540.1/540.1 kB 16.2 MB/s eta 0:00:00
    Downloading slicer-0.0.8-py3-none-any.whl (15 kB)
    Installing collected packages: slicer, shap
    Successfully installed shap-0.46.0 slicer-0.0.8
```

```python
df = pd.read_csv('/content/Reviews.csv')
print(df.shape)
df = df.head(500)
print(df.shape)
```

```
(29941, 10)
(500, 10)
```

```python
df.head()
```

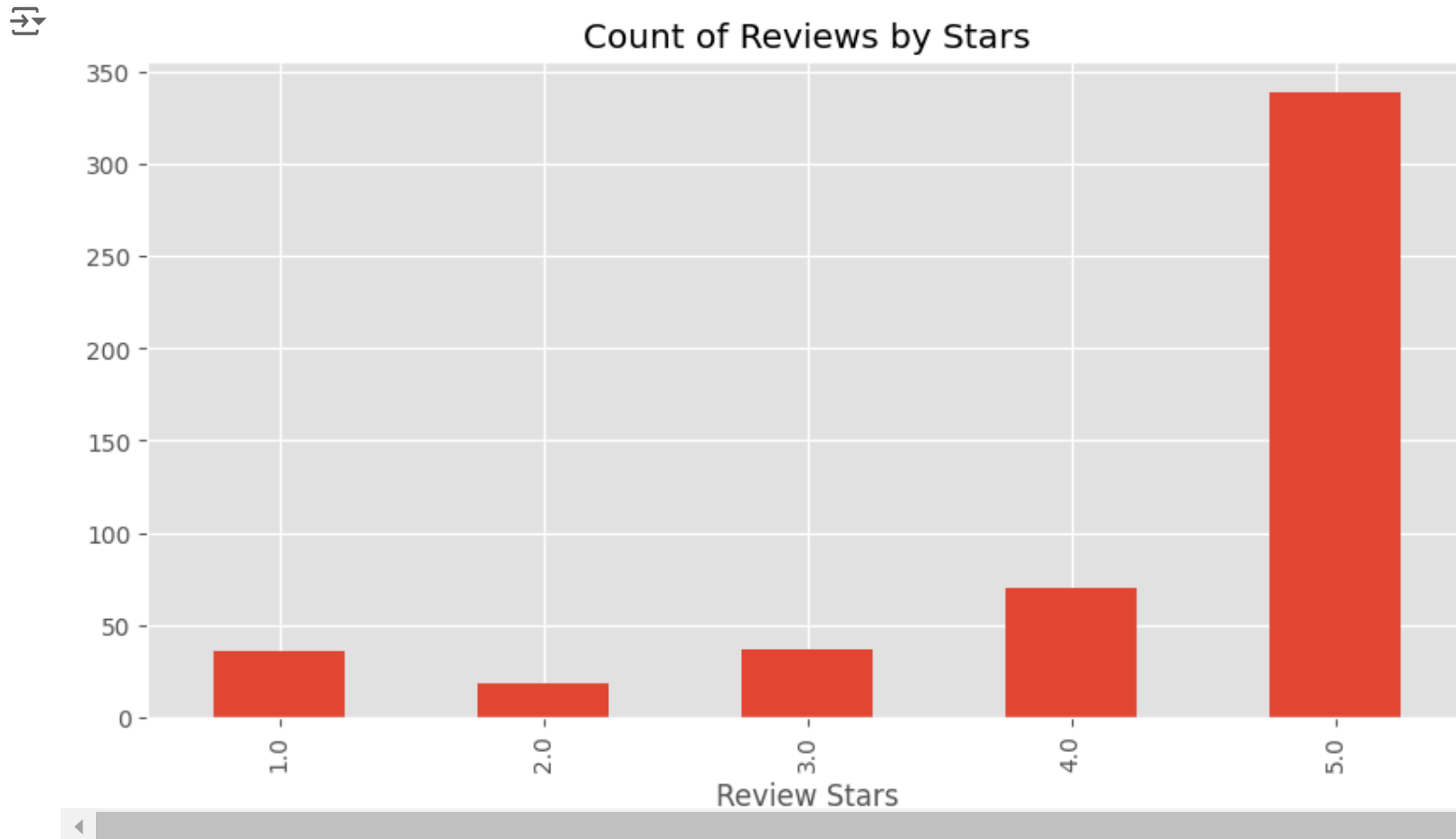| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | Time | Summary | Te |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1.0 | 1.0 | 5.0 | 1.303862e+09 | Good Quality Dog Food | I ha boug several t Vita cann c |
| 1 | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0.0 | 0.0 | 1.0 | 1.346976e+09 | Not as Advertised | Prod arriv labeled Jum Salt Peanu |
| 2 | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1.0 | 1.0 | 4.0 | 1.219018e+09 | "Delight" says it all | This i confecti that h be around fe |
| 3 | 4 | B000UA0QIQ | A395BORC6FGVXV | Karl | 3.0 | 3.0 | 2.0 | 1.307923e+09 | Cough Medicine | If you looki for t sec ingredi |
| 4 | 5 | B006K2ZZ7K | A1UQRSCLF8GW1T | Michael D. Bigham "M. Wassir" | 0.0 | 0.0 | 5.0 | 1.350778e+09 | Great taffy | Gr taffy a gre pri The wa wi |

Next steps:   Generate code with df      ◉ View recommended plots      New interactive sheet

```python
ax = df['Score'].value_counts().sort_index() \
    .plot(kind='bar',
          title='Count of Reviews by Stars',
          figsize=(10, 5))
ax.set_xlabel('Review Stars')
plt.show()
```



```python
example = df['Text'][50]
print(example)
```

Show hidden output

## SETP :1 **VADER Seniment Scoring**

We will use NLTK's SentimentIntensityAnalyzer to get the neg/neu/pos scores of the text.

```
import nltk

# Download the 'averaged_perceptron_tagger' resource
nltk.download('averaged_perceptron_tagger')

# Now, you can use the pos_tag function
tagged = nltk.pos_tag(tokens)
tagged[:10]
```

Show hidden output

Generate | print hello world using rot13 | 🔍 | Close

```
import nltk

# Download the 'words' resource
nltk.download('words')

# Download the 'maxent_ne_chunker' resource
nltk.download('maxent_ne_chunker')

# Now, you can use the ne_chunk function
entities = nltk.chunk.ne_chunk(tagged)
entities.pprint()
```

Show hidden output

```
import nltk
from nltk.sentiment import SentimentIntensityAnalyzer
from tqdm.notebook import tqdm

# Download the VADER lexicon
nltk.download('vader_lexicon')
```

```
sia = SentimentIntensityAnalyzer()
```

> [nltk_data] Downloading package vader_lexicon to /root/nltk_data...

```
sia.polarity_scores('I am so happy!')
```

> {'neg': 0.0, 'neu': 0.318, 'pos': 0.682, 'compound': 0.6468}

```
sia.polarity_scores('This is the worst thing ever.')
```

> {'neg': 0.451, 'neu': 0.549, 'pos': 0.0, 'compound': -0.6249}

```
sia.polarity_scores(example)
```

> {'neg': 0.22, 'neu': 0.78, 'pos': 0.0, 'compound': -0.5448}

```
res = {}
for i, row in tqdm(df.iterrows(), total=len(df)):
    text = row['Text']
    myid = row['Id']
    res[myid] = sia.polarity_scores(text)
```

> 100%                                            500/500 [00:00<00:00, 1006.15it/s]

```
vaders = pd.DataFrame(res).T
vaders = vaders.reset_index().rename(columns={'index': 'Id'})
vaders = vaders.merge(df, how='left')
```
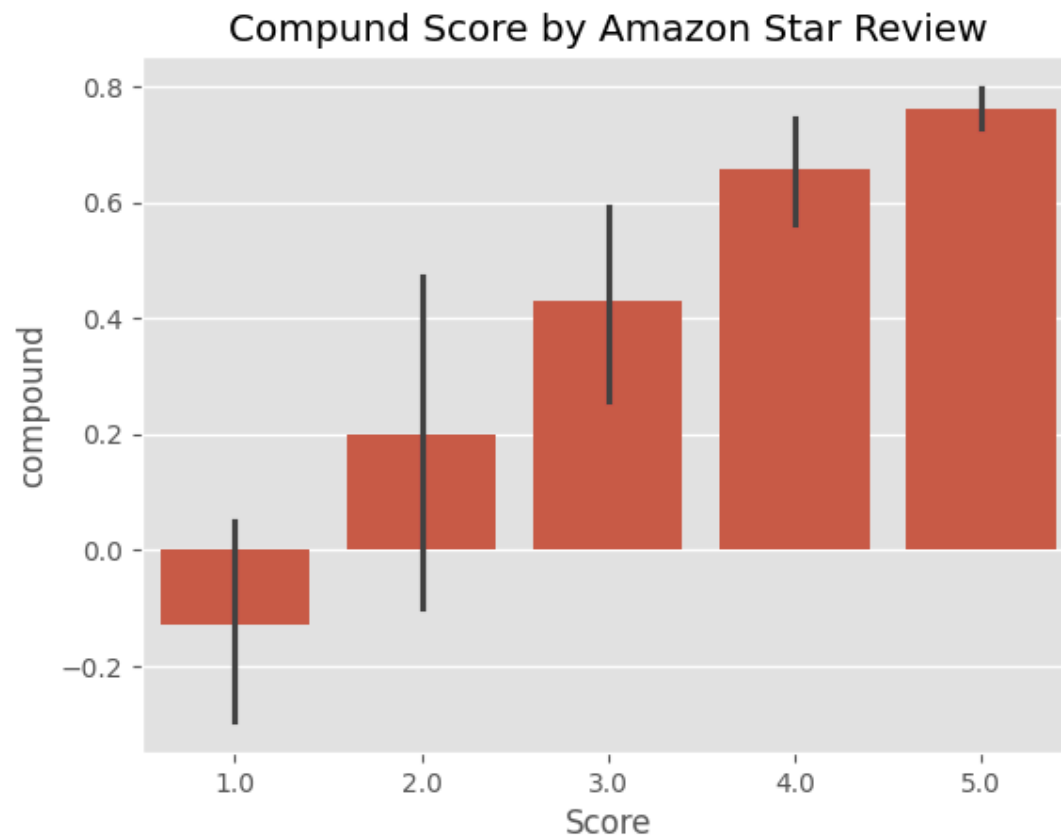
```
vaders.head()
```

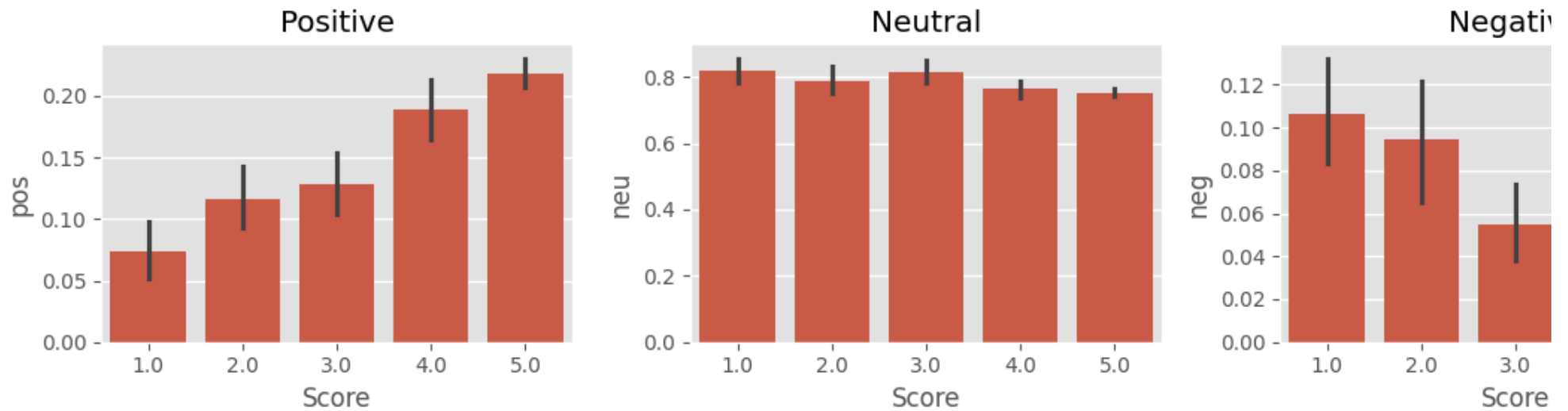| | Id | neg | neu | pos | compound | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0.000 | 0.695 | 0.305 | 0.9441 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1.0 | 1.0 | 5.0 | 1 |
| **1** | 2 | 0.138 | 0.862 | 0.000 | -0.5664 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0.0 | 0.0 | 1.0 | 1 |
| **2** | 3 | 0.091 | 0.754 | 0.155 | 0.8265 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1.0 | 1.0 | 4.0 | 1 |
| **3** | 4 | 0.000 | 1.000 | 0.000 | 0.0000 | B000UA0QIQ | A395BORC6FGVXV | Karl | 3.0 | 3.0 | 2.0 | 1 |
| **4** | 5 | 0.000 | 0.552 | 0.448 | 0.9468 | B006K2ZZ7K | A1UQRSCLF8GW1T | Michael D. Bigham "M. Wassir" | 0.0 | 0.0 | 5.0 | 1 |

Next steps:   Generate code with `vaders`      View recommended plots      New interactive sheet

```
ax = sns.barplot(data=vaders, x='Score', y='compound')
ax.set_title('Compund Score by Amazon Star Review')
plt.show()
```



Compund Score by Amazon Star Review

```
fig, axs = plt.subplots(1, 3, figsize=(12, 3))
sns.barplot(data=vaders, x='Score', y='pos', ax=axs[0])
sns.barplot(data=vaders, x='Score', y='neu', ax=axs[1])
sns.barplot(data=vaders, x='Score', y='neg', ax=axs[2])
axs[0].set_title('Positive')
axs[1].set_title('Neutral')
axs[2].set_title('Negative')
plt.tight_layout()
plt.show()
```

## Roberta Pretrained Model

```
from transformers import AutoTokenizer
from transformers import AutoModelForSequenceClassification
from scipy.special import softmax
```

```
MODEL = f"cardiffnlp/twitter-roberta-base-sentiment"
tokenizer = AutoTokenizer.from_pretrained(MODEL)
model = AutoModelForSequenceClassification.from_pretrained(MODEL)
```

Show hidden output

```
print(example)
sia.polarity_scores(example)
```

```
This oatmeal is not good. Its mushy, soft, I don't like it. Quaker Oats is the way to go.
{'neg': 0.22, 'neu': 0.78, 'pos': 0.0, 'compound': -0.5448}
```

```python
encoded_text = tokenizer(example, return_tensors='pt')
output = model(**encoded_text)
scores = output[0][0].detach().numpy()
scores = softmax(scores)
scores_dict = {
    'roberta_neg' : scores[0],
    'roberta_neu' : scores[1],
    'roberta_pos' : scores[2]
}
print(scores_dict)
```

⇥ {'roberta_neg': 0.97635514, 'roberta_neu': 0.020687465, 'roberta_pos': 0.0029573692}

```python
def polarity_scores_roberta(example):
    encoded_text = tokenizer(example, return_tensors='pt')
    output = model(**encoded_text)
    scores = output[0][0].detach().numpy()
    scores = softmax(scores)
    scores_dict = {
        'roberta_neg' : scores[0],
        'roberta_neu' : scores[1],
        'roberta_pos' : scores[2]
    }
    return scores_dict
```

```python
res = {}
for i, row in tqdm(df.iterrows(), total=len(df)):
    try:
        text = row['Text']
        myid = row['Id']
        vader_result = sia.polarity_scores(text)
        vader_result_rename = {}
        for key, value in vader_result.items():
            vader_result_rename[f"vader_{key}"] = value
        roberta_result = polarity_scores_roberta(text)
        both = {**vader_result_rename, **roberta_result}
        res[myid] = both
    except RuntimeError:
        print(f'Broke for id {myid}')
```

⮞  100%                                                          500/500 [04:27<00:00,   2.56it/s]
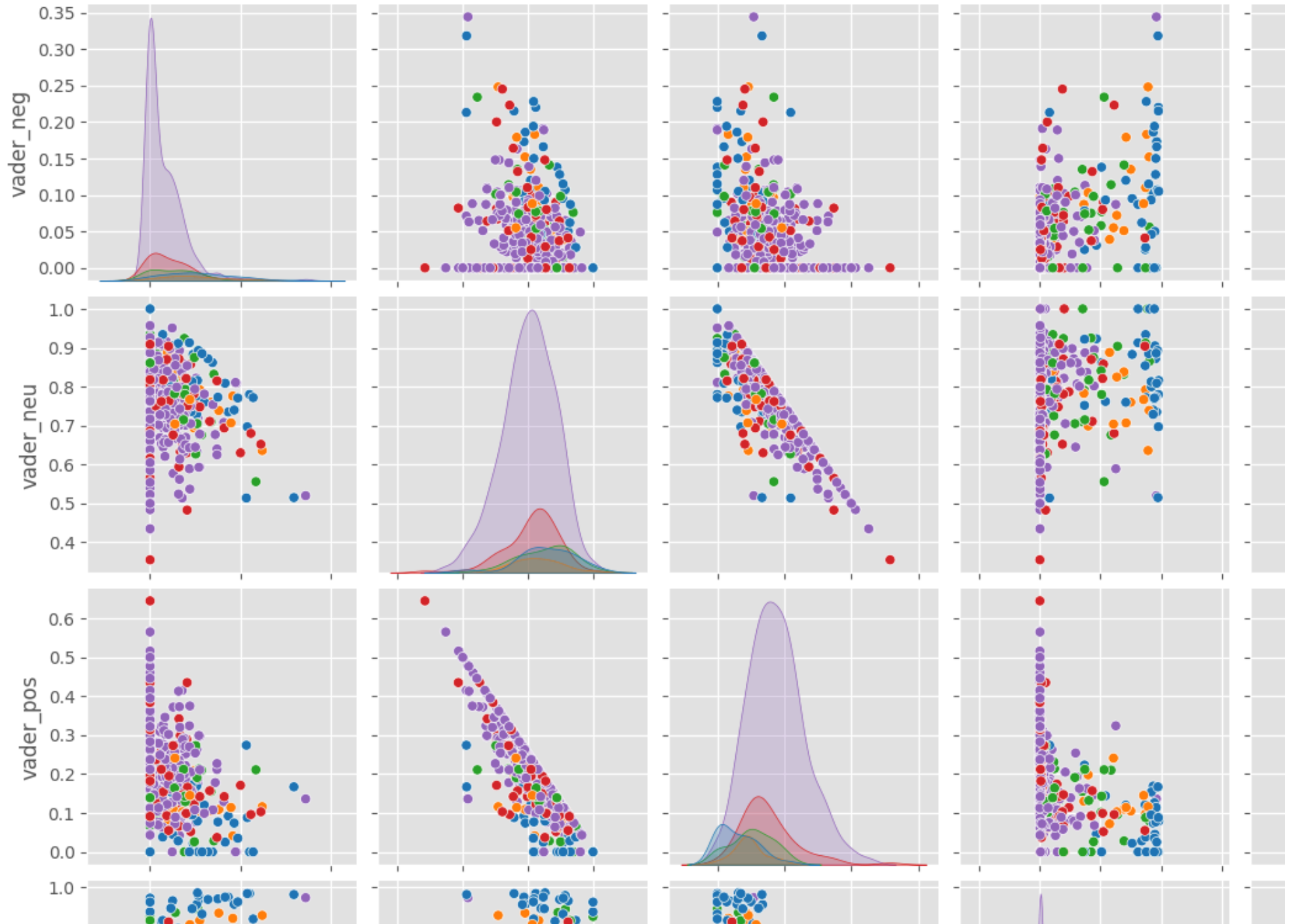
      Broke for id 83
      Broke for id 187

```
results_df = pd.DataFrame(res).T
results_df = results_df.reset_index().rename(columns={'index': 'Id'})
results_df = results_df.merge(df, how='left')
```
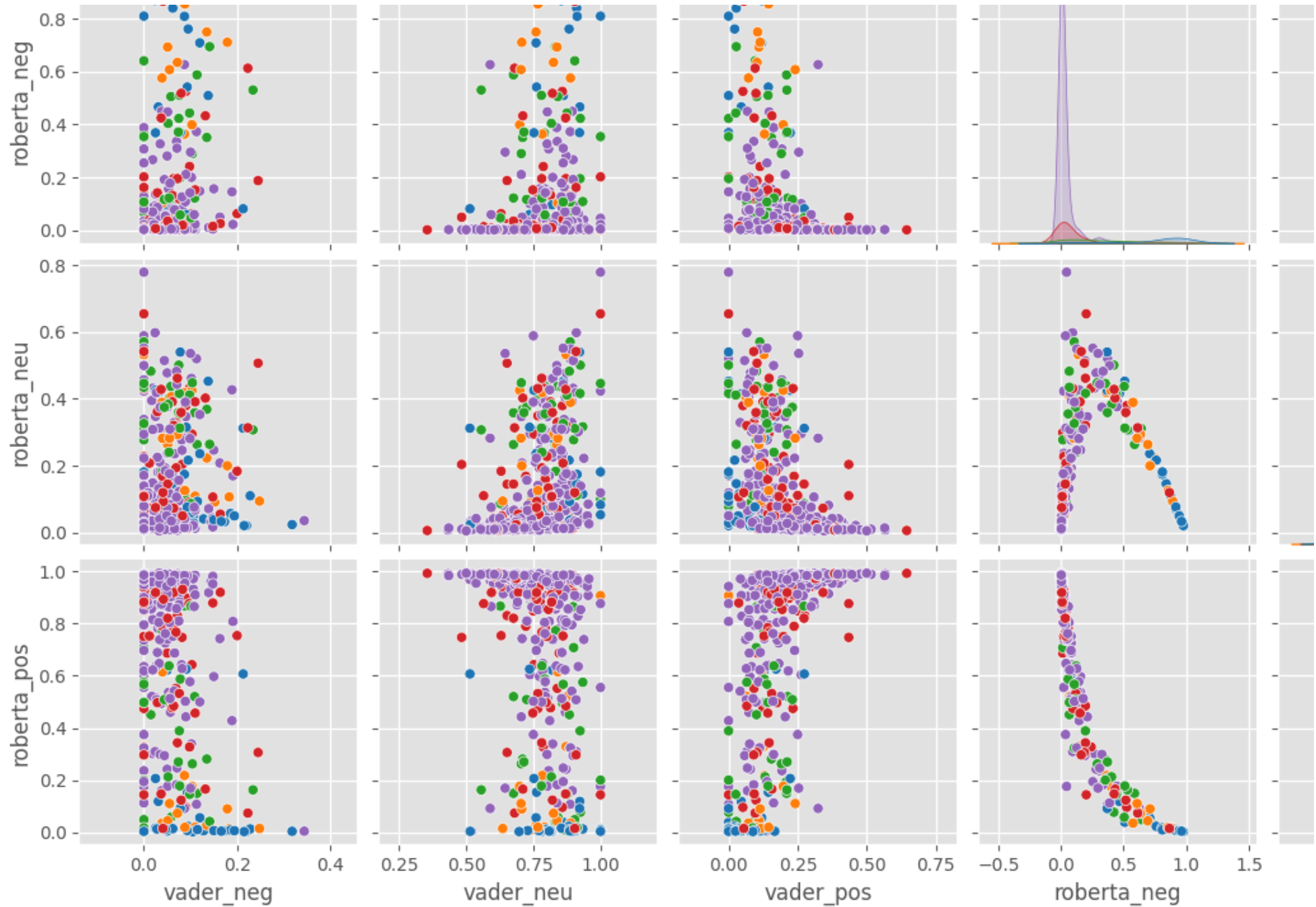
```
results_df.columns
```

⮞  Index(['Id', 'vader_neg', 'vader_neu', 'vader_pos', 'vader_compound',
             'roberta_neg', 'roberta_neu', 'roberta_pos', 'ProductId', 'UserId',
             'ProfileName', 'HelpfulnessNumerator', 'HelpfulnessDenominator',
             'Score', 'Time', 'Summary', 'Text'],
          dtype='object')

```
sns.pairplot(data=results_df,
             vars=['vader_neg', 'vader_neu', 'vader_pos',
                   'roberta_neg', 'roberta_neu', 'roberta_pos'],
            hue='Score',
            palette='tab10')
plt.show()
```

**Review Examples:**

```
results_df.query('Score == 1') \
    .sort_values('roberta_pos', ascending=False)['Text'].values[0]
```

→ 'I felt energized within five minutes, but it lasted for about 45 minutes. I paid $3.99 for this drink. I could have just drunk a cup o
   f coffee and saved my money.'

```
results_df.query('Score == 1') \
    .sort_values('vader_pos', ascending=False)['Text'].values[0]
```

→ 'So we cancelled the order.  It was cancelled without any problem.  That is a positive note...'

```
results_df.query('Score == 5') \
    .sort_values('roberta_neg', ascending=False)['Text'].values[0]
```

→ 'this was sooooo deliscious but too bad i ate em too fast and gained 2 pds! my fault'

```
results_df.query('Score == 5') \
    .sort_values('vader_neg', ascending=False)['Text'].values[0]
```

→ 'this was sooooo deliscious but too bad i ate em too fast and gained 2 pds! my fault'

**The Transformers Pipeline**

Quick & easy way to run sentiment predictions

```
from transformers import pipeline

sent_pipeline = pipeline("sentiment-analysis")
```

No model was supplied, defaulted to distilbert/distilbert-base-uncased-finetuned-sst-2-english and revision af0f99b (https://huggingface
Using a pipeline without specifying a model name and revision in production is not recommended.

config.json: 100%                                          629/629 [00:00<00:00, 14.7kB/s]

model.safetensors: 100%                                    268M/268M [00:01<00:00, 164MB/s]

tokenizer_config.json: 100%                                48.0/48.0 [00:00<00:00, 3.02kB/s]

vocab.txt: 100%                                            232k/232k [00:00<00:00, 9.40MB/s]

/usr/local/lib/python3.10/dist-packages/transformers/tokenization_utils_base.py:1601: FutureWarning: `clean_up_tokenization_spaces` was

```
sent_pipeline('I love sentiment analysis!')
```

[{'label': 'POSITIVE', 'score': 0.9997853636741638}]

```
sent_pipeline('ITS GOOD TO BE HERE')
```

[{'label': 'POSITIVE', 'score': 0.9998099207878113}]

```
sent_pipeline('I LOVE THIS THING')
```

[{'label': 'POSITIVE', 'score': 0.9998818635940552}]

```
sent_pipeline('The end')
```