

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 plt.style.use('ggplot')
6 import nltk

```

```

1 # Read in data
2 df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/Dataset/Sentimental Analysis_Amazon Reviews.csv')
3 print(df.shape)
4 df = df.head(500)
5 print(df.shape)

```

```

(40437, 10)
(500, 10)

```

```

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

```



```

1 example = df['Text'][50]
2 print(example)

```

This oatmeal is not good. Its mushy, soft, I don't like it. Quaker Oats is the way to go.

```

1 nltk.download('punkt')
2 nltk.download('averaged_perceptron_tagger')
3 nltk.download('maxent_ne_chunker')

```

```

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] /root/nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
[nltk_data] Downloading package maxent_ne_chunker to
[nltk_data] /root/nltk_data...
[nltk_data] Package maxent_ne_chunker is already up-to-date!
True

```

```

1 import nltk
2 example = "370 chars 'Got a free package ...'"
3 tokens = nltk.word_tokenize(example)
4 tokens[:10]

```

```
['370', 'chars', "'Got", 'a', 'free', 'package', '...']
```

```

1 tagged = nltk.pos_tag(tokens)
2 tagged[:10]

```

```
1 nltk.download('words')
2 entities = nltk.chunk.ne_chunk(tagged)
3 entities.pprint()
```

```
1 nltk.download('vader_lexicon')
2 from nltk.sentiment import SentimentIntensityAnalyzer
3 from tqdm.notebook import tqdm
4 sia = SentimentIntensityAnalyzer()
```

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

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

```
1 sia.polarity_scores(example)

{'neg': 0.0, 'neu': 0.548, 'pos': 0.452, 'compound': 0.5106}
```

100% 500/500 [00:00<00:00, 1610.24it/s]

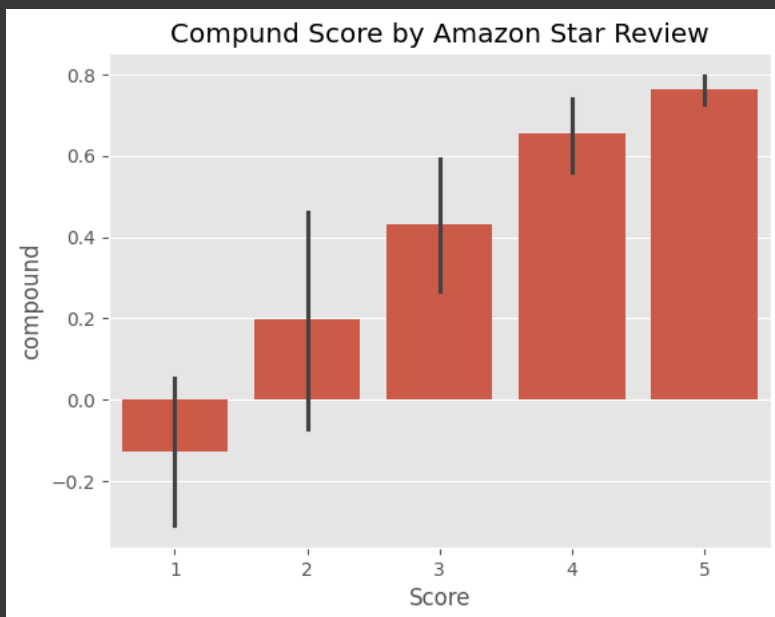
```
1 # Now we have sentiment score and metadata
2 vaders.head()
```

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

```

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

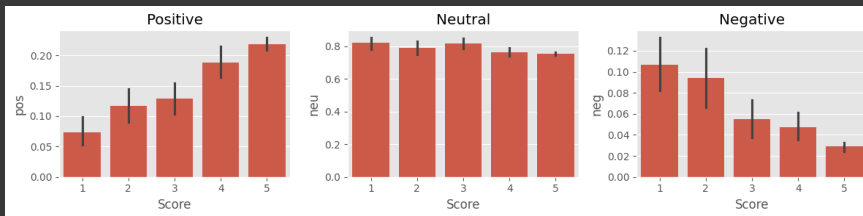
```



```

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

```



```
1 # Pre Trained Model
2 from transformers import AutoTokenizer
3 from transformers import AutoModelForSequenceClassification
4 from scipy.special import softmax
```

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

```
1 # VADER results on example
2 print(example)
3 sia.polarity_scores(example)
```

```
370 chars 'Got a free package ...
{'neg': 0.0, 'neu': 0.548, 'pos': 0.452, 'compound': 0.5106}
```

```
1 # Run for Roberta Model
2 encoded_text = tokenizer(example, return_tensors='pt')
3 output = model(**encoded_text)
4 scores = output[0][0].detach().numpy()
5 scores = softmax(scores)
6 scores_dict = {
7     'roberta_neg' : scores[0],
8     'roberta_neu' : scores[1],
9     'roberta_pos' : scores[2]
10 }
11 print(scores_dict)
```

```
{'roberta_neg': 0.015122626, 'roberta_neu': 0.7379595, 'roberta_pos': 0.24691787}
```

```
1 def polarity_scores_roberta(example):
2     encoded_text = tokenizer(example, return_tensors='pt')
3     output = model(**encoded_text)
4     scores = output[0][0].detach().numpy()
5     scores = softmax(scores)
6     scores_dict = {
7         'roberta_neg' : scores[0],
8         'roberta_neu' : scores[1],
9         'roberta_pos' : scores[2]
10     }
11     return scores_dict
```

```
1 res = {}
2 for i, row in tqdm(df.iterrows(), total=len(df)):
3     try:
4         text = row['Text']
5         myid = row['Id']
6         vader_result = sia.polarity_scores(text)
7         vader_result_rename = {}
8         for key, value in vader_result.items():
9             vader_result_rename[f"vader_{key}"] = value
10        roberta_result = polarity_scores_roberta(text)
11        both = {**vader_result_rename, **roberta_result}
12        res[myid] = both
13    except RuntimeError:
14        print(f'Broke for id {myid}')
```

100%   500/500 [02:12<00:00, 4.01it/s]

```
Broke for id 83
Broke for id 187
```

```

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

```

```

1 # Compare Scores between Models
2 results_df.columns

```

```

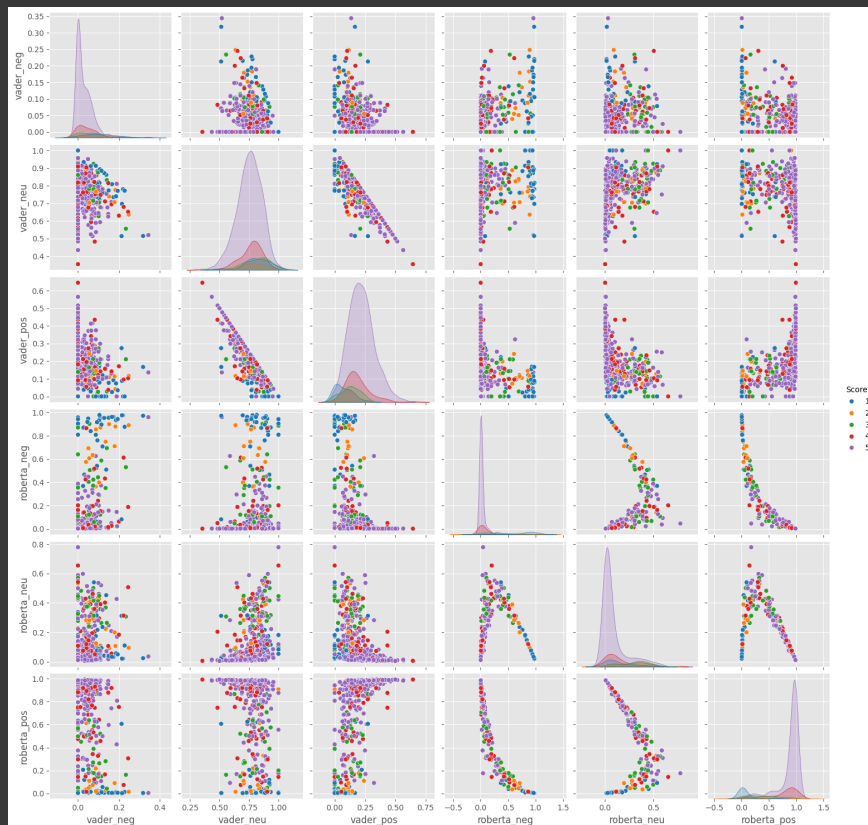
Index(['Id', 'vader_neg', 'vader_neu', 'vader_pos', 'vader_compound',
      'roberta_neg', 'roberta_neu', 'roberta_pos', 'ProductId', 'UserId',
      'ProfileName', 'HelpfulnessNumerator', 'HelpfulnessDenominator',
      'Score', 'Time', 'Summary', 'Text'],
      dtype='object')

```

```

1 # Combine and Compare
2 sns.pairplot(data=results_df,
3             vars=['vader_neg', 'vader_neu', 'vader_pos',
4                 'roberta_neg', 'roberta_neu', 'roberta_pos'],
5             hue='Score',
6             palette='tab10')
7 plt.show()

```



```
1 #Review Examples:Positive 1-Star and Negative 5-Star Reviews Lets look at some examples where the model scoring and review score diffe
2 results_df.query('Score == 1') \
3     .sort_values('roberta_pos', ascending=False)['Text'].values[0]
```

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

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

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

```
1 from transformers import pipeline
2
3 sent_pipeline = pipeline("sentiment-analysis")
```

config.json: 100% 629/629 [00:00<00:00, 19.9kB/s]

model.safetensors:	268M/268M [00:03<00:00,
--------------------	-------------------------

tokenizer\_config.json: [REDACTED] 48.0/48.0 [00:00&lt;00:00]