

Amazon Customer Review Analysis using Web Scraping And Sentimental Analysis

Importing Required Modules

```
In [1]: import requests # To send HTTP/1.1 requests
from bs4 import BeautifulSoup #convert the data return from request to HTML for parsing

import pandas as pd #to convert the reviews to dataframe

import string #remove punctuations
import nltk
from nltk.corpus import stopwords
from nltk import PorterStemmer

import matplotlib.pyplot as plt
import seaborn as sns

from collections import Counter # to count the most common 50 words
```

```
In [2]: pip install wordcloud

Requirement already satisfied: wordcloud in c:\users\hp\anaconda3\lib\site-packages (1.8.1)
Requirement already satisfied: numpy>=1.6.2 in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (3.3.4)
Requirement already satisfied: scipy<1.2.0,!=1.2.1,!=2.0.3 in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (1.2.0)
Requirement already satisfied: pillow in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (8.2.0)
Requirement already satisfied: pytz==2020.4,!=2.1.2,!=2.1.3 in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (2020.4)
Requirement already satisfied: kiwicrawler==1.0.1 in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (1.0.1)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (2020.12.5)
Requirement already satisfied: idna<3,>=2.5 in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (2.10)
Requirement already satisfied: chardet<3.0.2 in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (3.0.0)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (1.26.0)
Requirement already satisfied: six in c:\users\hp\anaconda3\lib\site-packages (from certifi>=2017.4.17>wordcloud) (1.15.0)
Note: you may need to restart the kernel to use updated packages.
```

```
In [3]: pip install plotly

Requirement already satisfied: plotly in c:\users\hp\anaconda3\lib\site-packages (5.3.1)
Requirement already satisfied: tenacity<2.0.0 in c:\users\hp\anaconda3\lib\site-packages (from plotly) (8.0.1)
Requirement already satisfied: six in c:\users\hp\anaconda3\lib\site-packages (from plotly) (1.15.0)
Note: you may need to restart the kernel to use updated packages.
```

```
In [4]: import plotly.express as px
from wordcloud import WordCloud
```

```
In [5]: pip install vaderSentiment

Requirement already satisfied: vaderSentiment in c:\users\hp\anaconda3\lib\site-packages (3.3.2)
Requirement already satisfied: requests in c:\users\hp\anaconda3\lib\site-packages (from vaderSentiment) (2.25.1)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\hp\anaconda3\lib\site-packages (from requests->vaderSentiment) (2020.12.5)
Requirement already satisfied: idna<3,>=2.5 in c:\users\hp\anaconda3\lib\site-packages (from requests->vaderSentiment) (2.10)
Requirement already satisfied: chardet<3.0.2 in c:\users\hp\anaconda3\lib\site-packages (from requests->vaderSentiment) (3.0.0)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\hp\anaconda3\lib\site-packages (from requests->vaderSentiment) (1.26.0)
Note: you may need to restart the kernel to use updated packages.
```

```
In [6]: from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
```

Defining Functions to Read Amazon Data and Extract Customer Reviews from it

```
In [7]: #Create a header that contains your request cookies, without cookies can not scrape amazon data it always shows error
def AmazonData(url):
    r = requests.get(url, headers=header)
    return r.text

def ToHTML(url):
    htmldata = AmazonData(url)
    soup = BeautifulSoup(htmldata, 'html.parser')
    return soup
```

```
In [8]: url ="https://www.amazon.com/Harry-Potter-Sorcerers-Stone-Rowling/dp/059035342X/ref=as_li_ss_til?ie=UTF8&linkCode=urllib.parse.urljoin(url, review_link)
```

```
In [9]: header = {'User-Agent':
'Mozilla/5.0 (Windows NT 10.0; Win64; x64) \
 AppleWebKit/537.36 (KHTML, like Gecko) \
 Chrome/90.0.4430.212 Safari/537.36',
'Accept-Language': 'en-US, en;q=0.5'}
```

```
In [10]: soup = ToHTML(url) # RAW HTML data
```

Extracting Customer Reviews

Extracting Customer Name

```
In [11]: def customer_name(soup):
    data_str = ""
    cus_list = []
    for item in soup.findall("span", class_ = "a-profile-name"): # the customer name is under HTML tag span, a-profile-name
        data_str = data_str + item.get_text()
        cus_list.append(data_str)
    data_str = ""
    return cus_list
```

```
In [12]: Customer_List = customer_name(soup)
```

```
In [13]: print(Customer_List)
```

```
[8]: ['Werkema', 'lilly', 'Allison Geitner', 'Hellogoodbuy', 'Kat', 'Sara Bertrand', 'Linda D. Combs', 'Talulah 3', 'Samires Dias', 'MJH', 'Anthony R', 'Mariana Nascimento Ferreira', 'abcde']
```

Extracting Customer Review

```
In [14]: def customer_review(soup):
    data_str = ""
    for item in soup.findall("div", class_ = "a-expander-content reviewText review-text-content a-expander-part"):
        data_str += data_str + item.get_text()
    review = data_str.split("\n")
    return review
```

```
In [15]: Cust_Rev = customer_review(soup)
Customer_Review_list = []
for reviews in Cust_Rev:
    if reviews == "":
        pass
    else:
        Customer_Review_list.append(reviews)
```

```
In [16]: # get link of see all reviews
review_link = soup.find("a", {"data-hook": "see-all-reviews-link-foot"})
```

```
In [17]: full_review_link = "https://www.amazon.com/*review_link['href']"
```

```
In [18]: full_review_link
```

```
Out[18]: 'https://www.amazon.com/Harry-Potter-Sorcerers-Stone-Rowling/product-reviews/059035342X/ref=cm_cr_dp_d_show_all_btm?ie=UTF8&reviewerType=all_reviews'
```

```
In [19]: reviews_list = []
for k in range(50):
    page = requests.get(full_review_link+"&pageNumber='"+str(k)+"",headers=header)
    soup2 = BeautifulSoup(page.content)
    for i in soup2.findAll("span", {"data-hook": "review-body"}):
        reviews_list.append(i.text)
```

```
In [20]: len(reviews_list)
```

```
Out[20]: 500
```

Converting the Reviews into a Dictionary

```
In [21]: Dict_reviews = {'reviews':reviews_list}
reviews_df = pd.DataFrame.from_dict(Dict_reviews)
```

```
In [22]: reviews_df.head()
```

```
Out[22]: reviews
```

```
[0]: \n\nI'm probably the only person on earth who...
```

```
[1]: \n\nI got this book, Harry Potter and the So...
```

```
[2]: \n\nThis book is a great value! If you can a...
```

```
[3]: \n\nSeen the movies a handful of times but decid...
```

```
[4]: \n\nPerfect condition! Had to replace my fir...
```

Cleaning the Reviews

```
In [23]: cleaned_reviews = reviews_df['reviews']
```

Removing the '\n' from the reviews

```
In [24]: reviews_df['reviews'] = reviews_df['reviews'].apply(lambda x : x.strip('\n'))
```

```
In [25]: reviews_df.head()
```

```
Out[25]: reviews
```

```
[0]: I'm probably the only person on earth who ha...
```

```
[1]: I got this book, Harry Potter and the Sorcer...
```

```
[2]: This book is a great value! If you can afford...
```

```
[3]: Seen the movies a handful of times but decided...
```

```
[4]: Perfect condition! Had to replace my first o...
```

```
In [26]: reviews_df.shape
```

```
Out[26]: (500, 1)
```

Removing any reviews which is not in English

```
In [27]: reviews_df = reviews_df[reviews_df.reviews.map(lambda x: x.isascii())]
```

```
In [28]: reviews_df.shape
```

```
Out[28]: (465, 1)
```

Converting All Reviews To LowerCase

```
In [29]: to_lower = lambda x: x.lower()
cleaned_reviews = cleaned_reviews.apply(to_lower)
```

Removing Punctuations

```
In [30]: #Translate will transform each character according to the Translation table
remove_punctuations = lambda x: x.translate(str.maketrans('!','!',string.punctuation))
#Creating a Translation Table by maketrans-the 1st argument will be replaced with 2nd argument(in this case none)
cleaned_reviews = cleaned_reviews.apply(remove_punctuations)
```

Remove Emojis

```
In [31]: remove_emojis = lambda x: x.encode('ascii', 'ignore').decode('ascii')
cleaned_reviews = cleaned_reviews.apply(remove_emojis)
```

Remove extra whitespaces

```
In [32]: cleaned_reviews = cleaned_reviews.str.strip()
```

Removing Stop words

```
In [33]: cleaned_reviews
```

```
Out[33]: 0 I'm probably the only person on earth who ha...
```

```
[1]: I got this book, Harry Potter and the Sorcer...
```

```
[2]: This book is a great value! If you can afford...
```

```
[3]: Seen the movies a handful of times but decided...
```

```
[4]: Perfect condition! Had to replace my first o...
```

```
In [34]: #getting stop words from nltk module
nltk.download('stopwords')
Stop_words = stopwords.words("english") #list of stopwords
# removing stopwords
remove_words = lambda x: ' '.join([word for word in x.split() if word not in Stop_words])
cleaned_reviews = cleaned_reviews.apply(remove_words)
```

```
[nltk_data] Downloading package stopwords to C:\Users\HP\AppData\Roaming\nltk_data...
```

```
[nltk_data] Package stopwords is already up-to-date!
```

```
In [35]: cleaned_reviews
```

```
Out[35]: 0 I'm probably the only person on earth never read harry pott...
```

```
[1]: I got this book, Harry Potter and the Sorcer...
```

```
[2]: seen movies a handful of times but decided...
```

```
[3]: perfect condition replace first one since stol...
```

```
[4]: recipient is delighted seeing the film first ...
```

```
[495]: awesome book great ending finally started...
```

```
[497]: havent read these since i was a kid probably around 9 decide...
```

```
[498]: condition book amazing story addictive captiva...
```

```
[499]: awesome
```

```
Name: reviews, Length: 500, dtype: object
```

```
In [20]: len(reviews_list)
```

```
Out[20]: 500
```

Converting the Reviews into a Dictionary

```
In [21]: Dict_reviews = {'reviews':reviews_list}
reviews_df = pd.DataFrame.from_dict(Dict_reviews)
```

```
In [22]: reviews_df.head()
```

```
Out[22]: reviews
```

```
[0]: \n\nI'm probably the only person on earth who...
```

```
[1]: \n\nI got this book, Harry Potter and the So...
```

```
[2]: \n\nThis book is a great value! If you can a...
```

```
[3]: \n\nSeen the movies a handful of times but decid...
```

```
[4]: \n\nPerfect condition! Had to replace my fir...
```

Cleaning the Reviews

```
In [23]: cleaned_reviews = reviews_df['reviews']
```

Removing the '\n' from the reviews

```
In [24]: reviews_df['reviews'] = reviews_df['reviews'].apply(lambda x : x.strip('\n'))
```

```
In [25]: reviews_df.head()
```

```
Out[25]: reviews
```

```
[0]: I'm probably the only person on earth who ha...
```

```
[1]: I got this book, Harry Potter and the Sorcer...
```

```
[2]: This book is a great value! If you can afford...
```

```
[3]: Seen the movies a handful of times but decided...
```

```
[4]: Perfect condition! Had to replace my first o...
```

```
In [26]: reviews_df.shape
```

```
Out[26]: (500, 1)
```

Removing any reviews which is not in English

```
In [27]: reviews_df = reviews_df[reviews_df.reviews.map(lambda x: x.isascii())]
```

```
In [28]: reviews_df.shape
```

```
Out[28]: (465, 1)
```

Converting All Reviews To LowerCase

```
In [29]: to_lower = lambda x: x.lower()
cleaned_reviews = cleaned_reviews.apply(to_lower)
```

Removing Punctuations

```
In [30]: #Translate will transform each character according to the Translation table
remove_punctuations = lambda x: x.translate(str.maketrans('!','!',string.punctuation))
#Creating a Translation Table by maketrans-the 1st argument will be replaced with 2nd argument(in this case none)
cleaned_reviews = cleaned_reviews.apply(remove_punctuations)
```

Remove Emojis

```
In [31]: remove_emojis = lambda x: x.encode('ascii', 'ignore').decode('ascii')
cleaned_reviews = cleaned_reviews.apply(remove_emojis)
```

Remove extra whitespaces

```
In [32]: cleaned_reviews = cleaned_reviews.str.strip()
```

Removing Stop words

```
In [33]: cleaned_reviews
```

```
Out[33]: 0 I'm probably the only person on earth who ha...
```

```
[1]: I got this book, Harry Potter and the Sorcer...
```

```
[2]: This book is a great value! If you can afford...
```

```
[3]: Seen the movies a handful of times but decided...
```

```
[4]: Perfect condition! Had to replace my first o...
```

```
In [34]: #getting stop words from nltk module
nltk.download('stopwords')
Stop_words = stopwords.words("english") #list of stopwords
# removing stopwords
remove_words = lambda x: ' '.join([word for word in x.split() if word not in Stop_words])
cleaned_reviews = cleaned_reviews.apply(remove_words)
```

```
[nltk_data] Downloading package stopwords to C:\Users\HP\AppData\Roaming\nltk_data...
```

```
[nltk_data] Package stopwords is already up-to-date!
```

```
In [35]: cleaned_reviews
```

```
Out[35]: 0 I'm probably the only person on earth never read harry pott...
```

```
[1]: I got this book, Harry Potter and the Sorcer...
```