# Aspect/Modifier Classification Analysis

## Project Links   ¶

Below is the link for the GitHub project page.

Github link (https://github.com/ddey117/ABSA_Project_4)

Import research papers for developement of parser logic. Includes pdf links for spaCy research paper as well as VADER sentiment intensity analyzer. Much work has been done on aspect based sentiment analysis. Please feel free to check out some previous work in the link below.

Research Papers (https://github.com/ddey117/ABSA_Project_4/tree/main/research_papers)

For another way to navigate of my overall project, please feel free to check out a HTML version of my project overview at the link below. There is also a link in this directory to see an example of what exactly a Turk worker was looking at when they were labeling data for this project.

html project directory (https://github.com/ddey117/ABSA_Project_4/tree/main/html)

## Overview

```
The target for this project is an established e-commerce business w
ith a large amount of review data, such as Amazon.com or other onli
ne retailers. The goal of this project is to take advantage of tech
nology and models provided by Spacy combined with a pretrained sent
iment intensity classifier provided by the NLTK toolkit in order to
perform more fine grained sentiment analysis at scale in an efficie
nt manner. This project takes advantage of the parsing and part of
speech tagging capabilites of Spacy's pipeline in order to extract
 aspect/opinion/sentiment triplets. After the aspects are identifie
d, they can be grouped using unsupervised machine learning clusteri
ng techniques; in this case k-means clustering for model speed and
 simplicity. The buisness can use the finished product to quickly t
ransform a large amount of informal review data (text data from rev
iews that may ramble for pages) and transform it into helpful graph
s in order to tune into a small number of categories and help funne
l resources into areas where they are most needed. Amazon Turk was
 taken advantage of to crowd source human labels to analyze the per
formance of the model.
```

Data Exploration Notebook

Author: Dylan Dey

The Author can reached by email: ddey2985@gmail.com (mailto:ddey2985@gmail.com)

## Buisness Problem

**Aspect Category**     **Sentiment**



Customer Service



Value



quality



design

Sentiment analysis involves computationally identifying and categorizing the sentiment expressed by an author in a body of text. It has a wide range of applications in industry from stock speculation using sentiment expressed in news and blogs, to identifying customer satisfaction from their reviews and social media posts.

Today, most e-commerce website designs include a section where their customers can post reviews for products or services. Customers are free to write how they feel about fine grained aspects of a product at length. From a business perspective, very valuable information can be extracted from this section, such as customers' opinion on a product, understanding of a product, etc..

On Amazon.com the rating can be between 1 and 5 where 1 is the worst and 5 is the best. A customer can leave as lengthy of a review as they wish about a product to explain why a given rating was posted.  For example, a customer may give a product a low rating because they didn't like someone they spoke to in customer service but liked everything else about the product.  In typical sentiment analysis, these kinds of nuances would be missed since it could only be determined  if the overall body of the review contained positive, neutral, or negative sentiment. Valuable information would be left on the table.

There is potentially a disconnect from the amazon review ratings, and the overall sentiment of the body text explaining the review, especially if you begin to break down the text into smaller aspects. Thus, Aspect Based Sentiment Analysis (ABSA) was chosen to see if a deeper understanding of each product can be gained by breaking down each review into aspect categories to be paired with predicted sentiment, which will then be compared with the overall rating (1-5).

It is often difficult to efficiently get useful data from a large collection of text data. A lot of e-commerce websites have thousands of reviews and more incoming all of the time. Thousands of reviews with hundreds of words of mostly unhelpful information seems fairly unmanageable to most companies. While the reviews are rather informal, if they are carefully broken down there is information worth saving before generalizing again for efficiency. Aspect Based Sentiment Analysis can transform a messy collection of thousands of informal reviews into a neat and manageable collection of a few aspect categories, in this case 4 different categories using the out of box  Aspect/Opinion/Sentiment Triplet Extractor. Each category will have an associated degree of sentiment related to it, and therefore graphics can easily be prepared and presented to digest more precisely what it is that customers do and do not like about a product in a quickly digestible format in real time. By breaking it down into these categories, say for example Product Design, Value, Quality, and Customer Support, the mass of text data has now been transformed into a numerical representation of sentiment towards broad categories of a product that can be directly improved upon by the company. If a product scores very high sentiment for value and design b

```
                                      ...y...a...product...scores...very...high...sentiments...for...value...and...design...b
                   ut lower scores for customer support, then a company knows it doesn
                   t need to invest more money into improving the product and actually
                   needs to focus on improving how its forward facing employees intera
                   ct with customers.
```

# The Data

Helpful links:

ReadMe file for Amazon Product Reviews (https://s3.amazonaws.com/amazon-reviews-pds/readme.html) | MetaData (https://s3.amazonaws.com/amazon-reviews-pds/tsv/index.txt)

The Amazon Customer Reviews (Product Reviews) contains over 130+ million customer reviews available to researchers in TSV files in the amazon-reviews-pds S3 bucket in AWS US East Region, as per the provided readme file. The reviews were collected from 1995 to 2015. See the provided link for associated metadata. This project focuses on the dataset given by pulling "https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon_reviews_us_Electronics_v1_00.tsv.gz" (https://s3.amazonaws.com/amazon-reviews-pds/tsv/amazon_reviews_us_Electronics_v1_00.tsv.gz%E2%80%9D) from the S3 bucket.

Product_id "B0001FTVEK" (https://www.amazon.com/Sennheiser-RS120-Wireless-Headphones-Charging/dp/B0001FTVEK) was chosen to showcase the triplet extractor as it had a large amount of verified reviews and a pair of headphones seemed like a reasonable choice for aspect based sentiment analysis.

### *Clean Data*

Text data trends towards exponential growth with increasing dataset size. Therefore, text cleaning and preprocessing was a major considersation of this project. Please refer to my Text Preprocessing Toolset (https://github.com/ddey117/preprocess_ddey117) that I created to use for this and other projects that involve text data preprocessing.

### Unlabeled Data Created Through Unsupervised Learning

This project showcases an out of box product for extracting opinon/aspect/sentiment triplets from a large amount of messy text data and converting it into a neat set of categories for analysis. To do this, however, it takes advantage of some simple clustering techniques from the sklearn cluster library. For this project, kmeans clustering was chosen for speed and simplicity. Error analysis will be discussed later in more detail in regards to how the model performs with clustering the reviews appropriately into categories and what issues it may run into when parsing internet language. Error analysis for the SentimentIntensityClassifier (https://www.nltk.org/howto/sentiment.html) offered by the Natural Language ToolKit (NLTK library) will be tested against this 'newly' generated data from my unsupervised learning will be performed by comparing to a seperate set of hand labeled aspect/modifier pairs by humans in an expiremental setting.

### *experimental setup*

Using the following Turk_Form_HTML (html/Turk_Instructions.html) I crowdsourced some labels from humans using Amazon Mechanical Turk to compare to my model using the SentimentIntensityAnalyzer for each aspect/modifier pair extracted from the Amazon reviews. Amazon Mechanical Turk works by quickly dispersing large amounts of data to a large number of people in order to complete simple tasks for a reward. This experiment was set up to reward a penny for each aspect/modifier pair labeled for sentiment from very negative to very positive with an option for NA from a drop down menu (see html above for reference). In total, 410 workers submitted 6107 non-null aspect/opinion pairs for sentiment intensity pertaining to 1438 unique aspects. Duplicate pairs of aspect/opinion pairs were included to inspect variance of submission from human labels and machine labels for each opinion pair. No qualifications or screening was put in place before the workers were chosen, but I did review sections of the data and accept or reject what seemed reasonable.

All labels were generated using my triplet extractor on the dataset describing Product_id "B0001FTVEK" (https://www.amazon.com/Sennheiser-RS120-Wireless-Headphones-Charging/dp/B0001FTVEK) and randomized for different aspect/modifier pairs before sending out to humans for rating for sentiment.

---

Amazon Product Reviews ReadME (https://s3.amazonaws.com/amazon-reviews-pds/readme.html) | Amazon Product Reviews MetaData (https://s3.amazonaws.com/amazon-reviews-pds/tsv/index.txt) | Sennheiser-RS120-Wireless-Headphones (https://www.amazon.com/Sennheiser-RS120-Wireless-Headphones-Charging/dp/B0001FTVEK)

```
In [1]:
 1  #import necessary libraries
 2
 3  import pandas as pd
 4  import numpy as np
 5  # import dataframe_image as dfi
 6  import string
 7  import re
 8  from matplotlib import pyplot as plt
 9  import seaborn as sns
10  from nltk.corpus import stopwords
11  from nltk import FreqDist, word_tokenize
12  # from nltk.tokenize import TweetTokenizer
13  from nltk.stem import WordNetLemmatizer
14  # from wordcloud import WordCloud, STOPWORDS
15  import re,string
16  import unidecode
17  import html
18
19  import preprocess_ddey117 as pp
20
21
22  import requests
23  import os
24  import csv
25  import urllib.request
26  import gzip
27  import sys
28  import spacy
29  import json
30  # import boto3
31  # from boto.s3.connection import S3Connection
32
33  from collections import defaultdict
34  from sklearn import cluster
35  import seaborn as sns
36
37  import nltk
38  # nltk.download('vader_lexicon')
39
40  import spacy
41  nlp = spacy.load("en_core_web_lg")
42
43  from nltk.sentiment.vader import SentimentIntensityAnalyzer
44  sid = SentimentIntensityAnalyzer()
45
46  # from nltk.collocations import *
47  # bigram_measures = nltk.collocations.BigramAssocMeasures()
48  # trigram_measures = nltk.collocations.TrigramAssocMeasures()
49  # fourgram_measures = nltk.collocations.QuadgramAssocMeasures()
```

```
/Users/dylandey/anaconda3/envs/learn-env/lib/python3.6/site-packages/
tensorflow/python/framework/dtypes.py:517: FutureWarning: Passing (ty
pe, 1) or '1type' as a synonym of type is deprecated; in a future ver
sion of numpy, it will be understood as (type, (1,)) / '(1,)type'.
  _np_qint8 = np.dtype([("qint8", np.int8, 1)])
```

```
/Users/dylandey/anaconda3/envs/learn-env/lib/python3.6/site-packages/
tensorflow/python/framework/dtypes.py:518: FutureWarning: Passing (ty
pe, 1) or '1type' as a synonym of type is deprecated; in a future ver
sion of numpy, it will be understood as (type, (1,)) / '(1,)type'.
  _np_quint8 = np.dtype([("quint8", np.uint8, 1)])
/Users/dylandey/anaconda3/envs/learn-env/lib/python3.6/site-packages/
tensorflow/python/framework/dtypes.py:519: FutureWarning: Passing (ty
pe, 1) or '1type' as a synonym of type is deprecated; in a future ver
sion of numpy, it will be understood as (type, (1,)) / '(1,)type'.
  _np_qint16 = np.dtype([("qint16", np.int16, 1)])
/Users/dylandey/anaconda3/envs/learn-env/lib/python3.6/site-packages/
tensorflow/python/framework/dtypes.py:520: FutureWarning: Passing (ty
pe, 1) or '1type' as a synonym of type is deprecated; in a future ver
sion of numpy, it will be understood as (type, (1,)) / '(1,)type'.
```

A large collection of amazon reviews that fall under the "electronics" category. For this project, product_id "B0001FTVEK" (https://www.amazon.com/Sennheiser-RS120-Wireless-Headphones-Charging/dp/B0001FTVEK) was chosen as it had a large amount of verified reviews and a pair of headphones seemed like a reasonable choice for aspect based sentiment analysis.

In [2]:
```
1  # df = pd.read_csv('data/df_electronics.tsv', sep='\t')
2  # df.groupby('product_id').count().sort_values(by='star_rating').tail(2
```

In [4]:
```
1  df = pd.read_csv('data/df_electronics.tsv', sep='\t')
2  df.groupby('product_id').count().sort_values(by='star_rating').tail(20)
```

Out[4]:

| product_id | marketplace | customer_id | review_id | product_parent | product_title | product_category | star_rating | helpful_votes | total_votes | vine | verified_purcl |
|---|---|---|---|---|---|---|---|---|---|---|---|
| B004LTEUDO | 3997 | 3997 | 3997 | 3997 | 3997 | 3997 | 3997 | 3997 | 3997 | 3997 | |
| B004HHICKC | 4213 | 4213 | 4213 | 4213 | 4213 | 4213 | 4213 | 4213 | 4213 | 4213 | |
| B0 click to scroll output; double click to hide 3 | | 4773 | 4773 | 4773 | 4773 | 4773 | 4773 | 4773 | 4773 | 4773 | |
| B001TH7GSW | 4866 | 4866 | 4866 | 4866 | 4866 | 4866 | 4866 | 4866 | 4866 | 4866 | |
| B008KVUAGU | 5015 | 5015 | 5015 | 5015 | 5015 | 5015 | 5015 | 5015 | 5015 | 5015 | |
| B003WGRUQQ | 5072 | 5072 | 5072 | 5072 | 5072 | 5072 | 5072 | 5072 | 5072 | 5072 | |
| B002MAPT7U | 5295 | 5295 | 5295 | 5295 | 5295 | 5295 | 5295 | 5295 | 5295 | 5295 | |
| B001GTT0VO | 5580 | 5580 | 5580 | 5580 | 5580 | 5580 | 5580 | 5580 | 5580 | 5580 | |
| B0052SCU8U | 5756 | 5756 | 5756 | 5756 | 5756 | 5756 | 5756 | 5756 | 5756 | 5756 | |
| B00316263Y | 5813 | 5813 | 5813 | 5813 | 5813 | 5813 | 5813 | 5813 | 5813 | 5813 | |
| B00D5Q75RC | 6062 | 6062 | 6062 | 6062 | 6062 | 6062 | 6062 | 6062 | 6062 | 6062 | |
| B004QK7HI8 | 6536 | 6536 | 6536 | 6536 | 6536 | 6536 | 6536 | 6536 | 6536 | 6536 | |
| B00F5NE2KG | 6688 | 6688 | 6688 | 6688 | 6688 | 6688 | 6688 | 6688 | 6688 | 6688 | |
| B0019EHU8G | 7586 | 7586 | 7586 | 7586 | 7586 | 7586 | 7586 | 7586 | 7586 | 7586 | |
| B000WYVBR0 | 7835 | 7835 | 7835 | 7835 | 7835 | 7835 | 7835 | 7835 | 7835 | 7835 | |
| B0001FTVEK | 8793 | 8793 | 8793 | 8793 | 8793 | 8793 | 8793 | 8793 | 8793 | 8793 | |
| B0012S4APK | 9359 | 9359 | 9359 | 9359 | 9359 | 9359 | 9359 | 9359 | 9359 | 9359 | |
| B003EM8008 | 9766 | 9766 | 9766 | 9766 | 9766 | 9766 | 9766 | 9766 | 9766 | 9766 | |
| B0002L5R78 | 11166 | 11166 | 11166 | 11166 | 11166 | 11166 | 11166 | 11166 | 11166 | 11166 | 1 |
| B003L1ZYYM | 15334 | 15334 | 15334 | 15334 | 15334 | 15334 | 15334 | 15334 | 15334 | 15334 | 15 |

In [3]:

```python
#testing other dataframes not shown in this notebook



# df_hp1 = df.loc[df['product_id'] == 'B003EM8008'].copy()
# df_hp1.reset_index(drop=True, inplace=True)
# df_hp2 = df.loc[df['product_id'] == 'B0001FTVEK'].copy()
# df_hp2.reset_index(drop=True, inplace=True)
# df_hp3 = df.loc[df['product_id'] == 'B004RKQM8I'].copy()
# df_hp3.reset_index(drop=True, inplace=True)
# df_hp4 = df.loc[df['product_id'] == 'B0038W0K2K'].copy()
# df_hp4.reset_index(drop=True, inplace=True)

# df_sb1 = df.loc[df['product_id'] == 'B00D5Q75RC'].copy()
# df_sb1.reset_index(drop=True, inplace=True)
# df_sb2 = df.loc[df['product_id'] == 'B00F5NE2KG'].copy()
# df_sb2.reset_index(drop=True, inplace=True)

# df_mp1 = df.loc[df['product_id'] == 'B00020S7XK'].copy()
# df_mp1.reset_index(drop=True, inplace=True)
# df_mp2 = df.loc[df['product_id'] == 'B002MAPT7U'].copy()
# df_mp2.reset_index(drop=True, inplace=True)

# #dropping uneccessary columns

# columns_td = ['marketplace', 'customer_id',  'product_id',
#                 'product_parent', 'product_title', 'product_category',
#                 'helpful_votes', 'total_votes', 'vine', 'verified_purch
#                 'review_headline', 'review_date']

# df2.drop(columns=columns_td, inplace=True)

#saving dataframe with chosen product for quick loading time

# df2.to_csv('data/df_electronics_example.csv', index=False)
```

In [4]:
```python
 1  #display important information about dataset
 2  #Almost 9000 reviews, mostly 4 or 5 stars.
 3
 4  df = pd.read_csv("data/df_electronics_example.csv")
 5  display(df.head())
 6  print("\n")
 7  print('star rating value counts')
 8  display(df.star_rating.value_counts())
 9  display(df.star_rating.hist())
10  df.info()
11  print("\n")
12  print("distribution of star rating")
```

| | review_id | star_rating | review_body |
|---|---|---|---|
| **0** | R17U6AU06HR16Q | 5 | Great product. |
| **1** | R31Y01GPXH7P64 | 4 | work great sounds amazing |
| **2** | RG40FO7CNWOG7 | 5 | Works tremendously well. |
| **3** | R2CI6TXIGZR6RU | 5 | THEY WORK GREAT |
| **4** | R37KF8VRBUZNQD | 5 | Works fine |

```
star rating value counts

5    4715
4    1924
1     829
3     724
2     601
Name: star_rating, dtype: int64

<AxesSubplot:>

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8793 entries, 0 to 8792
Data columns (total 3 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   review_id    8793 non-null   object
 1   star_rating  8793 non-null   int64
 2   review_body  8793 non-null   object
dtypes: int64(1), object(2)
memory usage: 206.2+ KB


distribution of star rating
```

# Function Definition

```python
In [5]:  1  #I wrote a library for taking care of many common tasks for data langua
         2  #However, since this project must load spaCy for parsing the aspect/opi
         3  #Some of the major processing will be handled by spaCy
         4
         5  #simple processing such as removing emails, html tags, urls
         6  #accented characters and counting all words in the dataset
         7  #are handled with my own toolset
         8
         9  def master_preprocess(df):
        10      df['wordcounts'] = df['review_body'].apply(lambda x: pp.get_wordcou
        11      df['review_body'] = df['review_body'].apply(lambda x: pp.remove_ema
        12      df['review_body'] = df['review_body'].apply(lambda x: pp.remove_url
        13      df['review_body'] = df['review_body'].apply(lambda x: pp.remove_htm
        14      df['review_body'] = df['review_body'].apply(lambda x: pp.remove_acc
        15      return df
```

```python
In [6]:  1  df = master_preprocess(df)
```

```
In [7]:    1  #Looking at the statistics for wordcounts below
           2  #The median length for a amazon review
           3  #for this particular set of headphones is about 50 words
           4  #ranging from 1 to 1565 words for each review
           5  #Thus showing the wide range of informal data that is
           6  #allowed in the review body and how unmanageable it can
           7  #become
           8  #With this dataset the total amount of words
           9  #has already reached over half of a million words
          10
          11
          12  display(df.describe())
          13  df.wordcounts.sum()
```

|       | star_rating | wordcounts   |
|-------|-------------|--------------|
| count | 8793.000000 | 8793.000000  |
| mean  | 4.034346    | 74.089730    |
| std   | 1.318899    | 91.416481    |
| min   | 1.000000    | 1.000000     |
| 25%   | 4.000000    | 23.000000    |
| 50%   | 5.000000    | 47.000000    |
| 75%   | 5.000000    | 92.000000    |
| max   | 5.000000    | 1565.000000  |

Out[7]:  651471

# Explaining The Parser

# MODIFIER/ASPECT/OPINION EXTRACTOR

**01**

input: raw amazon review text

**Step One: Aspect Extraction**

dependency parsing rules

acomp          nsubj

I'm happy about the sound quality.

(happy, sound)

sum up sentiment for each cluster

## sound

**02**

**Step Two: Clustering and Sentiment**

bad sound

not great bass

perfect pitch

**03**

**Step Three: Visualization**

simplify results with visulations

## Clustering and Polarity

A large number of amazon reviews produce a large number of aspect-modifier pairs. These pairs ultimately seemed to diverge to common topics, and therefore it would make sense to use machine learning to automatically figure out these categories for us. This leads to a better summation of insight from the total pool of customers who were kind enough to leave a review. Polarity scores are also averaged out of every cluster to give a quantifiable explanation to opinion to distinct categories of a given product.

### Word Vectors and Clustering

In order to work with any amazon review data, first the text data must be converted into something that a machine can recognize. The most famous implementation of words vectors is the word2vec project. However, spaCy vectorization was chosen for this projec as it provides fast and easy access to over a million unique word vectors, and its multi-task CNN model is trained on 'web' data and not 'newspaper' data as in other libraries like NLTK.

The word vectors were then grouped using K-Means clustering algorithm in Scikit-Learn. Other clustering algorithms such as DBSCAN were tested. However, K-Means gave optimal results with four clusters. The clusers were labeled with input from a user after suggesting the top most common word for each cluster.
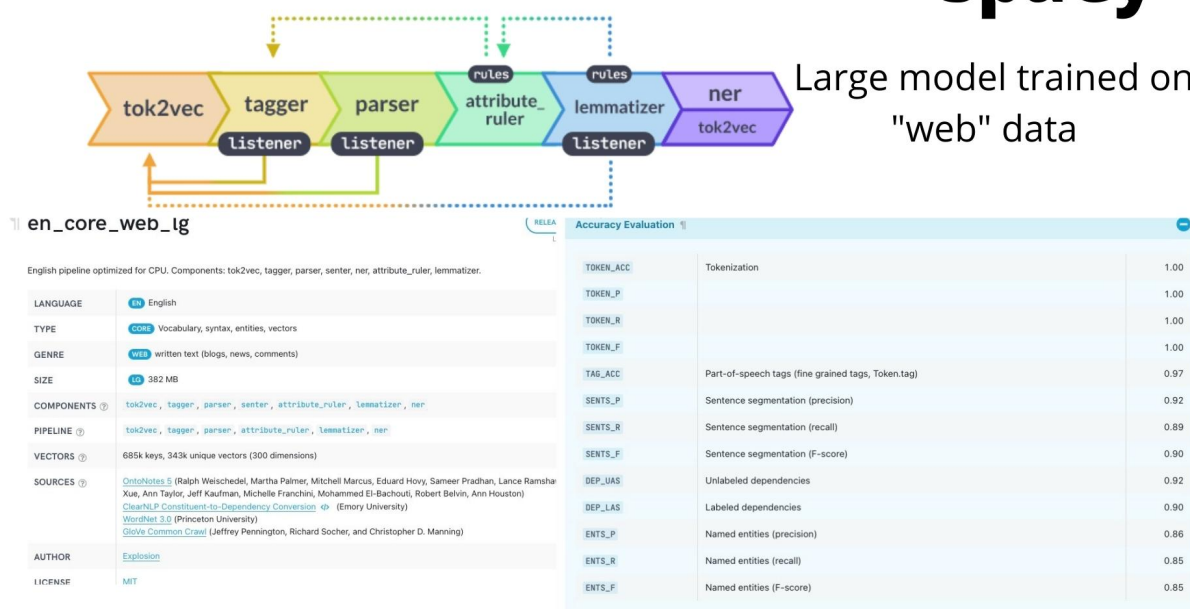
Below Is the pipeline design for spaCy and a description for the size and sources for the model loaded to run this project and parse the amazon reviews.

en_core_web_large (https://spacy.io/models/en#en_core_web_lg)

## CPU pipeline design

# spaCy

Large model trained on "web" data

### en_core_web_lg

English pipeline optimized for CPU. Components: tok2vec, tagger, parser, senter, ner, attribute_ruler, lemmatizer.

| LANGUAGE | EN English |
|---|---|
| TYPE | CORE Vocabulary, syntax, entities, vectors |
| GENRE | WEB written text (blogs, news, comments) |
| SIZE | LG 382 MB |
| COMPONENTS | tok2vec, tagger, parser, senter, attribute_ruler, lemmatizer, ner |
| PIPELINE | tok2vec, tagger, parser, attribute_ruler, lemmatizer, ner |
| VECTORS | 685k keys, 343k unique vectors (300 dimensions) |
| SOURCES | OntoNotes 5 (Ralph Weischedel, Martha Palmer, Mitchell Marcus, Eduard Hovy, Sameer Pradhan, Lance Ramsha Xue, Ann Taylor, Jeff Kaufman, Michelle Franchini, Mohammed El-Bachouti, Robert Belvin, Ann Houston) |
| | ClearNLP Constituent-to-Dependency Conversion (Emory University) |
| | WordNet 3.0 (Princeton University) |
| | GloVe Common Crawl (Jeffrey Pennington, Richard Socher, and Christopher D. Manning) |
| AUTHOR | Explosion |
| LICENSE | MIT |

| Accuracy Evaluation | | |
|---|---|---|
| TOKEN_ACC | Tokenization | 1.00 |
| TOKEN_P | | 1.00 |
| TOKEN_R | | 1.00 |
| TOKEN_F | | 1.00 |
| TAG_ACC | Part-of-speech tags (fine grained tags, Token.tag) | 0.97 |
| SENTS_P | Sentence segmentation (precision) | 0.92 |
| SENTS_R | Sentence segmentation (recall) | 0.89 |
| SENTS_F | Sentence segmentation (F-score) | 0.90 |
| DEP_UAS | Unlabeled dependencies | 0.92 |
| DEP_LAS | Labeled dependencies | 0.90 |
| ENTS_P | Named entities (precision) | 0.86 |
| ENTS_R | Named entities (recall) | 0.85 |
| ENTS_F | Named entities (F-score) | 0.85 |

"spaCy uses the terms **head** and **child** to describe the words connected by a single arc in the

dependency tree. The term **dep** is used for the arc label, which describes the type of syntactic relation that connects the child to the head. As with other attributes, the value of .dep is a hash value. You can get the string value with .dep_." Navigating The Parse Tree (https://spacy.io/usage/linguistic-features#navigating)

**First Rule of Dependency Parser:** The Aspect (A) token is a subject noun with a child modifier (M) that has a relation of amod (adjectival modifier). This just means that the aspect and opinion share a simple adjective/noun relationship that can be extracted. However, there are certain caveats that need to be kept in mind when parsing the tree for this rule.

- First, it is important to check to see if there is an additional adverbial modifier that could adjust the intensity of the sentiment implied by the adjective and adverb combination in regards to the subject/aspect. This is important to keep in mind as we are taking advantage of NLTK vader sentiment intensity analyzer which can make use of additional adverbs to get a better understanding of sentiment.

- Another important thing to keep in mind when parsing for this rule is to be aware of the possibility of negating the adjective with 'no' as a determiner.

**First Rule Examples**

**Example1:** The comfortable headphones.

**Example2:** The most comfortable headphones.

**Example3:** No comfortable features.

- det = determiner
- A = aspect
- M = modifier
- amod = adjectival modifier

**Second Rule of Dependency Parser:** The aspect (A) is a child of something with a relation of nominal subject (nsubj.) while the modifier (M) is a child of the same something with a relationship of direct object. In this case, the adjective would be acting as the determiner of the clause. For simplicity's sake, it was determined to assume that each verb will have only one NSUBJ and DOBJ. This is a fair assumption for the application of this project, because even if there are multiple subjects, they will both be reviewing the same thing and will likely share the same opinion as it is written as a single review. For example, if an author were to say "My wife and I bought the awesome headphones", we still only want to extract the keywords 'awesome' and 'headphones.' If this sounds confusing, hopefully the example below will help clarify.

**Second Rule Example**

**Example:** I bought the awesome headphones.

- nsubj = nominal subject
- dobj =headphones
- det= awesome

**Third Rule of Dependency Parser:** The modifier (M) is a child of something with a relation of an adjectival complement (acomp), while the aspect (A) is a child of that same something with a relation of nominal subject (nsubj).

- This rule needs to handle special cases in which the child is tagged as a modal verb with an auxiliary dependency. This would flag for phrases such as "the sound of the speakers could be better." For special cases like this, the parser will add a negative prefix before scoring the aspect/modifier pairs for sentiment.

**Third Rule Examples**

**Example1:** Barb is happy about the sound quality.

**Example2:** This could be better.

Example2 would be extracted as A= "this" and M= "not better"

- A = aspect
- M = modifier

**Fourth Rule of Dependency Parser:** The aspect (A) is a child of something with a relationship of passive nominal subject (nsubjpass) while the modifier (M) is a child of that same something with a relationship of adverbial modifier (advmod). In other words, the modifier is an adverbial modifier to a passive verb.

- nsubjpass: A passive nominal subject is a noun phrase which is the syntactic subject of a passive clause.
- This step of the parser will also check to add a negative prefix before extracting and scoring for sentiment if necessary

**Fourth Rule Examples**

**Example1:** The headphones died quickly.

- A = aspect
- M = modifier

**Fifth Rule of Dependency Parser:** The aspect (A) is a child of the modifier with a relationship

of nominal subject, while the modifier has a child with a relation of copula(cop). Here the parser is looking for the complement of a copular verb. An often used copula verb is the word "is," as in the phrase "Bill is big."

- Assumption - A verb will have only one NSUBJ and DOBJ
- cop: copula A copula is the relation between the complement of a copular verb and the copular verb. (We normally take a copula as a dependent of its complement.

**Fifth Rule Example**

**Example1:** The sound is awesome.

- A = aspect
- M = modifier

**Sixth Rule of Dependency Parser:** Aspect/modifier are children of an interjection

- NTJ (interjections like bravo, great etc)

**Sixth Rule Example**

**Example1:** Bravo, headphones.

- A = aspect
- M = modifier

**Seventh Rule of Dependency Parser:** This rule is similar to rule 5, but makes use of the attr (attribute) tag instead. It seems to function similarly, in which an attribute is considered a noun phrase following a copular verb

- ATTR - link between a verb like 'is/seem/appear/became' and its complement

**Seventh Rule Example**

**Example1:** This is garbage.

- A = aspect
- M = modifier

**For all Parsing:** SpaCy has a large library of named entities it can recoginize and tag. This logic is added for each step in the model.

In [8]:
```python
1  spacy.explain('acomp')
```

Out[8]: `'adjectival complement'`

Please feel free to review the following sections above under the spaCy documentation for pos_tagging if you would like to get an understanding of how the parser was designed in spaCy. Each link should be a direct link to the appropriate topic.

Dependecy Parsing with spaCy (https://spacy.io/usage/linguistic-features#dependency-parse)

Navigating The Tree (https://spacy.io/usage/linguistic-features#navigating)

Named Entity Recognition (https://spacy.io/usage/linguistic-features#named-entities)

Hutto, C.J. & Gilbert, E.E. (2014). VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. Eighth International Conference on Weblogs and Social Media (ICWSM-14). Ann Arbor, MI, June 2014.

VADER sentiment (https://www.researchgate.net/publication/275828927_VADER_A_Parsimonious_Rule-based_Model_for_Sentiment_Analysis_of_Social_Media_Text)

The reasearch paper was published on release of the VADER intensity sentiment analyzer. Please feel free to read to get a better understanding of how this tool was developed before being taken advantage of in this project.

## Detecting Product Aspects

In [9]:
```python
1  display(spacy.explain('nsubjpass'))
2  display(spacy.explain('cop'))
3  display(spacy.explain('INTJ'))
4  spacy.explain('attr')
```

`'nominal subject (passive)'`

`'copula'`

`'interjection'`

Out[9]: `'attribute'`

TABLE OF REFERENCE:
**AMOD**
adjectival modifier
**ADVMOD**
adverbial modifier
example: Genetically Modified Food, Less often
**NSUBJ**

"Nominal subject (nsubj) is a nominal which is the syntactic subject and the proto-agent of a clause. That is, it is in the position that passes typical grammatical test for subjecthood, and this argument is the more agentive, the do-er, or the proto-agent of the clause. This nominal may be headed by a noun, or it may be a pronoun or relative pronoun or, in ellipsis contexts, other things such as an adjective." Taken from the documentation.

example: Genetically Modified Food, Less often

**DOBJ**

The direct object of a VP is the noun phrase which is the (accusative) object of the verb

**DET**

Determiner. "The English DET covers most cases of Penn Treebank DT, PDT, WDT. However, when a Penn Treebank word with one of these tags stands alone as a noun phrase rather than modifying another word, then it becomes PRON." Taken from the documentation.

**ACOMP**

Adjective complement. A phrase that modifies an adjective.

**cop**

"A cop (copula) is the relation of a function word used to link a subject to a nonverbal predicate, including the expression of identity predication (e.g. sentences like "Kim is the President"). It is often a verb but nonverbal (pronominal) copulas are also frequent in the world's languages. Verbal copulas are tagged AUX, not VERB. Pronominal copulas are tagged PRON or DET." From the documentation.

**INTJ**

interjection. An interjection is a word that is used most often as an exclamation or part of an exclamation.

In [10]:
```python
def apply_extraction(row,nlp=nlp,sid=sid):
    review_body = row['review_body']
#     review_id = row['review_id']
#     review_marketplace = row['marketplace']
#     customer_id = row['customer_id']
#     product_id = row['product_id']
#     product_parent = row['product_parent']
#     product_title = row['product_title']
#     product_category = row['product_category']
#     date = str(row['review_date'])
#     star_rating = row['star_rating']
#     url = add_amazonlink(product_id)



    doc=nlp(review_body)


    ## FIRST RULE OF DEPENDANCY PARSE -
    ## M - Sentiment modifier || A - Aspect
    ## RULE = M is child of A with a relationshion of amod(adjectival
    ner_heads = {ent.root.idx: ent for ent in doc.ents}
    rule1_pairs = []
    for token in doc:
        A = "999999"
        M = "999999"
        if token.dep_ == "amod" and not token.is_stop:
            M = token.text
            if token.head in ner_heads:
                A = ner_heads[token.head].text
            else:
                A = token.head.text

            # add adverbial modifier of adjective (e.g. '*most* comfor
            M_children = token.children
            for child_m in M_children:
                if(child_m.dep_ == "advmod"):
                    M_hash = child_m.text
                    M = M_hash + " " + M
                    break

            # negation in adjective, the "no" keyword is a 'det' of th
            A_children = token.head.children
            for child_a in A_children:
                if(child_a.dep_ == "det" and child_a.text == 'no'):
                    neg_prefix = 'not'
                    M = neg_prefix + " " + M
                    break

        if(A != "999999" and M != "999999"):
            rule1_pairs.append((A, M,sid.polarity_scores(token.text)['

    ## SECOND RULE OF DEPENDANCY PARSE -
    ## M - Sentiment modifier || A - Aspect
    #Direct Object - A is a child of something with relationship of ns
    # M is a child of the same something with relationship of dobj
```

```python
57         #Assumption - A verb will have only one NSUBJ and DOBJ
58  #       ner_heads = {ent.root.idx: ent for ent in doc.ents}
59       rule2_pairs = []
60       for token in doc:
61           children = token.children
62           A = "999999"
63           M = "999999"
64           add_neg_pfx = False
65           for child in children :
66               if(child.dep_ == "nsubj" and not child.is_stop):
67                   if child.idx in ner_heads:
68                       A = ner_heads[child.idx].text
69                   else:
70                       A = child.text
71                   # check_spelling(child.text)
72
73               if((child.dep_ == "dobj" and child.pos_ == "ADJ") and not
74                   M = child.text
75                   #check_spelling(child.text)
76
77               if(child.dep_ == "neg"):
78                   neg_prefix = child.text
79                   add_neg_pfx = True
80
81       if (add_neg_pfx and M != "999999"):
82          M = neg_prefix + " " + M
83
84          if(A != "999999" and M != "999999"):
85              rule2_pairs.append((A, M,sid.polarity_scores(M)['compound'
86
87
88       ## THIRD RULE OF DEPENDANCY PARSE -
89       ## M - Sentiment modifier || A - Aspect
90       ## Adjectival Complement - A is a child of something with relation
91       ## M is a child of the same something with relationship of acomp
92       ## Assumption - A verb will have only one NSUBJ and DOBJ
93       ## "The sound of the speakers would be better. The sound of the sp
94
95
96       rule3_pairs = []
97
98       for token in doc:
99
100          children = token.children
101          A = "999999"
102          M = "999999"
103          add_neg_pfx = False
104          for child in children :
105              if(child.dep_ == "nsubj" and not child.is_stop):
106                  if child.idx in ner_heads:
107                      A = ner_heads[child.idx].text
108                  else:
109                      A = child.text
110                  # check_spelling(child.text)
111
112              if(child.dep_ == "acomp" and not child.is_stop):
113                  M = child.text
```

```python
114
115                # example - 'this could have been better' -> (this, not be
116            if(child.dep_ == "aux" and child.tag_ == "MD"):
117                neg_prefix = "not"
118                add_neg_pfx = True
119
120            if(child.dep_ == "neg"):
121                neg_prefix = child.text
122                add_neg_pfx = True
123
124        if (add_neg_pfx and M != "999999"):
125            M = neg_prefix + " " + M
126                #check_spelling(child.text)
127
128        if(A != "999999" and M != "999999"):
129            rule3_pairs.append((A, M, sid.polarity_scores(M)['compound
130
131    ## FOURTH RULE OF DEPENDENCY PARSE -
132    ## M - Sentiment modifier || A - Aspect
133
134    #Adverbial modifier to a passive verb - A is a child of something
135    # M is a child of the same something with relationship of advmod
136
137    #Assumption - A verb will have only one NSUBJ and DOBJ
138
139    rule4_pairs = []
140    for token in doc:
141
142
143        children = token.children
144        A = "999999"
145        M = "999999"
146        add_neg_pfx = False
147        for child in children :
148            if((child.dep_ == "nsubjpass" or child.dep_ == "nsubj") an
149                if child.idx in ner_heads:
150                    A = ner_heads[child.idx].text
151                else:
152                    A = child.text
153                # check_spelling(child.text)
154
155            if(child.dep_ == "advmod" and not child.is_stop):
156                M = child.text
157                M_children = child.children
158                for child_m in M_children:
159                    if(child_m.dep_ == "advmod"):
160                        M_hash = child_m.text
161                        M = M_hash + " " + child.text
162                        break
163                #check_spelling(child.text)
164
165            if(child.dep_ == "neg"):
166                neg_prefix = child.text
167                add_neg_pfx = True
168
169        if (add_neg_pfx and M != "999999"):
170            M = neg_prefix + " " + M
```

```python
171
172            if(A != "999999" and M != "999999"):
173                rule4_pairs.append((A, M,sid.polarity_scores(M)['compound'
174
175
176        ## FIFTH RULE OF DEPENDANCY PARSE -
177        ## M - Sentiment modifier || A - Aspect
178
179        #Complement of a copular verb - A is a child of M with relationshi
180        # M has a child with relationship of cop
181
182        #Assumption - A verb will have only one NSUBJ and DOBJ
183
184        rule5_pairs = []
185        for token in doc:
186            children = token.children
187            A = "999999"
188            buf_var = "999999"
189            for child in children :
190                if(child.dep_ == "nsubj" and not child.is_stop):
191                    if child.idx in ner_heads:
192                        A = ner_heads[child.idx].text
193                    else:
194                        A = child.text
195
196                    # check_spelling(child.text)
197
198                if(child.dep_ == "cop" and not child.is_stop):
199                    buf_var = child.text
200                    #check_spelling(child.text)
201
202            if(A != "999999" and buf_var != "999999"):
203                rule5_pairs.append((A, token.text,sid.polarity_scores(toke
204
205
206        ## SIXTH RULE OF DEPENDENCY PARSE -
207        ## M - Sentiment modifier || A - Aspect
208        ## INTJ (interjections like bravo, great etc)
209
210
211        rule6_pairs = []
212        for token in doc:
213            children = token.children
214            A = "999999"
215            M = "999999"
216            if(token.pos_ == "INTJ" and not token.is_stop):
217                for child in children :
218                    if(child.dep_ == "nsubj" and not child.is_stop):
219                        M = token.text
220                        if child.idx in ner_heads:
221                            A = ner_heads[child.idx].text
222                        else:
223                            A = child.text
224                        # check_spelling(child.text)
225
226            if(A != "999999" and M != "999999"):
227                rule6_pairs.append((A, M,sid.polarity_scores(M)['compound'
```

```
228
229
230          ## SEVENTH RULE OF DEPENDENCY PARSE -
231          ## M - Sentiment modifier || A - Aspect
232          ## ATTR - link between a verb like 'be/seem/appear' and its comple
233          ## Example: 'this is garbage' -> (this, garbage)
234
235          rule7_pairs = []
236          for token in doc:
237              children = token.children
238              A = "999999"
239              M = "999999"
240              add_neg_pfx = False
241              for child in children :
242                  if(child.dep_ == "nsubj" and not child.is_stop):
243                      if child.idx in ner_heads:
244                          A = ner_heads[child.idx].text
245                      else:
246                          A = child.text
247                      # check_spelling(child.text)
248
249                  if((child.dep_ == "attr") and not child.is_stop):
250                      M = child.text
251                      #check_spelling(child.text)
252
253                  if(child.dep_ == "neg"):
254                      neg_prefix = child.text
255                      add_neg_pfx = True
256
257              if (add_neg_pfx and M != "999999"):
258                  M = neg_prefix + " " + M
259
260              if(A != "999999" and M != "999999"):
261                  rule7_pairs.append((A, M,sid.polarity_scores(M)['compound'
262
263
264
265          aspects = []
266
267          aspects = rule1_pairs + rule2_pairs + rule3_pairs +rule4_pairs +ru
268          prod_pronouns = ['it', 'this', 'they']
269
270          # replace all instances of "it", "this" and "they" with "product"
271          aspects = [(A,M,P,r) if A not in prod_pronouns else ("product",M,P
272
273 #      dic = {"review_id" : review_id , "aspect_pairs" : aspects, "revi
274 #      , "customer_id" : customer_id, "product_id" : product_id, "produ
275 #      "product_title" : product_title, "product_category" : product_ca
276
277
278          return aspects
279
280
281
282 def remove_digits(x):
283     return " ".join([t for t in x.split() if not t.isdigit()])
284
```

```python
285
286
287   def get_word_vectors(unique_aspects, nlp=nlp):
288       asp_vectors = []
289       for aspect in unique_aspects:
290           # print(aspect)
291           token = nlp(aspect)
292           asp_vectors.append(token.vector)
293       return asp_vectors
294
295
296   def get_aspect_freq_map(aspects):
297       aspect_freq_map = defaultdict(int)
298       for asp in aspects:
299           aspect_freq_map[asp] += 1
300       return aspect_freq_map
301
302
303
304   NUM_CLUSTERS = 4
305
306   def get_word_cluster_labels(unique_aspects, nlp=nlp):
307       # print("Found {} unique aspects for this product".format(len(uniq
308       asp_vectors = get_word_vectors(unique_aspects, nlp)
309       # n_clusters = min(NUM_CLUSTERS,len(unique_aspects))
310       if len(unique_aspects) <= NUM_CLUSTERS:
311           # print("Too few aspects ({}) found. No clustering required...
312           return list(range(len(unique_aspects)))
313
314       # print("Running k-means clustering...")
315       n_clusters = NUM_CLUSTERS
316       kmeans = cluster.KMeans(n_clusters=n_clusters)
317       kmeans.fit(asp_vectors)
318       labels = kmeans.labels_
319       # dbscan = cluster.DBSCAN(eps = 0.2, min_samples = 2).fit(asp_vect
320       # labels = dbscan.labels_
321
322       # print("Finished running k-means clustering with {} labels".forma
323       # print(labels)
324       return labels
325
326
327
328   def get_cluster_names_map(asp_to_cluster_map, aspect_freq_map):
329       cluster_id_to_name_map = defaultdict()
330       # cluster_to_asp_map = defaultdict()
331       clusters = set(asp_to_cluster_map.values())
332       for i in clusters:
333           this_cluster_asp = [k for k,v in asp_to_cluster_map.items() if
334           filt_freq_map = {k:v for k,v in aspect_freq_map.items() if k i
335           filt_freq_map = sorted(filt_freq_map.items(), key = lambda x:
336           cluster_id_to_name_map[i] = filt_freq_map
337
338           # cluster_to_asp_map[i] = this_cluster_asp
339
340       # print(cluster_to_asp_map)
341       # print(cluster_id_to_name_map)
```

```python
342        return cluster_id_to_name_map
343
344
345
346  #Two master functions below for applying above functions
347
348
349  def extract_aspect_sentiment_tuples(df):
350      df = master_preprocess(df)
351      df = df.loc[df['wordcounts'] > 10].copy()
352      df.reset_index(drop=True, inplace=True)
353      df['aspect_tups'] = df.apply(apply_extraction, axis=1)
354      df = df.explode('aspect_tups').copy()
355      df.dropna(inplace=True)
356      df['asp'] = df['aspect_tups'].apply(lambda x: x[0])
357      df['modifier'] = df['aspect_tups'].apply(lambda x: x[1])
358      df['modifier_sentiment'] = df['aspect_tups'].apply(lambda x: x[2])
359      df['rule_number'] = df['aspect_tups'].apply(lambda x: x[3])
360      return df
361
362
363  #This function will work with input from user
364  #to find best possible label name for clusters
365  def get_cluster_name_inputs(df):
366      print('loading....')
367      aspect_freq_map = get_aspect_freq_map(df['asp'].values)
368      unique_asp_array = df['asp'].unique()
369      mapped_labels = get_word_cluster_labels(unique_asp_array)
370      asp_labels_map = dict(zip(unique_asp_array, mapped_labels))
371      label_names_map = get_cluster_names_map(asp_labels_map, aspect_fre
372
373      df['asp_cluster_label'] = df['asp'].map(asp_labels_map)
374
375      print("write misc if low counts and special characters")
376      print("the top word is usually the best fit")
377
378      display(label_names_map[0][:10])
379      print("Pick a category for above words: ")
380      clust_0 = input()
381
382      display(label_names_map[1][:10])
383      print("Pick a category for above words: ")
384      clust_1 = input()
385
386      display(label_names_map[2][:10])
387
388      print("Pick a category for above words: ")
389      clust_2 = input()
390
391      display(label_names_map[3][:10])
392      print("Pick a category for above words: ")
393      clust_3 = input()
394
395      clusters = [clust_0] + [clust_1] + [clust_2] + [clust_3]
396
397
398      name_clust_dict = {0: clusters[0],
```

```python
399                             1: clusters[1],
400                             2: clusters[2],
401                             3: clusters[3]
402                             }
403
404
405
406         df['cluster_name'] = df['asp_cluster_label'].map(name_clust_dict)
407
408         fig = plt.figure(figsize=(10,8))
409
410         df_senti.groupby(by='cluster_name')['modifier_sentiment'].sum().pl
411
412         plt.title("Clustered Aspect Sentiment Totals")
413
414         plt.savefig('images/extractor_example1.jpg')
415
416         return df
417
418
419
420
421     #function for visualizing variance in opinion for the same
422     #opinion/aspect couple among voters
423
424     def variance_visualizer(df, tup):
425         graph_s = df[df['tup_pair'] == tup]['sentiment'].value_counts()
426         fig = plt.figure(figsize=(6,4))
427
428
429         ax = graph_s.plot.bar()
430
431         #ADD TOTAL NUMBER OF VOTES
432         fig.text(0.9,
433                 0.9,
434                 s=f'Total Votes: {(graph_s.sum())} \n tuple: {tup}',
435                 ha='center',
436                 va='center',
437                 transform=ax.transAxes,
438                 size=12
439                 )
440
441         plt.savefig('images/variance_graph_example.jpg');
442
443
444
445
446     #function for visualizing variance in opinion for the same
447     #opinion/aspect couple among voters
448
449     def sum_squares(df, tup):
450         s = df[df['tup_pair'] == tup]['sentiment'].value_counts()
451     #    print(f'{tup} Residual sum of squares: {(np.sum((s-s.mean())**2)
452         sum_s = np.sum((s-s.mean())**2)
453         return(tup,sum_s)
```

# Running The Extractor

In [11]:
```
1  #aspect extractor ran in order to extract tuples from
2  #entire amazon review text body
3  df_senti = extract_aspect_sentiment_tuples(df)
4  df_senti.head()
```

Out[11]:

| | review_id | star_rating | review_body | wordcounts | aspect_tups | asp | modifier | modif |
|---|---|---|---|---|---|---|---|---|
| **0** | RA8R84N9JMZLD | 3 | My Dad loves this, only issue is charging, it ... | 32 | (connection, proper, 0.0, 1) | connection | proper | |
| **2** | R1KL6IIDB77A2O | 4 | I would not have paid the original price for t... | 58 | (price, original, 0.3182, 1) | price | original | |
| **2** | R1KL6IIDB77A2O | 4 | I would not have paid the original price for t... | 58 | (bit, little, 0.0, 1) | bit | little | |
| **2** | R1KL6IIDB77A2O | 4 | I would not have paid the original price for t... | 58 | (star, how durable, 0.0, 1) | star | how durable | |
| **2** | R1KL6IIDB77A2O | 4 | I would not have paid the original price for t... | 58 | (price, lower, -0.296, 1) | price | lower | |

In [12]:
```python
#keep only columns of interest
df_senti_predicted = df_senti.loc[:, ['modifier', 'asp', 'modifier_sent


#convert sentiment to Categorical to compare to human labels more easil
df_senti_predicted.loc[df_senti_predicted['modifier_sentiment'] < -.5,
df_senti_predicted.loc[(df_senti_predicted['modifier_sentiment'] >= -.5
df_senti_predicted.loc[df_senti_predicted['modifier_sentiment'] == 0, '
df_senti_predicted.loc[(df_senti_predicted['modifier_sentiment'] > 0) &
df_senti_predicted.loc[df_senti_predicted['modifier_sentiment'] > 0.5,

#convert aspects back to pair insted of triplet
#for easy comparison to human labels
df_senti_predicted['tup_pair'] = list(zip(df_senti_predicted['modifier'
#drop duplicate aspect pairs
df_to_turk2 = df_senti_predicted.drop_duplicates(subset='tup_pair')

df_to_turk2.sort_values(by='asp')
```

Out[12]:

|  | modifier | asp | modifier_sentiment | sentiment | tup_pair |
|---|---|---|---|---|---|
| 1072 | Able | | 0.0000 | Neutral | (Able, ) |
| 7660 | Though bulky | | 0.0000 | Neutral | (Though bulky, ) |
| 1900 | Only negative | | -0.5719 | Very Negative | (Only negative, ) |
| 6809 | Really disappointed | | -0.4767 | Negative | (Really disappointed, ) |
| 6810 | Otherwise excellent | | 0.5719 | Very Positive | (Otherwise excellent, ) |
| ... | ... | ... | ... | ... | ... |
| 5475 | fake | youout | -0.4767 | Negative | (fake, youout) |
| 3941 | past | yrs | 0.0000 | Neutral | (past, yrs) |
| 6229 | dead | zone | -0.6486 | Very Negative | (dead, zone) |
| 197 | sweet | zone | 0.4588 | Positive | (sweet, zone) |
| 7059 | extra | ~$20 | 0.0000 | Neutral | (extra, ~$20) |

17675 rows × 5 columns

In [13]:
```python
#sending data to create template
#for workers to submit labels
#using Amazon Turk
df_to_turk2 = df_senti_predicted.drop_duplicates(subset='tup_pair')
df_to_turk2.loc[:, ['modifier', 'asp']].to_csv('data/turk_data2.csv', i
```

Below is an example of the extractor being run to cluster the aspects and
return a bar graph of total sentiment (total summation of negative and
positive from the VADER sentiment intensity analyzer) as well as a
DataFrame with cluster names for further analyses. You can see that for
the product_id "B0001FTVEK" (https://www.amazon.com/Sennheiser-RS120-

Wireless-Headphones-Charging/dp/B0001FTVEK), which are RS120 Wireless
Headphones, there is a lot of positive sentiment for the value and
sound_quality categories as compared to the headphone_design and
hiss/tech_diff categories. The hiss category is low enough that it should
be the major focus for the company to funnel resources in response to
customer demand for improving their product. If they can focus first on
fixing the hiss mentioned in many amazon reviews, they can also put a few
extra resources into impoving some other design aspects of the headphones,
such as the batteries or cradle.

In [14]:
```
1  df_senti = get_cluster_name_inputs(df_senti)
```

```
loading....
write misc if low counts and special characters
the top word is usually the best fit

[('headphones', 2382),
 ('noise', 460),
 ('batteries', 432),
 ('headset', 425),
 ('headphone', 289),
 ('signal', 259),
 ('static', 228),
 ('unit', 213),
 ('cradle', 206),
 ('phones', 188)]

Pick a category for above words:
design

[('hiss', 227),
 (' ', 186),
```

In [15]:
```python
print('\nheadphone_design\n')
display(df_senti.loc[df_senti['cluster_name'] == 'design', 'asp'].value
print('\nheadphone_design\n')
df_senti.loc[df_senti['cluster_name'] == 'design', 'asp'].value_counts(
```

headphone_design

```
headphones     2382
noise           460
batteries       432
headset         425
headphone       289
signal          259
static          228
unit            213
cradle          206
phones          188
base            156
fit             150
audio           135
light           133
station         118
output          107
headsets         97
tuning           96
setup            84
device           82
frequency        81
earphones        80
Headphones       78
pads             78
transmitter      76
controls         75
channels         74
speakers         65
screen           65
battery          63
Name: asp, dtype: int64
```

headphone_design

## Loading Turk Data

 The DataFrames below are batch results loaded as csv files submitted from Amazon Mechanical Turk. As stated before, there are a total of 6107 non-null entries with information for 1438 unqiue aspects, submitted from a total of410 different workers from around the globe. Duplicate pairs of aspect/opinion pairs were included to inspect variance of submission from human labels and machine labels for each opinion pair.

In [16]:
```python
#Read all Turk Data

df1 = pd.read_csv("data/batch1_results.csv")
df2 = pd.read_csv("data/batch2_results.csv")
df3 = pd.read_csv("data/batch3_results.csv")
df4 = pd.read_csv("data/batch4_results.csv")
# df_lab = pd.read_csv("data/locally_sourced_labels.csv")

display(df1.info())
display(df2.info())
display(df3.info())
display(df4.info())
# df_lab.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1020 entries, 0 to 1019
Data columns (total 32 columns):
 #   Column                      Non-Null Count  Dtype
---  ------                      --------------  -----
 0   HITId                       1020 non-null   object
 1   HITTypeId                   1020 non-null   object
 2   Title                       1020 non-null   object
 3   Description                 1020 non-null   object
 4   Keywords                    1020 non-null   object
 5   CreationTime                1020 non-null   object
 6   MaxAssignments              1020 non-null   int64
 7   RequesterAnnotation         1020 non-null   object
 8   AssignmentDurationInSeconds 1020 non-null   int64
 9   AutoApprovalDelayInSeconds  1020 non-null   int64
 10  Expiration                  1020 non-null   object
 11  NumberOfSimilarHITs         0 non-null      float64
 12  LifetimeInSeconds           0 non-null      float64
 13  AssignmentId                1020 non-null   object
```

In [17]:
```python
1   #concact turk data into one dataframe
2   #keep only necessary columns
3
4   df_turk = pd.concat([df1,df2,df3, df4])
5
6   num_turk = len(df_turk.WorkerId.unique())
7   print(f'Total number of Turk Workers: {num_turk}\n')
8
9   df_turk = df_turk.loc[:, ['Input.Modifier', 'Input.Aspect',  'Answer.se
10  df_turk.dropna(inplace=True)
11
12  num_aspects = len(df_turk["Input.Aspect"].unique())
13  print(f'Total number of Unique Aspects: {num_aspects}\n')
14
15  display(df_turk.info())
16
17  #rename to differentiate human labels
18  df_turk.rename({'Input.Modifier': 'modifier',
19                 'Input.Aspect': 'asp',
20                 'Answer.sentiment.label': 'sentiment_h'
21                 }, axis=1, inplace=True)
22
23
24
25  df_turk.sort_values(by='asp')
```

```
Total number of Turk Workers: 410

Total number of Unique Aspects: 1438

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6107 entries, 0 to 3228
Data columns (total 3 columns):
 #   Column                  Non-Null Count   Dtype
---  ------                  --------------   -----
 0   Input.Modifier          6107 non-null    object
 1   Input.Aspect            6107 non-null    object
 2   Answer.sentiment.label  6107 non-null    object
dtypes: object(3)
memory usage: 190.8+ KB

None
```

Out[17]:

|      | modifier     | asp    | sentiment_h   |
|------|--------------|--------|---------------|
| 32   | extra        | $      | Very Positive |
| 115  | more missing | 1/2    | Negative      |
| 646  | single       | 1/8    | Neutral       |
| 1659 | stereo       | 1/8\\  | Neutral       |
| 589  | new          | 120s   | Positive      |
| ...  | ...          | ...    | ...           |
| 3042 | pair         | years  | Neutral       |

|      | modifier | asp | sentiment_h |
| --- | --- | --- | --- |
| **1031** | plus | years | Positive |
| **1878** | past | years | Positive |
| **2296** | fake | youout | Very Negative |
| **299** | sweet | zone | Positive |

6107 rows × 3 columns

In [18]:
```python
#create aspect/opinion pair to
#compare machine label vs human
df_turk['tup_pair'] = list(zip(df_turk['modifier'],df_turk['asp']))
display(df_turk.info())
df_senti_predicted.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 6107 entries, 0 to 3228
Data columns (total 4 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   modifier     6107 non-null   object
 1   asp          6107 non-null   object
 2   sentiment_h  6107 non-null   object
 3   tup_pair     6107 non-null   object
dtypes: object(4)
memory usage: 238.6+ KB

None

<class 'pandas.core.frame.DataFrame'>
Int64Index: 35355 entries, 0 to 7875
Data columns (total 5 columns):
 #   Column              Non-Null Count  Dtype
---  ------              --------------  -----
 0   modifier            35355 non-null  object
 1   asp                 35355 non-null  object
 2   modifier_sentiment  35355 non-null  float64
 3   sentiment           35355 non-null  object
 4   tup_pair            35355 non-null  object
dtypes: float64(1), object(4)
memory usage: 1.6+ MB
```

# Results

```python
In [19]:   1  #create merged DataFrame
           2  #in order to compare human
           3  #labeled aspect/opinion sentiment
           4  #vs model aspect/opinion sentiment
           5  #can check accuracy and precision
           6
           7  #how much do Turk workers agree with
           8  #my model?
           9
          10  #How much do Turk workers agree with
          11  #each other?
          12
          13  df_results = pd.merge(df_turk,
          14                        df_senti_predicted,
          15                        on='tup_pair',
          16                        how='left'
          17                        )
          18
          19  df_results.drop_duplicates(subset='tup_pair', inplace=True)
          20  df_results.head()
```

Out[19]:

| | modifier_x | asp_x | sentiment_h | tup_pair | modifier_y | asp_y | modifier_sentiment |
|---|---|---|---|---|---|---|---|
| **0** | old | father | Neutral | (old, father) | old | father | 0.0 |
| **16** | wireless | headphones | Neutral | (wireless, headphones) | wireless | headphones | 0.0 |
| **748** | charging | rack | Neutral | (charging, rack) | charging | rack | 0.0 |
| **752** | also available | headphones | Very Positive | (also available, headphones) | also available | headphones | 0.0 |
| **753** | impaired | husband | Negative | (impaired, husband) | impaired | husband | 0.0 |

```python
In [20]:   1  display(df_results.sentiment_h.value_counts())
           2  df_results.sentiment.value_counts()
```

```
Positive          1860
Neutral           1521
Very Positive      790
Negative           692
Very Negative      223
Name: sentiment_h, dtype: int64
```

Out[20]:
```
Neutral           3309
Positive           763
Very Positive      523
Negative           397
Very Negative       88
Name: sentiment, dtype: int64
```

In [21]:
```python
df_results2 = df_results.reset_index(drop=True)
df_results2.dropna(inplace=True)

matching_predictions = df_results2.loc[df_results2['sentiment_h'] == df
mismatching_predictions = df_results2.loc[df_results2['sentiment_h'] !=

display(matching_predictions.sentiment.value_counts())
mismatching_predictions.sentiment.value_counts()
```

```
Neutral          1198
Positive          394
Very Positive     187
Negative          141
Very Negative      21
Name: sentiment, dtype: int64
```
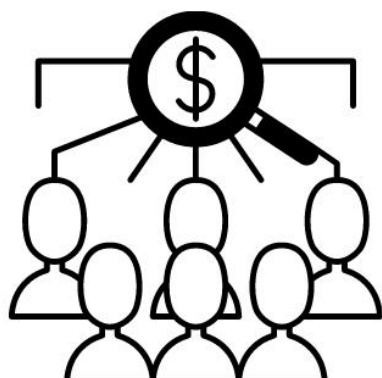
Out[21]:
```
Neutral          2111
Positive          369
Very Positive     336
Negative          256
Very Negative      67
Name: sentiment, dtype: int64
```

In [22]:
```python
#create table to see accuracy and precision of extractor
#AS COMPARED TO MY AMAZON TURK HUMAN LABELS
#ARE MY HUMAN LABELS RELIABLE????

agree = matching_predictions.sentiment.value_counts().values
disagree = mismatching_predictions.sentiment.value_counts().values
total_arr = agree + disagree
precision = agree/total_arr
denom = 5*[5080]
accuracy = agree/denom
columns = matching_predictions.sentiment.value_counts().index
df_table = pd.DataFrame(data=[precision, accuracy], columns=columns)
df_table = df_table.apply(lambda x: round(x,4))
df_table.rename({0: 'precision', 1: 'accuracy'},
                inplace = True
               )

df_table.to_csv("data/acc_prec_table.csv", index=False)

df_table.head()
```

Out[22]:

|           | Neutral | Positive | Very Positive | Negative | Very Negative |
|-----------|---------|----------|---------------|----------|---------------|
| precision | 0.3620  | 0.5164   | 0.3576        | 0.3552   | 0.2386        |
| accuracy  | 0.2358  | 0.0776   | 0.0368        | 0.0278   | 0.0041        |

# Comparison of Turk Data to Extractor Data

**VS**

Do we agree?

|  | Neutral | Positive | Very Positive | Negative | Very Negative |
|---|---|---|---|---|---|
| Precision | 36.20% | 51.64% | 35.76% | 35.52% | 23.86% |
| Accuracy | 23.58% | 7.76% | 3.68% | 2.78% | 0.41% |

The precision values explain how many times the machine label correctly line up with the human labels for each aspect/opinion pair when guessing for that appropriate sentiment. That is, if the machine were to guess only for negative sentiment, it was in agreement with humans 35% of the time. The accuracy scores represent how many times the machine were right overall for the entire dataset when predicted a certain class. Due to the nature of the dataset, accuracy scores will be low for a lot of the classes strictly because of class imbalances. Therefore, precision is a better metric for judging the performance of my model. However, there were some major issues with the overal expiremental setup that need to be discussed and anaylzed.

**Reliability of Experimental Setup**

To quickly visualize the variance among the different Amazon Turk workers assigned to labeling the aspect/opinion pairs, a function was utilized to create a bargraph that displays the different labels each worker voted for each aspect/opinion set that appeared in the dataset more than 10 times. A more statistical approach will be carried out using sum of squares residuals further down in the notebook. It is very apparent that the Amazon turks had a lot of trouble coming to any agreement on sentiment. The task was setup to reward workers to label data as quickly as possible without much safeguard to the quality of the work being submitted other than a quick overview by myself.I am just one poorly funded individual. Without proper funding from a research grant it was dificult to setup a reliable expirement. This cast a large shadow of doubt on the reliabilty of this data to be used as a way to reliably test the accuarcy of my model. In comparison, you can see that the extractor chooses the same sentiment for a unique aspect/opinion pair every single time it shows

up in the dataset. While the human data has been revealed to be severely flawed, it has brought up a shining example as to why machine learning may be a better substitute for labeling large amounts of tedius data in the first place.

```
In [23]:    1   #create dataframe for tuple_pairs with more than 10 votes from turk dat
            2
            3   multi_votes = df_turk['tup_pair'].value_counts()[df_turk['tup_pair'].va
            4
            5   df_turk2 = df_turk[df_turk["tup_pair"].isin(multi_votes)].copy()
            6
            7   df_turk2.rename({'sentiment_h': 'sentiment'},
            8                    axis=1,
            9                    inplace=True
           10                    )
           11
           12   unique_tup = df_turk2.tup_pair.unique()
           13
           14
           15
           16   #create dataframe for tuple_pairs with more than 10 votes from extracto
           17
           18   multi_votes_e = df_senti_predicted['tup_pair'].value_counts()[df_senti_
           19
           20   df_ex = df_senti_predicted[df_senti_predicted["tup_pair"].isin(multi_vc
           21
           22   unique_tup_e = df_ex.tup_pair.unique()
```
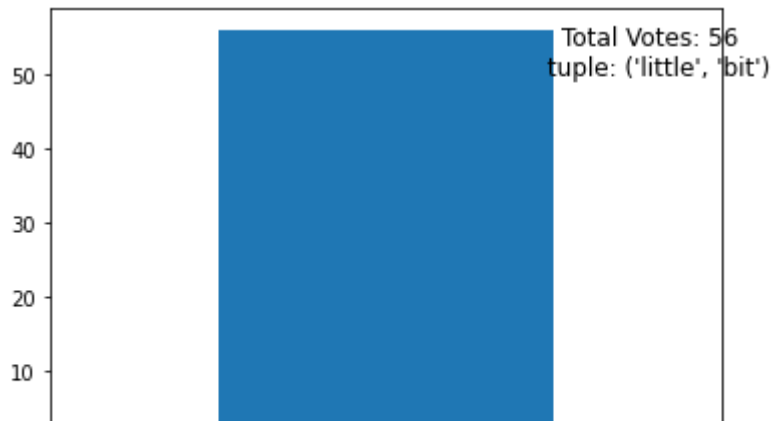
# Partition Sum Of Squares to Measure Disperson of Turk Data

Partition Sum of Squares (https://en.wikipedia.org/wiki/Partition_of_sums_of_squares)

For simplicity, sum of square statistical measurements were calculated for each aspect in the human labeled data as well as the machine labeled data that had more than 10 votes. This was chosen over entropy as a calculation for dispersion as the "sum of squares" for these data is a close enough estimate for the variance in categorical data for this expirement. Every single human labeled aspect with multiple aspects had a non-zero value for residual sum of squares, while every machine labeled aspect had zero residual error. This shows the difference in variance statistically and mathmatically very clearly between the two and shows the unreliability of the turk data. See the figure below for the total range in the residual sum of square values for the tested aspects in each group.

```
In [24]:    1  for tup in unique_tup_e:
            2      variance_visualizer(df_ex, tup)
```

```
/Users/dylandey/anaconda3/envs/learn-env/lib/python3.6/site-packages/
ipykernel_launcher.py:426: RuntimeWarning: More than 20 figures have
been opened. Figures created through the pyplot interface (`matplotli
b.pyplot.figure`) are retained until explicitly closed and may consum
e too much memory. (To control this warning, see the rcParam `figure.
max_open_warning`).
```



```
In [25]:    1  for tup in unique_tup:
            2      variance_visualizer(df_turk2, tup)
```



```
In [26]:    1  for tup in unique_tup_e:
            2      sum_squares(df_ex, tup)
```

In [27]:
```python
emp_list = []
for tup in unique_tup:
    emp_list.append(sum_squares(df_turk2, tup))


df_turk_var = pd.DataFrame(emp_list, columns=['asp', 'sum_of_squares'])
df_turk_var.to_csv('data/turk_var1.csv')

emp_list = []
for tup in unique_tup_e:
    sum_squares(df_ex, tup)

df_ex_var = pd.DataFrame(emp_list, columns=['asp', 'sum_of_squares'])
df_ex_var.to_csv('data/human_var1.csv')
```

In [28]:
```python
df_turk_var
```

Out[28]:

|    | asp | sum_of_squares |
|----|-----|----------------|
| 0 | (wireless, headphones) | 232.750000 |
| 1 | (second, set) | 4.750000 |
| 2 | (clear, sound) | 26.000000 |
| 3 | (sound, quality) | 870.000000 |
| 4 | (great, headphones) | 19.000000 |
| 5 | (second, pair) | 11.200000 |
| 6 | (good, sound) | 60.750000 |
| 7 | (rechargeable, batteries) | 4.666667 |
| 8 | (good, range) | 16.666667 |
| 9 | (great, quality) | 58.750000 |
| 10 | (Sound, quality) | 4.500000 |
| 11 | (great, product) | 34.750000 |
| 12 | (great, range) | 18.666667 |
| 13 | (great, sound) | 60.666667 |
| 14 | (good, quality) | 83.000000 |
| 15 | (long, time) | 11.200000 |

In [29]:
```python
# all values sum to zero

df_ex_var
```

Out[29]:

| asp | sum_of_squares |
|-----|----------------|

### Further Discussion and Future Work

1) Due to a lack of funding from grants and some other issues in expiremental design, I feel the most appropriate next step would be to repeat the experiment with more carefully collected labeled data. Increasing the reward per label and the requirements for submission (such as proving a proficiency in english) I believe cam substantiely improve the quality of the human labeled data and reduce the variance in this data significantly. Sourcing reliable test data is the number one priority in continuing future work for this project.

2) Integrate into AWS for scaling

3) integrate into a SQL server for data management

4) incoroprate a flash based UI for viewing results at scale

# THANK YOU

[Blog (https://dev.to/ddey117)](https://dev.to/ddey117) | [GitHub (https://github.com/ddey117/ABSA_Project_4)](https://github.com/ddey117/ABSA_Project_4) | [PreProcess Github (https://github.com/ddey117/preprocess_ddey117)](https://github.com/ddey117/preprocess_ddey117)

Dylan Dey

email: [ddey2985@gmail.com (mailto:ddey2985@gmail.com)](mailto:ddey2985@gmail.com)



# Appendix

## Text Data Preprocessing and Cleaning toolkit

### Word Counts

```
In [30]:    1  ls data
```

```
acc_prec_table.csv                human_var1.csv
batch1_results.csv                locally_sourced_labels.csv
batch2_results.csv                modifier_aspect_unlabeled.csv
batch3_results.csv                turk_data2
batch4_results.csv                turk_data2.csv
df_electronics_example.csv        turk_var1.csv
```

```
In [31]:    1  df = df.sample(100)
```

```
In [32]:    1  len('this is text'.split())
```

Out[32]: 3

```
In [33]:    1  df['word_counts'] = df['review_body'].apply(lambda x: len(str(x).split(
```

```
In [34]:    1  df.sample(5)
```

Out[34]:

|       | review_id        | star_rating | review_body                                      | wordcounts | word_counts |
|-------|------------------|-------------|--------------------------------------------------|------------|-------------|
| **1090** | R2EV6UHK8ZZKPP | 5           | 92 year old dad LOVED these! Said it was the ... | 21         | 21          |
| **3360** | R2F8Z2G2RSUZ8  | 3           | I thought when I ordered this would be the sam... | 69         | 69          |
| **5394** | R2BG44IAP6EQH9 | 5           | It opened a whole new world of voice and music... | 23         | 23          |
| **378**  | R1YMTKNTUYLCJZ | 5           | My RS 120s are my third set of headphones. My... | 106        | 106         |
| **123**  | R1BUZY42O4J4D6 | 5           | Great Product! Keeps your partner happy when w... | 18         | 18          |

```
In [35]:    1  df['word_counts'].max()
```

Out[35]: 474

```
In [36]:    1  df['word_counts'].min()
```

Out[36]: 1

In [37]:
```python
1  df[df['word_counts']==1]
```

Out[37]:

| | review_id | star_rating | review_body | wordcounts | word_counts |
|---|---|---|---|---|---|
| **2073** | R1T3D9ZJKPQF9V | 4 | good | 1 | 1 |

## Characters Count

In [38]:
```python
1  len('this is')
```

Out[38]: 7

In [39]:
```python
1  def char_counts(x):
2      s = x.split()
3      x = ''.join(s)
4      return len(x)
```

In [40]:
```python
1  char_counts('this is')
```

Out[40]: 6

In [42]:
```python
1  df['char_counts'] = df['review_body'].apply(lambda x: char_counts(str(x
```

In [43]:
```python
1  df.sample(5)
```

Out[43]:

| | review_id | star_rating | review_body | wordcounts | word_counts | char_counts |
|---|---|---|---|---|---|---|
| **5702** | R11YXEXQVPM3LJ | 5 | In the living room with the TV in the middle o... | 86 | 86 | 345 |
| **5850** | R2J9ADIZPDEMFR | 1 | This item has never worked. All I get is stati... | 34 | 34 | 147 |
| **2345** | RE17VN0I9T0X1 | 4 | Earphones works great. I'd give it 5 stars, ex... | 38 | 38 | 171 |
| **1578** | R3O1FPZB3J46D6 | 5 | I have had other head sets. they aren't even c... | 36 | 36 | 152 |
| **8506** | R36THE4Z3X3797 | 5 | I searched everywhere to find the best headpho... | 209 | 209 | 938 |

## Average Word Length

In [44]:
```python
1 x = 'this is' # 6/2 = 3
2 y = 'thankyou guys' # 12/2 = 6
```

In [45]:
```python
1 df['avg_word_len'] = df['char_counts']/df['word_counts']
```

In [46]:
```python
1 df.sample(4)
```

Out[46]:

| | review_id | star_rating | review_body | wordcounts | word_counts | char_counts | avg_wo |
|---|---|---|---|---|---|---|---|
| **800** | R3KVSFIW1PCEZ4 | 5 | works great after initially adjusting tuning f... | 7 | 7 | 48 | 6.8 |
| **5236** | R7PV556J6UHKF | 1 | Was really hoping these would work well but th... | 68 | 68 | 299 | 4.3 |
| **7581** | R3QZYC0TDVCQYE | 4 | The overall iem is good. The fit is a little l... | 17 | 17 | 71 | 4.1 |
| **3683** | R1EHM1VMSNBDKI | 4 | THE ONLY CRITICISM I HAVE IS THAT THE EAR PHON... | 26 | 26 | 118 | 4.5 |

## Stop Words Count

In [47]:
```python
1 print(stopwords)
```

```
<WordListCorpusReader in '.../corpora/stopwords' (not loaded yet)>
```

In [48]:
```python
1 len(stopwords)
```

```
---------------------------------------------------------------------------
TypeError                                 Traceback (most recent call last)
<ipython-input-48-0e227a118e01> in <module>
----> 1 len(stopwords)

TypeError: object of type 'LazyCorpusLoader' has no len()
```

```
In [49]:    1  x = 'this is the text data'
```

```
In [50]:    1  x.split()
```

```
Out[50]: ['this', 'is', 'the', 'text', 'data']
```

```
In [51]:    1  [t for t in x.split() if t in stopwords]
```

```
---------------------------------------------------------------------------
TypeError                                 Traceback (most recent call last)
<ipython-input-51-d6f9f3c59f41> in <module>
----> 1 [t for t in x.split() if t in stopwords]

<ipython-input-51-d6f9f3c59f41> in <listcomp>(.0)
----> 1 [t for t in x.split() if t in stopwords]

TypeError: argument of type 'LazyCorpusLoader' is not iterable
```

```
In [ ]:    1  len([t for t in x.split() if t in stopwords])
```

```
In [ ]:    1  df['stop_words_len'] = df['tweets'].apply(lambda x: len([t for t in x.s
```

```
In [ ]:    1  df.sample(5)
```

## Count #HashTags and @Mentions

```
In [52]:    1  x = 'this is #hashtag and this is @mention'
```

```
In [53]:    1  x.split()
```

```
Out[53]: ['this', 'is', '#hashtag', 'and', 'this', 'is', '@mention']
```

```
In [54]:    1  [t for t in x.split() if t.startswith('@')]
```

```
Out[54]: ['@mention']
```

```
In [55]:    1  len([t for t in x.split() if t.startswith('@')])
```

```
Out[55]: 1
```

```
In [58]:  1  df['hashtags_count'] = df['review_body'].apply(lambda x: len([t for t i
```

```
In [59]:  1  df['mentions_count'] = df['review_body'].apply(lambda x: len([t for t i
```

```
In [60]:  1  df.sample(5)
```

Out[60]:

| | review_id | star_rating | review_body | wordcounts | word_counts | char_counts | avg_wor |
|---|---|---|---|---|---|---|---|
| **6936** | R1BJ04XPMLMDXU | 3 | I really liked almost everything about this un... | 248 | 248 | 997 | 4.0: |
| **1052** | RTV3GR33VC179 | 4 | I like these headphones, I get good range and ... | 96 | 96 | 409 | 4.2( |
| **5926** | R93XGUJB1286 | 1 | When I first received this product, the chargi... | 27 | 27 | 132 | 4.8i |
| **5702** | R11YXEXQVPM3LJ | 5 | In the living room with the TV in the middle o... | 86 | 86 | 345 | 4.0 |
| **7985** | R1SC5Z6J9E9OYG | 1 | I got this product on April 1 and when I tried... | 128 | 128 | 522 | 4.0: |

## If numeric digits are present in tweets

```
In [61]:  1  x = 'this is 1 and 2'
```

```
In [62]:  1  x.split()
```

Out[62]:  ['this', 'is', '1', 'and', '2']

```
In [63]:  1  x.split()[3].isdigit()
```

Out[63]:  False

In [64]:
```python
1  [t for t in x.split() if t.isdigit()]
```

Out[64]: ['1', '2']

In [65]:
```python
1  df['numerics_count'] = df['review_body'].apply(lambda x: len([t for t i
```

In [66]:
```python
1  df.sample(5)
```

Out[66]:

| | review_id | star_rating | review_body | wordcounts | word_counts | char_counts | avg_wo |
|---|---|---|---|---|---|---|---|
| 4935 | R2CZX0A7L8IWG2 | 5 | Best headphones I ever owned. The sound qualit... | 40 | 40 | 181 | 4.5 |
| 504 | R2JGAAP0GCZUKN | 3 | Sound is pretty good. Charging base is nice. ... | 100 | 100 | 410 | 4.1 |
| 976 | RBXAL93CXFM2V | 5 | Exactly what I needed for a fair price with ti... | 11 | 11 | 50 | 4.5 |
| 5850 | R2J9ADIZPDEMFR | 1 | This item has never worked. All I get is stati... | 34 | 34 | 147 | 4.3 |
| 48 | R3UH6JAA5NRDWK | 5 | These are great, no more batteries and we can ... | 16 | 16 | 66 | 4.1 |

## UPPER case words count

In [67]:
```python
1  x = 'I AM HAPPY'
2  y = 'i am happy'
```

In [68]:
```python
1  [t for t in x.split() if t.isupper()]
```

Out[68]: ['I', 'AM', 'HAPPY']

In [69]:
```python
1  df['upper_counts'] = df['review_body'].apply(lambda x: len([t for t in
```

```
In [70]:    1  df.sample(5)
```

Out[70]:

|  | review_id | star_rating | review_body | wordcounts | word_counts | char_counts |
|---|---|---|---|---|---|---|
| **4423** | R18HP3Z44HMBC4 | 4 | It is good but, sometimes have to adjust to ge... | 27 | 27 | 117 |
| **4381** | R10VX4AU06SV1B | 4 | I use these headphones all the time. The sound... | 21 | 21 | 106 |
| **976** | RBXAL93CXFM2V | 5 | Exactly what I needed for a fair price with ti... | 11 | 11 | 50 |
| **2946** | R4HM1MIPN5BOI | 5 | Wish I knew about this headset before going th... | 58 | 58 | 251 |
| **8300** | R3K0Z5RQKPVXT5 | 5 | I bought the wireless headphones(sennheiser 12... | 39 | 39 | 182 |

## Lower Case Conversion

```
In [ ]:    1  x = 'this is Text'
```

```
In [ ]:    1  x.lower()
```

```
In [ ]:    1  x = 45.0
           2  str(x).lower()
```

```
In [72]:   1  df['review_body'] = df['review_body'].apply(lambda x: str(x).lower())
```

In [73]:
```
1  df.sample(5)
```

Out[73]:

| | review_id | star_rating | review_body | wordcounts | word_counts | char_counts | avg_wor |
|---|---|---|---|---|---|---|---|
| **5595** | R2G6GVNS929E0X | 2 | the little dial used to "tune" this headset in... | 64 | 64 | 278 | 4.34 |
| **4423** | R18HP3Z44HMBC4 | 4 | it is good but, sometimes have to adjust to ge... | 27 | 27 | 117 | 4.33 |
| **6027** | R1YZZ822BSFSVF | 1 | the product is excelent, but this refurbished ... | 22 | 22 | 100 | 4.54 |
| **4935** | R2CZX0A7L8IWG2 | 5 | best headphones i ever owned. the sound qualit... | 40 | 40 | 181 | 4.52 |
| **6961** | R26HOZU0XT9M06 | 5 | after reviewing many different headphones in t... | 75 | 75 | 333 | 4.44 |

## Contraction to Expansion

```python
#created using wikipedia

contractions = {
"a'ight": "alright",
"ain't": "am not",
"amn't": "am not",
"arencha": "are not you",
"aren't": "are not",
"'bout": "about",
"can't": "cannot",
"cap'n": "captain",
"'cause": "because",
"'cept" : "except",
"could've": "could have",
"couldn't": "could not",
"couldn't've": "could not have",
"cuppa": "cup of",
"dammit": "damn it",
"daren't": "dared not",
"daresn't": "dare not",
"dasn't": "dare not",
"didn't": "did not",
"doesn't": "does not",
"don't": "do not",
"dunno" : "do not know",
"d'ye": "do you",
"e'en": "even",
"e'er": "ever",
"'em": "them",
"everybody's": "everybody is",
"everyone's": "everyone is",
"finna": "fixing to",
"fo'c'sle": "forecastle",
"'gainst": "against",
"g'day": "good day",
"gimme": "give me",
"giv'n": "given",
"gi'": "give us",
"gonna": "going to",
"gon't": "go not" ,
"gotta": "got to",
"hadn't": "had not",
"had've": "had have",
"hasn't": "has not",
"haven't": "have not",
"he'd":  "he would",
"he'll": "he will",
"helluva": "hell of a",
"he's": "he is",
"here's": "here is",
"how'd" : "how would",
"howdy" : "how do you fare",
"how'll": "how will",
"how're": "how are",
"how's": "how is",
"i'd": "i would",
```

```
 57    "i'd've": "i would have",
 58    "i'd'nt": "i would not",
 59    "i'd'nt've": "i would not have",
 60    "i'll": "i will",
 61    "i'm": "i am",
 62    "imma": "i am about to",
 63    "i'm'o": "i am going to",
 64    "innit": "is it not it",
 65    "i've": "i have",
 66    "isn't": "is not",
 67    "it'd": "it would",
 68    "it'll": "it will",
 69    "it's": "it is",
 70    "idunno": "i do not know",
 71    "kinda" : "kind of",
 72    "let's": "let us",
 73    "loven't": "love not",
 74    "ma'am": "madam",
 75    "mayn't": "may not",
 76    "may've": "may have",
 77    "methinks":" i think",
 78    "mightn't": "might not",
 79    "might've": "might have",
 80    "mustn't": "must not",
 81    "mustn't've": "must not have",
 82    "must've": "must have",
 83    "'neath": "beneath",
 84    "needn't": "need not",
 85    " nal ": " and all ",
 86    "ne'er": "never",
 87    " o'er" : " over ",
 88    " ol '": " old ",
 89    "oughtn'": "ought not",
 90    "'round": "around",
 91    "shalln't": "shall not",
 92    "shan't": "shall not",
 93    "she'd": "she would",
 94    "she'll": "he will",
 95    "she's": "she is",
 96    "should've": "should have",
 97    "shouldn't": "should not",
 98    "shouldn't've": "should not have",
 99    "somebody's": "somebody is",
100    "someone's": "someone is",
101    "something's": "something is",
102    "so're": "so are",
103    "so's": "so is",
104    "so've": "so have",
105    "that'll": "that shall",
106    "that're": "that are",
107    "that's": "that is",
108    "that'd": "that would",
109    "there'd": "there had",
110    "there'll": "there shall",
111    "there're": "there are",
112    "there's": "there is",
113    "these're": "these are",
```

```
114    "these've": "these have",
115    "they'd": "they would",
116    "they'll": "they will",
117    "they're": "they are",
118    "they've": "they have",
119    "this's": "this is",
120    "those're": "those are",
121    "those've": "those have",
122    "'thout": "without",
123    "'til": "until",
124    "'tis": "it is",
125    "to've": "to have",
126    "'twas": "it was",
127    "'tween": "between",
128    "'twere": "it were",
129    "w'all": "we all",
130    "wanna": "want to",
131    "wasn't": "was not",
132    " we'd ": " we would ",
133    "we'd've": "we would have",
134    "we'll": "we will",
135    "we're": "we are",
136    "we've": "we have",
137    "weren't": "were not",
138    "whatcha": "what are you",
139    "what'd": "what did",
140    "what'll": "what will",
141    "what're": "what are",
142    "what's": "what is",
143    "what've": "what have",
144    "when's": "when is",
145    "where'd": "where did",
146    "where'll": "where will",
147    "where're": "where are",
148    "where's": "where is",
149    "where've": "where have",
150    "which'd": "which would",
151    "which'll": "which will",
152    "which're": "which are",
153    "which's": "which is",
154    "which've": "which have",
155    "who'd": "who would",
156    "who'd've": "who would have",
157    "who'll": "who will",
158    "who're":" who are",
159    "who's": "who is ",
160    "who've": "who have",
161    "why'd": "why did",
162    "why're": "why are",
163    "why's": "why has ",
164    "willn't": "will not",
165    "won't": "will not",
166    "wonnot": "will not",
167    "would've": "would have",
168    "wouldn't": "would not",
169    "wouldn't've": "would not have",
170    "y'all": "you all ",
```

```
171   "y'all'd've": "you all would have",
172   "y'all'd'n't've": "you all would not have",
173   "y'all're": "you all are" ,
174   "y'all'ren't": "you all are not",
175   " y'at ": " you at ",
176   "yes'm": "yes madam",
177   "y'know": "you know",
178   " yessir ": " yes sir ",
179   "you'd": "you would",
180   "you'll": "you will",
181   "you're": "you are",
182   "you've": "you have",
183   "when'd": "when did",
184   "willn't": "will not",
185   #append extra web lingo
186   " u ": " you ",
187   " ur ": " your ",
188   " n ": " and ",
189   " dis ": " this ",
190   " bak ": " back ",
191   " brng ": " bring ",
192   "sux": "sucks",
193   "gr8": "great",
194   " amzing ": " amazing ",
195   " k ": "ok",
196   " kk ": "ok",
197   " il ": "i will"}
```

In [78]:
```
1   x = "i'm don't he'll" # "i am do not he will"
```

In [79]:
```
1   def cont_to_exp(x):
2       if type(x) is str:
3           for key in contractions:
4               value = contractions[key]
5               x = x.replace(key, value)
6           return x
7       else:
8           return x
9
```

In [80]:
```
1   cont_to_exp(x)
```

Out[80]:  'i am do not he will'

In [81]:
```
1   %%timeit
2   df['review_body'] = df['review_body'].apply(lambda x: cont_to_exp(x))
```

8.42 ms ± 94.7 µs per loop (mean ± std. dev. of 7 runs, 100 loops each)

## Count and Remove Emails

## Count and Remove Emails

```python
In [82]:   1  import re
```

```
In [83]:   1  df[df['review_body'].str.contains('hotmail.com')]['treview_body'].apply
```

```
---------------------------------------------------------------
--
KeyError                                      Traceback (most recent call las
t)
~/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/indexe
s/base.py in get_loc(self, key, method, tolerance)
   2645            try:
-> 2646                return self._engine.get_loc(key)
   2647            except KeyError:

pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.get_loc()

pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.get_loc()

pandas/_libs/hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjec
tHashTable.get_item()

pandas/_libs/hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjec
tHashTable.get_item()

KeyError: 'treview_body'

During handling of the above exception, another exception occurred:

KeyError                                      Traceback (most recent call las
t)
<ipython-input-83-ef9c4442bc0c> in <module>
----> 1 df[df['review_body'].str.contains('hotmail.com')]['treview_body']
.apply(print)

~/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/frame.
py in __getitem__(self, key)
   2798            if self.columns.nlevels > 1:
   2799                return self._getitem_multilevel(key)
-> 2800            indexer = self.columns.get_loc(key)
   2801            if is_integer(indexer):
   2802                indexer = [indexer]

~/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/indexe
s/base.py in get_loc(self, key, method, tolerance)
   2646                return self._engine.get_loc(key)
   2647            except KeyError:
-> 2648                return self._engine.get_loc(self._maybe_cast_inde
xer(key))
   2649        indexer = self.get_indexer([key], method=method, toleranc
e=tolerance)
   2650        if indexer.ndim > 1 or indexer.size > 1:

pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.get_loc()

pandas/_libs/index.pyx in pandas._libs.index.IndexEngine.get_loc()

pandas/_libs/hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjec
tHashTable.get_item()
```

```
pandas/_libs/hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjec
tHashTable.get_item()
```

`KeyError: 'treview_body'`

In [ ]: `1  x = '@securerecs arghh me please  markbradbury_16@hotmail.com'`

In [ ]: `1  re.findall(r'([a-z0-9+._-]+@[a-z0-9+._-]+\.[a-z0-9+_-]+)', x)`

In [ ]: `1  df['emails'] = df['review_body'].apply(lambda x: re.findall(r'([a-z0-9+`

In [ ]: `1  df['emails_count'] = df['emails'].apply(lambda x: len(x))`

In [ ]: `1  df[df['emails_count']>0]`

In [ ]: `1  re.sub(r'([a-z0-9+._-]+@[a-z0-9+._-]+\.[a-z0-9+_-]+)',"", x)`

In [ ]: `1  df['review_body'] = df['review_body'].apply(lambda x: re.sub(r'([a-z0-9`

In [ ]: `1  df[df['emails_count']>0]`

## Count URLs and Remove it

In [ ]: `1  x = 'hi, thanks to watching it. for more visit https://youtube.com/notr`

In [ ]: `1  #shh://git@git.com:username/repo.git=riif?%`

In [ ]: `1  re.findall(r'(http|https|ftp|ssh)://([\w_-]+(?:(?:\.[\w_-]+)+))([\w.,@?`

In [ ]: `1  df['url_flags'] = df['review_body'].apply(lambda x: len(re.findall(r'(h`

In [ ]: `1  df[df['url_flags']>0].sample(5)`

```
In [ ]:    1  x
```

```
In [ ]:    1  re.sub(r'(http|https|ftp|ssh)://([\w_-]+(?:(?:\.[\w_-]+)+))([\w.,@?^=%&
```

```
In [ ]:    1  df['review_body'] = df['review_body'].apply(lambda x: re.sub(r'(http|ht
```

```
In [ ]:    1  df.sample(5)
```

## Remove RT

```
In [ ]:    1  df[df['review_body'].str.contains('rt ')]
```

```
In [ ]:    1  x = "beautifully smart and simple idea rt @madebyma"
```

```
In [ ]:    1  re.sub(r'rt\ ', '', x).strip()
```

```
In [ ]:    1  df['review_body'] = df['review_body'].apply(lambda x: re.sub(r'rt\ ', '
```

## Special Chars removal or punctuation removal

```
In [ ]:    1  df.sample(10)
```

```
In [ ]:    1  x = '@mention (cnnmoney) for #sxsw 2011, any comput...'
```

```
In [ ]:    1  re.sub(r'[^\w ]+', "", x)
```

```
In [ ]:    1  df['review_body'] = df['review_body'].apply(lambda x: re.sub(r'[^\w ]+'
```

```
In [ ]:    1  df.sample(5)
```

## Remove multiple spaces `"hi    hello      "`

```
In [84]:  1  x = 'i     like      my          space     but       how are you'
```

```
In [85]:  1  ' '.join(x.split())
```

Out[85]: 'i like my space but how are you'

```
In [86]:  1  df['review_body'] = df['review_body'].apply(lambda x: ' '.join(x.split(
```

## Remove HTML tags

```
In [87]:  1  from bs4 import BeautifulSoup
```

```
In [88]:  1  x = '<html><h1> i love html so much </h1></html>'
```

```
In [89]:  1  # x.replace('<html><h1>', '').replace('</h1></html>', '') #don't
```

```
In [90]:  1  BeautifulSoup(x, 'lxml').get_text().strip()
```

Out[90]: 'i love html so much'

```
In [91]:  1  %%time
          2  df['review_body'] = df['review_body'].apply(lambda x: BeautifulSoup(x,
```

```
CPU times: user 23.8 ms, sys: 2.4 ms, total: 26.2 ms
Wall time: 25.2 ms
```

## Remove Accented Chars

```
In [92]:  1  x = 'Áccěntěd těxt'
```

```
In [93]:  1  import unicodedata
```

```
In [94]:  1  def remove_accented_chars(x):
          2      x = unicodedata.normalize('NFKD', x).encode('ascii', 'ignore').deco
          3      return x
```

```
In [95]:    1  remove_accented_chars(x)
```

```
Out[95]:  'Accented text'
```

```
In [96]:    1  df['review_body'] = df['review_body'].apply(lambda x: remove_accented_c
```

## Remove Stop Words

```
In [97]:    1  df['review_no_stop'] = df['review_body'].apply(lambda x: ' '.join([t fc
```

```
---------------------------------------------------------------------
--
TypeError                                 Traceback (most recent call las
t)
<ipython-input-97-ee168402147b> in <module>
----> 1 df['tweets_no_stop'] = df['review_body'].apply(lambda x: ' '.join
([t for t in x.split() if t not in stopwords]))

~/anaconda3/envs/learn-env/lib/python3.6/site-packages/pandas/core/serie
s.py in apply(self, func, convert_dtype, args, **kwds)
   3846                else:
   3847                    values = self.astype(object).values
-> 3848                    mapped = lib.map_infer(values, f, convert=convert
_dtype)
   3849
   3850            if len(mapped) and isinstance(mapped[0], Series):

pandas/_libs/lib.pyx in pandas._libs.lib.map_infer()

<ipython-input-97-ee168402147b> in <lambda>(x)
----> 1 df['tweets_no_stop'] = df['review_body'].apply(lambda x: ' '.join
([t for t in x.split() if t not in stopwords]))

<ipython-input-97-ee168402147b> in <listcomp>(.0)
----> 1 df['tweets_no_stop'] = df['review_body'].apply(lambda x: ' '.join
([t for t in x.split() if t not in stopwords]))

TypeError: argument of type 'LazyCorpusLoader' is not iterable
```

## Convert into base or root form of word

In [98]:
```python
1  nlp = spacy.load('en_core_web_sm')
```

```
---------------------------------------------------------------------------
OSError                                   Traceback (most recent call las
t)
<ipython-input-98-12102ff9e1a9> in <module>
----> 1 nlp = spacy.load('en_core_web_sm')

~/anaconda3/envs/learn-env/lib/python3.6/site-packages/spacy/__init__.py
 in load(name, vocab, disable, exclude, config)
     50      """
     51      return util.load_model(
---> 52          name, vocab=vocab, disable=disable, exclude=exclude, conf
ig=config
     53      )
     54

~/anaconda3/envs/learn-env/lib/python3.6/site-packages/spacy/util.py in l
oad_model(name, vocab, disable, exclude, config)
    425      if name in OLD_MODEL_SHORTCUTS:
    426          raise IOError(Errors.E941.format(name=name, full=OLD_MODE
L_SHORTCUTS[name]))  # type: ignore[index]
--> 427      raise IOError(Errors.E050.format(name=name))
    428
    429

OSError: [E050] Can't find model 'en_core_web_sm'. It doesn't seem to be
 a Python package or a valid path to a data directory.
```

In [99]:
```python
1  x = 'i having a craved chocolates. what is times? these dogs likes ball
```

In [100]:
```python
 1  def make_to_base(x):
 2      x = str(x)
 3      x_list = []
 4      doc = nlp(x)
 5
 6      for token in doc:
 7          lemma = token.lemma_
 8          if lemma == '-PRON-' or lemma == 'be':
 9              lemma = token.text
10
11          x_list.append(lemma)
12      return ' '.join(x_list)
```

In [101]:
```python
1  make_to_base(x)
```

Out[101]: 'I have a crave chocolate . what is time ? these dog like ball'

In [102]:
```
1  df['review_body'] = df['review_body'].apply(lambda x: make_to_base(x))
```

In [103]:
```
1  df.sample(5)
```

Out[103]:

|  | review_id | star_rating | review_body | wordcounts | word_counts | char_counts | avg_wor |
|---|---|---|---|---|---|---|---|
| **2345** | RE17VN0I9T0X1 | 4 | earphone work great . I would give it 5 star ,... | 38 | 38 | 171 | 4.5( |
| **1983** | R11U8YGNFAH57P | 5 | these head phone are wonderful ! my husband ca... | 34 | 34 | 120 | 3.5; |
| **927** | R2EQGEKSVGSSXF | 5 | this is the 3rd set I have own . 1st one the b... | 87 | 87 | 345 | 3.9( |
| **800** | R3KVSFIW1PCEZ4 | 5 | work great after initially adjust tuning frequ... | 7 | 7 | 48 | 6.8; |
| **4449** | RGEUGUCHQ9BU4 | 4 | fit and sound well . I can move around the hou... | 22 | 22 | 89 | 4.0( |

## Common words removal

In [104]:
```
1  text = ' '.join(df['review_body'])
```

In [105]:
```
1  len(text)
```

Out[105]: 40530

In [106]:
```
1  text = text.split()
```

In [107]:
```
1  len(text)
```

Out[107]: 8514

In [108]:
```
1  freq_comm = pd.Series(text).value_counts()
```

```
In [109]:   1  f20 = freq_comm[:20]
```

```
In [110]:   1  f20
```

```
Out[110]:  .              451
           the            442
           I              257
           to             239
           ,              232
           and            214
           is             149
           a              143
           it             125
           they           112
           have           105
           of              98
           not             95
           you             93
           my              90
           for             75
           on              73
           in              72
           headphone       72
           are             66
           dtype: int64
```

```
In [111]:   1  df['review_body'] = df['review_body'].apply(lambda x: ' '.join([t for t
```

In [112]:     1  df.sample(5)

Out[112]:

| | review_id | star_rating | review_body | wordcounts | word_counts | char_counts | avg_wor |
|---|---|---|---|---|---|---|---|
| **4761** | R16DZYDQG883HT | 5 | significant amount time elderly father been in... | 355 | 355 | 1641 | 4.62 |
| **6917** | RGHS0APBB60AY | 4 | like many wireless headset sound quality come ... | 259 | 259 | 1180 | 4.55 |
| **3026** | R1YI3CU3PJN55R | 4 | great sound price but do tend fall off when be... | 26 | 26 | 105 | 4.03 |
| **6044** | R3FPXDG1K7X8HC | 4 | were deliver very quickly quality function mer... | 31 | 31 | 159 | 5.12 |
| **7607** | R15BOFL2F3PHGF | 1 | do work base unit do work when arrive very poo... | 17 | 17 | 71 | 4.17 |

## Rare words removal

```
In [113]:    1  rare40 = freq_comm.tail(40)
             2  rare40
```

```
Out[113]:  extra           1
           fool            1
           cx              1
           hit             1
           code            1
           shop            1
           selection       1
           laptop          1
           avoid           1
           else            1
           mids            1
           powerful        1
           discomfort      1
           minimum         1
           whatever        1
           less            1
           forward         1
           bear            1
           amazing         1
           shipment        1
           anymore         1
           yes             1
           tangle          1
           occasional      1
           apropriate      1
           aa              1
           acre            1
           improvement     1
           prompt          1
           retire          1
           nothing         1
           shopping        1
           between         1
           distinct        1
           digital         1
           space           1
           throw           1
           eart            1
           determine       1
           kill            1
           dtype: int64
```

```
In [114]:    1  df['review_body'] = df['review_body'].apply(lambda x: ' '.join([t for t
```

In [115]:
```
1  df.sample(5)
```

Out[115]:

| | review_id | star_rating | review_body | wordcounts | word_counts | char_counts | avg_wor |
|---|---|---|---|---|---|---|---|
| **7158** | R35YV3AB2DV5IV | 4 | only gripe about these that instruction set up... | 101 | 101 | 474 | 4.6! |
| **7926** | RWUQH0N1HHZ0B | 4 | been very satisfied with cordless purchase hus... | 63 | 63 | 283 | 4.4! |
| **2679** | RBT5VQLABG9X3 | 5 | top quality recommend | 3 | 3 | 19 | 6.3: |
| **7255** | R2IHRCMRRVVZOY | 4 | buy this item few month ago because tend stay ... | 88 | 88 | 397 | 4.5 |
| **2893** | R38TG1Y9KFO7I | 5 | very very good | 3 | 3 | 12 | 4.0( |

## Word Cloud Visualization

In [116]:
```
1  # !pip install wordcloud
```

In [117]:
```
1  from wordcloud import WordCloud
2  import matplotlib.pyplot as plt
3  %matplotlib inline
```

In [118]:
```
1  text = ' '.join(df['review_body'])
```

In [119]:
```
1  len(text)
```

Out[119]: 29799

In [120]:
```python
wc = WordCloud(width=800, height=400).generate(text)
plt.imshow(wc)
plt.axis('off')
plt.show();
```



Other clustering techniques were expiremented with below.

```python
In [ ]:   1  # from gensim import corpora, models
          2  # tfidf = models.TfidfModel(bow_corpus)
          3  # corpus_tfidf = tfidf[bow_corpus]
          4  # from pprint import pprint
          5  # for doc in corpus_tfidf:
          6  #     pprint(doc)
          7  #     break
          8
          9
         10  # lda_model = gensim.models.LdaMulticore(bow_corpus, num_topics=4, id2w
         11
         12  # for idx, topic in lda_model.print_topics(-1):
         13  #     print('Topic: {} \nWords: {}'.format(idx, topic))
         14
         15
         16  # lda_model_tfidf = gensim.models.LdaMulticore(corpus_tfidf, num_topics
         17  # for idx, topic in lda_model_tfidf.print_topics(-1):
         18  #     print('Topic: {} Word: {}'.format(idx, topic))
         19
         20
         21  # from sklearn.cluster import DBSCAN
         22  # from sklearn import metrics
         23  # from sklearn.preprocessing import StandardScaler
         24
         25
         26
         27  # from sklearn.cluster import AffinityPropagation
         28
         29
         30  # # ##################################################################
         31  # # Compute Affinity Propagation
         32  # af = AffinityPropagation(random_state=42).fit(word_vecs_scaled)
         33  # cluster_centers_indices = af.cluster_centers_indices_
         34  # labels = af.labels_
         35
         36  # n_clusters_ = len(cluster_centers_indices)
         37
         38  # print("Estimated number of clusters: %d" % n_clusters_)
         39
         40
         41
         42
         43  # db = DBSCAN(eps=21, min_samples=3).fit(word_vecs_scaled)
         44  # core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
         45  # core_samples_mask[db.core_sample_indices_] = True
         46  # labels = db.labels_
         47
         48  # # Number of clusters in labels, ignoring noise if present.
         49  # n_clusters_ = len(set(labels)) - (1 if -1 in labels else 0)
         50  # n_noise_ = list(labels).count(-1)
         51
         52  # print("Estimated number of clusters: %d" % n_clusters_)
         53  # print("Estimated number of noise points: %d" % n_noise_)
         54  # print("Homogeneity: %0.3f" % metrics.homogeneity_score(labels_true, l
         55  # print("Completeness: %0.3f" % metrics.completeness_score(labels_true,
         56  # print("V-measure: %0.3f" % metrics.v_measure_score(labels_true, label
```

```
57   # print("Adjusted Rand Index: %0.3f" % metrics.adjusted_rand_score(labe
58   # print(
59   #       "Adjusted Mutual Information: %0.3f"
60   #       % metrics.adjusted_mutual_info_score(labels_true, labels)
61   # )
62   # print("Silhouette Coefficient: %0.3f" % metrics.silhouette_score(X, l
```

In [ ]:

```
1   # import gensim
2
3
4   # def get_aspect_gensim_dict(aspects):
5   #     aspect_freq_map = defaultdict(int)
6   #     for asp in aspects:
7   #         aspect_freq_map[asp] += 1
8   #     return aspect_freq_map
9
10
11  # from gensim.utils import simple_preprocess
12
13  # def preprocess(text):
14  #     result = []
15  #     for token in gensim.utils.simple_preprocess(text):
16  #         if token not in gensim.parsing.preprocessing.STOPWORDS and le
17  #             result.append(token)
18  #     return result
19
20
21  # df_single_aspect['gensim_tokens'] = df_single_aspect['asp'].apply([pr
22
23  # dictionary = gensim.corpora.Dictionary(df_single_aspect['gensim_token
24
25
26  # count = 0
27  # for k, v in dictionary.iteritems():
28  #     print(k, v)
29  #     count += 1
30  #     if count > 10:
31  #         break
32
33  # dictionary.filter_extremes(no_below=15, no_above=0.5, keep_n=10000)
34
35
36  # bow_corpus = [dictionary.doc2bow(doc) for doc in df_single_aspect['ge
37
38  # from gensim import corpora, models
39  # tfidf = models.TfidfModel(bow_corpus)
40  # corpus_tfidf = tfidf[bow_corpus]
41  # from pprint import pprint
42  # for doc in corpus_tfidf:
43  #     pprint(doc)
44  #     break
45
46  # lda_model = gensim.models.LdaMulticore(bow_corpus, num_topics=4, id2w
47
48  # for idx, topic in lda_model.print_topics(-1):
49  #     print('Topic: {} \nWords: {}'.format(idx, topic))
50
51
52  # lda_model_tfidf = gensim.models.LdaMulticore(corpus_tfidf, num_topics
53  # for idx, topic in lda_model_tfidf.print_topics(-1):
54  #     print('Topic: {} Word: {}'.format(idx, topic))
```

In [ ]:    1

In [ ]:    1

In [ ]:    1

In [ ]:    1