

Data Analysis on Product Review System: A Comprehensive Evaluation Model

Summary

E-commerce is growing at an unprecedented rate all over the world. In the process of purchasing goods, ratings and reviews play a vital reference role. Companies are pursuing a comprehensive and understandable analysis of online market data to craft greater success.

In this paper, we seek to devise several approaches to analyze the product evaluation by exploring ratings, text-based reviews and other related indicators.

We first perform exploratory data analysis by generating the data quality report and studying the distribution of the most critical indicators. Then we preprocess reviews by a series of steps, including removing punctuations, converting abbreviations, etc. Besides, we extract the frequent features of each product using **WordCloud**.

We then build PRMP, a framework that defines the patterns, relationships, measures, and parameters within and between ratings and reviews. We use **SentiWordNet** to obtain the sentiment of a review and normalize star ratings and helpfulness. Through **association analysis**, we visualize the relationship using a set of heat maps and draw conclusions from them.

After that, we propose a new approach to find a traceable measure for the product. We use **Entropy Weight Method** to obtain weights of the indicators. To identify reputation trends over time, we use the **ARIMA Model** to fit the reputation score. We select one of the hair dryer products as our main study object, calculate its score over time and give its most likely trending result. We use the **Non-linear Programming Model** to find the best combination of text-based and rating-based measures to indicate potential success and failure. We apply the **Hovland Persuasion Model** to build a decision model that describes the indicators that influence customer decisions, achieving the combination of the theory of social psychology and real-life context. And to analyze specific quality descriptors, we classify words into eight categories using the **NRC Emotion Lexicon**. Then we match the emotional intensity with star ratings and find the relationship between them.

Finally, based on the established model and detailed analysis, we put forward practical suggestions for Sunshine Company's marketing plan to improve its product competitiveness.

Contents

1	Introduction	2
1.1	Background	2
1.2	Problem Restatement and Our Work	2
2	Assumption and Data Exploration	3
2.1	Assumption	3
2.2	Data Quality Report	3
2.3	Review Preprocessing	3
2.4	Preliminary Insights Into the Data	4
3	PRMP Model	4
3.1	Review	5
3.1.1	Sentiment	5
3.1.2	Readability and Length	6
3.2	Star Rating	6
3.3	Helpfulness	6
3.4	PRMP Model and Result	7
4	Comprehensive Product Evaluation Model	8
4.1	Product Measurement	8
4.1.1	Entropy Weight Method to Obtain α_1 and α_2	9
4.1.2	Defining $\omega_1, \omega_2, \omega_3$	9
4.1.3	Scaling HVR	10
4.2	Reputation Trend Model using ARIMA	10
4.3	Non-linear Programming Model	12
4.4	Customer Decision Model	13
4.5	Correlation Between Emotional Descriptors and Star Rating	13
5	Sensitivity Analysis	15
6	Strengths and Weaknesses	15
6.1	Strengths	15
6.2	Weaknesses	16
	Letter	17
	References	17
	Appendices	19

1 Introduction

1.1 Background

Globally, more than 50% of e-commerce sales were made through online marketplaces in 2019, contributing \$1.7 trillion to the economy each year.[1] E-commerce is growing at an unprecedented rate. Amazon, as the top online marketplace, provides customers with an opportunity to review products and express their level of satisfaction by rating. At the same time, ratings and reviews can reduce potential annoyance on the buying experience.

Companies use these data to analyze the market demand, the timing of market participation, and product optimization. The description and evaluations of products are one of the most valuable references to customers when shopping online. An analysis of ratings and reviews for similar products from Amazon is critical to the company's sales strategies and can help increase product competitiveness.

1.2 Problem Restatement and Our Work

Sunshine Company is planning to introduce and sell three new products in the online marketplace. Therefore, online sales strategy, analysis of product reviews are needed to craft success in the future.

We are provided with data on ratings, reviews for microwave ovens, baby pacifiers and hairdryers sold in Amazon in recent years.

In this paper, we broke our work into sections as follows.

1. We analyze the three data sets thoroughly and generate the data quality report.
2. We preprocess text-based reviews to fit our natural language processing model better.
3. We use SentiWordNet to obtain the sentiment of the review and then find the relationships between and within star ratings, reviews, and helpfulness ratings.
4. Based on the product evaluation on Amazon, we create measures to analyze the reputation of a product and used the ARIMA model to see how it changes over time.
5. We use a non-linear programming model to explore the combinations of dimensions that best indicate a potentially successful or failing product.
6. Applying Hovland Persuasion Model, we perform a decision analysis on motivation for a person to change a review.
7. To analyze specific quality descriptors, we use the NRC Emotion Lexicon to obtain the intensity of eight categories of word sentiment, including anger, fear, trust, etc. Then we match the emotional intensity with star ratings to see if they were somehow related.
8. Finally, based on the established model and detailed analysis, we put forward practical suggestions for Sunshine Company's marketing director to improve their product competitiveness.

2 Assumption and Data Exploration

2.1 Assumption

- The product evaluations are all reasonable. That is, there is no extreme situation like low star rating and great review. We found it rare when checking the dataset.

2.2 Data Quality Report

The data quality report is the most important tool of the data exploration process. The tabular reports are accompanied by data visualizations that illustrate the distribution of the values in each feature in an **ABT**(analytical base table) [1].

Data quality report of **hair_dryer.tsv** is shown as **Table 1**.

Table 1: Analytical Base Table

	customer_id	product_parent	star_rating	helpful_votes	total_votes
<i>Count</i>	11470.00	11470.00	11470	11470	11470
<i>Mean</i>	28151220.00	484633800.00	4.116042	2.179076	2.563296
<i>Std</i>	152387700.00	287324000.00	1.300333	14.241304	15.38258
<i>Min</i>	12464.00	42396.00	1	0	0
<i>Lower Quartile</i>	14914410.00	235106000.00	4	0	0
<i>Media</i>	27071230.00	486774000.00	4	0	0
<i>Upper Quartile</i>	42336440.00	732252300.00	5	1	1
<i>Max</i>	53096370.00	999436600.00	5	499	575

From this report, it is clear that there is no continuous values are missing, so there is no need to fit any missing values. Star_rating is quite skewed and mostly concentrated on the 5. Meanwhile, there are also some extreme values in 'helpful_votes' and 'total_votes' that are worth exploring.

All data reports from the three data sets are similar, so we won't go into details here.

2.3 Review Preprocessing

For the categorical features, we selected four main features which are 'vine', 'verified_purchase', 'review_headline' and 'review_body' as our main targets. The first two are essentially boolean values but 'review_headline' and 'review_body' are long strings and need to be modified.

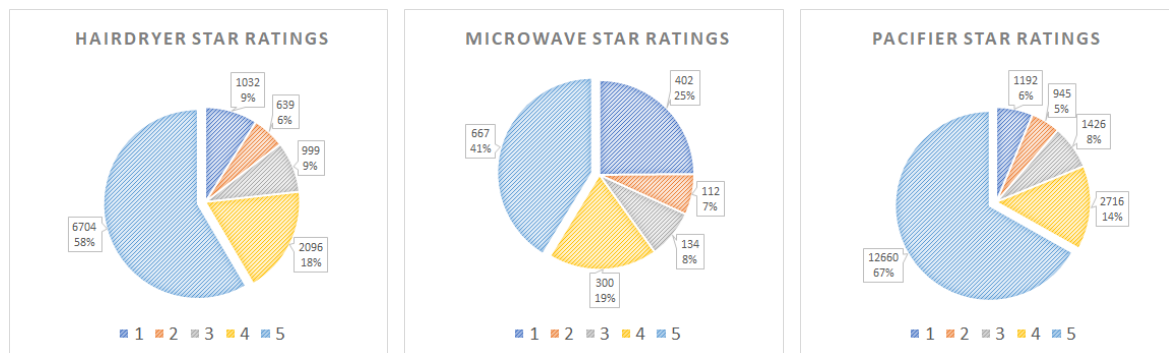
Here are our preprocessing steps for 'review_headline' and 'review_body':

1. Begin by removing dirty strings in the reviews, such as html tag '
'.
2. Convert English abbreviations to full expressions, such as changing 'isn't' into 'is not'.
3. Remove any punctuations, alpha-numeric and any other words in the limited set of special characters like ',', or '.' or '#' etc.

4. Cut a sentence into separate words, remove stopwords and convert the remain words to lowercase.
5. When processing review headlines, add a noun after each word in case the part of speech is misjudge .

2.4 Preliminary Insights Into the Data

Star Rating.After our statistical analysis, we can find that the 5-star ratings of the three categories of products are all in the majority. But the difference is that the star rating of the microwave oven product is polarized to 1 star and 5 stars. To avoid the potential extremes that may occur in a small number of evaluations, we selected MICROWAVE as a test target for our models.



(a) Star Distribution of Hair dryer (b) Star Distribution of Microwave (c) Star Distribution of Pacifier

Figure 1: Star Distribution of Each Product

Review.From the reviews of the three products, we can find the adjective words sifted by WordCloud. The first words that catch the eye are 'great', 'good', 'little' and 'new', then the adjective words related to the specific product. We also generate the feature word cloud of each product in the letter to Sunshine Company.



(a) Review Wordcloud of Hair dryer (b) Review Wordcloud of Microwave (c) Review Wordcloud of Pacifier

Figure 2: Review Wordcloud of Each Product

3 PRMP Model

The evaluation of a product has a significant impact on whether other customers purchase the product. Star ratings, reviews, and helpfulness ratings are the most valuable references to customers when shopping online, so we consider them as our primary measure factors.

3.1 Review

3.1.1 Sentiment

To analyze customer reviews, we mainly focused on analyzing opinions in the reviews and categorizing the reviews into positive or negative based on customer sentiment.

SentiWordNet [2] is the result of the automatic annotation of all the synsets of WORD-NET according to the notions of “positivity”, “negativity”, and “neutrality”. Each synset s is associated to three numerical scores $Pos(s)$, $Neg(s)$, and $Obj(s)$ which indicate how positive, negative, and “objective” (i.e., neutral) the terms contained in the synset are.[baccianella-et-al-2010-sentiwordnet]

We found whether a review is positive or negative in the following steps:

- 1.) Preprocessing reviews using methods shown in section 2.3 and extract adjectives and coordinating conjunctions, which reveal more emotions and opinions of review.
- 2.) Using the SentiWordNet to find the positive and negative values related to each word in a review and consider the impact of its neighboring words.
- 3.) Summing up the obtained positive and negative values to calculate a net positive and net negative values related to a review.

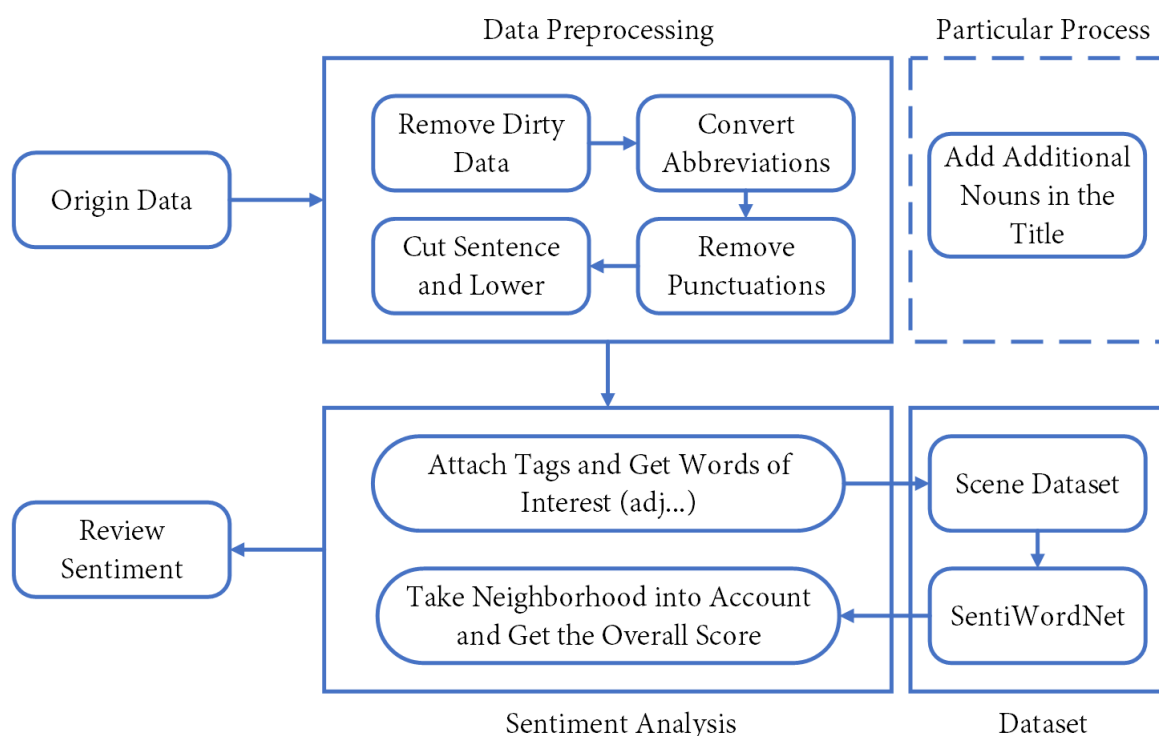


Figure 3: Flowchart of the Sentiment Processing

The sentiment of a review results in $[-1, 1]$, and we normalize it so that the result is between 0 and 1.

3.1.2 Readability and Length

Readability and length of the text are a direct expression of the intelligibility of the review text. In other words, intelligibility is judged by the difficulty of reading and understanding the relevant text.

The more understandable the review is, the more valuable it is. Therefore, intelligibility can be theoreticalized at the cognitive level, based on the cognitive adaptability of the review text to ordinary consumers.

Readability tests have been used to study the quality characteristics of several types of text in different areas of information science, and a number of readability indicators have been developed over the years. We cite the following two quantitative criteria: The Fog index measure complexity, and the ARI index measure reading ease [3].

$$FOG = 0.4 \times \left(\frac{Words}{Sentence} + 100 \times \left(\frac{N(complex_words)}{N(words)} \right) \right) \quad (1)$$

$$ARI = 4.71 \times \left(\frac{characters}{words} \right) + 0.5 \times \left(\frac{words}{sentence} \right) - 21.43 \quad (2)$$

Based on these, a quantitative study on the readability of text can be carried out. Due to the limited space of the article, we will not expand the description.

3.2 Star Rating

Star ratings tell people the quality of the product in an intuitive way, which may have an impact on customers' decisions. In general, a moderate star rating (like three stars) shows customers' neutral attitude and gives little information. What customers pay attention to are the extreme star ratings, which can convey more information.

Star ratings are originally a set of integers from 1 to 5. To better fit our model, we performed normalization on it.

Table 2: Star Normalized Value Table

star_ratings	1	2	3	4	5
after normalization	0	0.25	0.5	0.75	1

3.3 Helpfulness

Counting helpful votes is a measure of the quality of a review. With a certain total number of votes, the more helpful votes a review receives from other customers, the more valuable the review is.

We define the **Helpful Votes Ratio** of a review (**HVR**) as the number of helpful votes a review gets divided by the total number of votes.

$$HVR = \frac{helpful\ votes}{total\ votes} \quad (3)$$

If a review doesn't receive any votes at all, we define its HVR as 0.5, which means that this review was neither helpful nor unhelpful.

3.4 PRMP Model and Result

We propose a **PRMP** model that defines the patterns, relationships, measures, and parameters within and between star ratings, reviews, and helpfulness ratings.

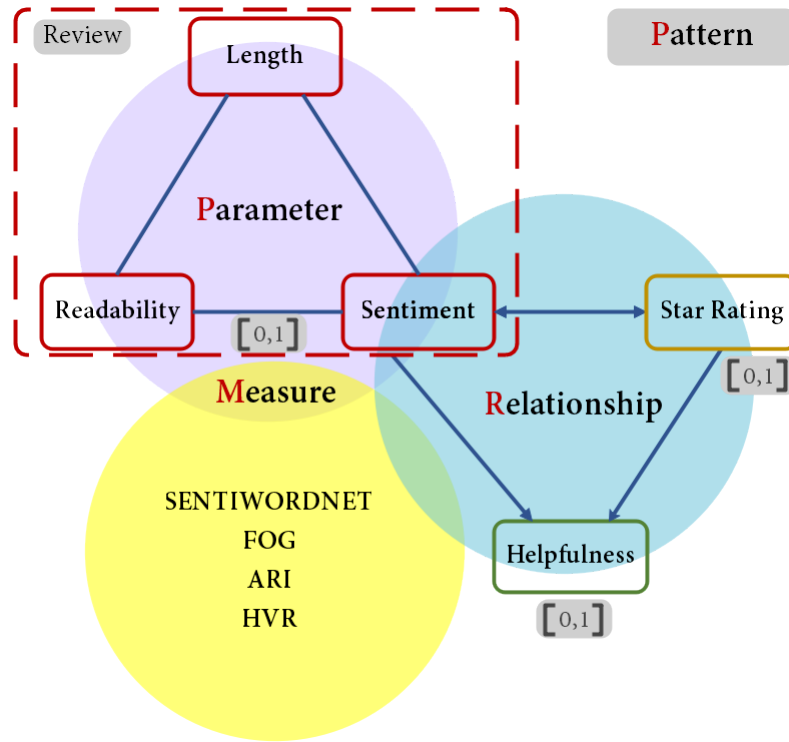


Figure 4: The Concept Graph of PRMP Model

Through the association analysis, we drew the following heat map. Star rating is closely related to the sentiment of a review and has an almost linear correlation.

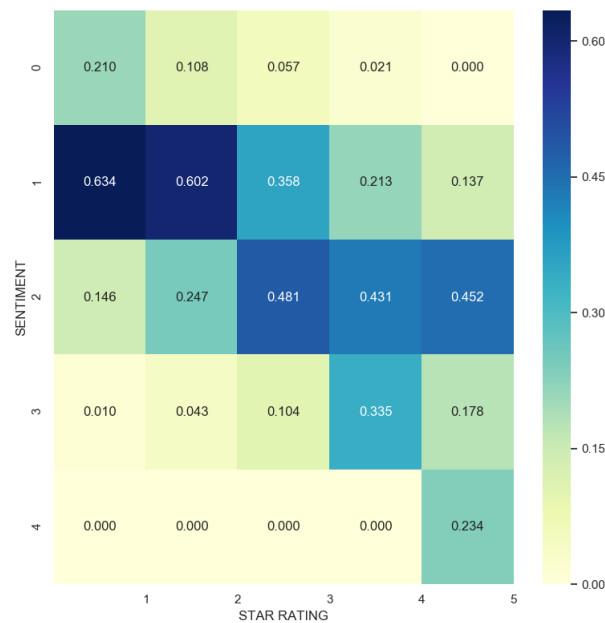
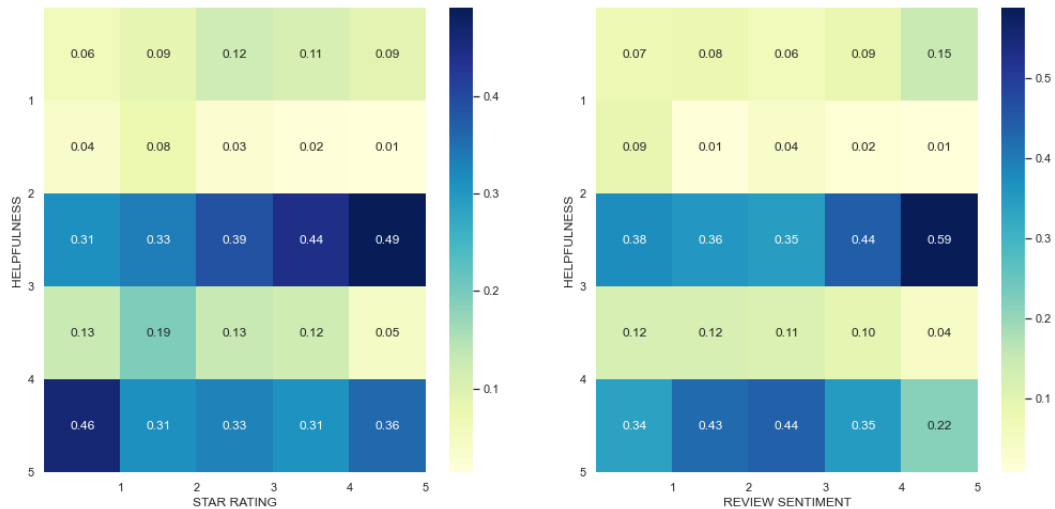


Figure 5: The Correlation Between Sentiment and Star Rating

The dark areas other than the diagonal in the heat map may have something to do with parameter adjustment.

When we performed the association analysis with helpfulness, we got different results. From the two heat maps below, the helpfulness and the other two parameters do not seem to be directly related. However, it is worth noting that in the helpfulness and star rating heat map, 1-star reviews seem to get more support from others than 5-star reviews. Low star reviews may be more helpful because of illustrating the actual problems with the product.



(a) The Correlation Between Helpfulness and Star Rating (b) The Correlation Between Helpfulness and Star Sentiment

Figure 6: The Correlation Between Analysis

Note: the reason why moderate helpfulness is the majority is that most reviews don't have any votes, so HVR equals 0.5.

4 Comprehensive Product Evaluation Model

4.1 Product Measurement

We believe that ratings and reviews give different amounts of information. By calculating the respective information entropy and assigning different weights, we get a score of a single product review.

But are all reviews equally important? Our conclusion is no. The importance of a review is also related to its helpfulness and authority. We used HVR to define the helpfulness and assign different weights for various combinations of 'vine' and 'verified', thus setting the authority.

Considering a star rating and a review can be potentially extreme, we cannot use one evaluation to judge the quality of a product. A short term overall score is needed for Sunshine Company to track.

The figure 7 shows our product measure model .

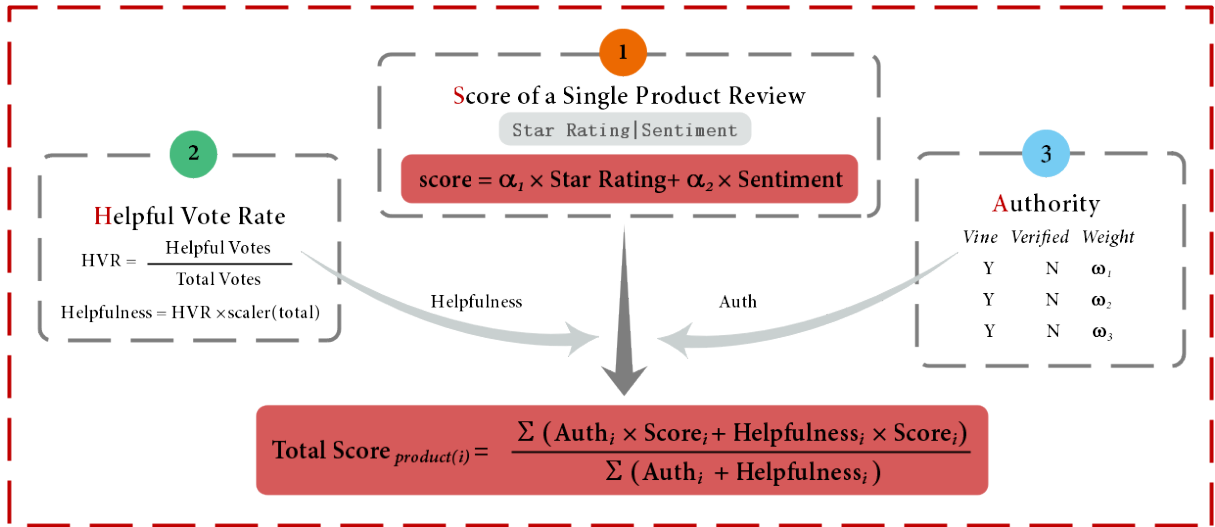


Figure 7: The Concept Graph of Product Measure Model

Since we used indicators with different units, normalization is needed to scale all values in the range [0,1]. Equation4 gives the form of data normalization.

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (4)$$

Where x_{max} and x_{min} are the maximum and minimum values of the indicator in the same unit.

4.1.1 Entropy Weight Method to Obtain α_1 and α_2

The weights of star rating and review sentiment are given using the entropy weight method.

According to the definition of information entropy in information theory, the information entropy E_j is equal to:

$$E_j = -\frac{1}{\ln(x)} \sum_{i=1}^n p_{ij} \ln(p_{ij}) \quad (5)$$

According to the calculation formula of information entropy, the information entropy of each indicator is calculated as E_1, E_2, \dots, E_k . From this we can get weights of each indicator:

$$W_j = \frac{1 - E_j}{k - \sum E_j} (i = 1, 2, \dots, k) \quad (6)$$

4.1.2 Defining $\omega_1, \omega_2, \omega_3$

Amazon Vine invites the most trusted reviewers on Amazon to post opinions about new and pre-release items to help their fellow customers make informed purchase decisions [4]. Vine Customers by Amazon are credible and professional reviewers, so such reviews are more convincing, and their evaluations are highly authentic. In addition, the rest of the reviewers are divided into two groups, depending on whether they have bought this product on Amazon. We have reason to believe that reviews written by costumers who make a purchase on Amazon are

more persuasive than those who don't. So we subjectively consider ω_1 to be 1, ω_2 to be 0.6 and ω_3 to be 0.4.

4.1.3 Scaling HVR

There are cases where there is only one vote and, at the same time being helpfulness vote. It results in a small number of total votes but a high support rate. So we adjust the **HVR** of a review to reduce the extreme impact of a small sample with high support.

We will divide the total number of votes in the short term as the denominator and divide each vote as the scaling number.

Based on the model mentioned above, we have given the following calculation formula to get a scoring formula for measuring a product in the short term.

$$TotalScore_{product(i)} = \frac{\sum(Auth_i \times Score_i + Helpfulness_i \times Score_i)}{\sum(Auth_i + Helpfulness_i)} \quad (7)$$

We calculate the overall score for andis 1875-watt hair dryer, using one month as a short term.

Table 3: The Information of the Specific Product Tested in This Paper

dataset	parent_id	product title
hairdryer.tsv	127343313	andis 1875-watt tourmaline ceramic ionic styling hair dryer

Table 4: The Total Score of Product.127343313

Date	14-Sep	14-Oct	14-Nov	14-Dec	15-Jan	15-Feb	15-Mar
Score	0.777895	0.583396	0.636676	0.559139	0.559205	0.731894	0.692061

4.2 Reputation Trend Model using ARIMA

Definition of the reputation of a company's products: Companies or consumers share various multimedia and comment on information about a product on the Internet. These discussions will affect the credibility of this product. This is the same as 'Internet word of mouth.' We believed that the reputation of a product related to many factors such as star rating, review and helpfulness rating, etc.

And we are looking for data related to the reputation of the product for Sunshine Company. In Section 4.1, we obtained the score of a product over a short period. We believe that the change of its score is related to the rise and fall of reputation, and the dimensions related to reputation are all mentioned in the 4.1 model. We correlate scores with time and build an ARIMA model.

We smooth the original time series. Here we notice that there are zero points in the obtained scores. We use cubic spline interpolation to complete the data, as the figure 8 shows.

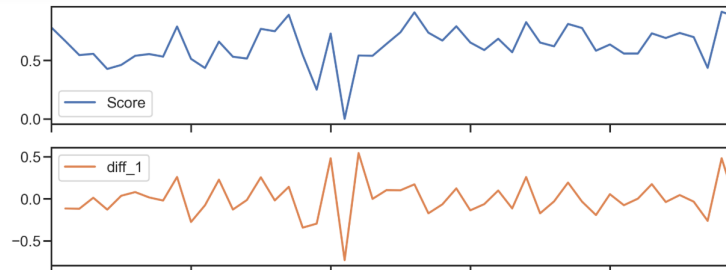


Figure 8: The Sequence of First Difference

The original line chart of the score shows that this is an unstable sequence. So we process first difference method on the sequence. Then we get $p < 0.05$, which means the data passes the test.

We then plot the ACF and PACF of this sequence to select the appropriate p and q of the ARIMA(p, q, d) model. Although there exists an autocorrelation value that exceeds the bounds, we may have p equals 1 due to accidentally exceeding the 95% confidence interval. Using the same method, according to Figure 9, we can get q equals 1.

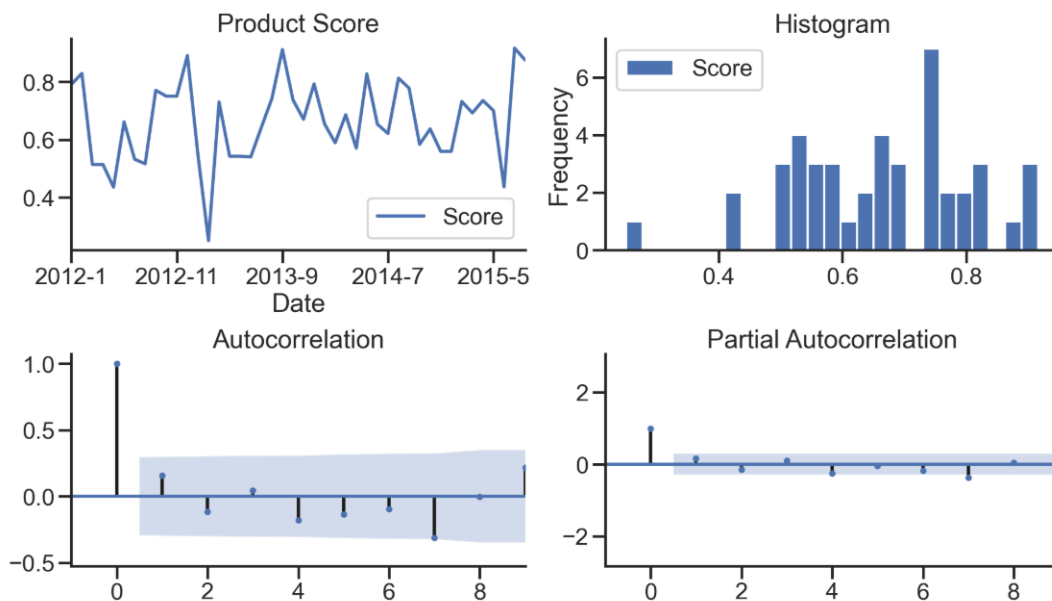


Figure 9: ACF and PACF

We use D-W test method to test this model and find it performs well.

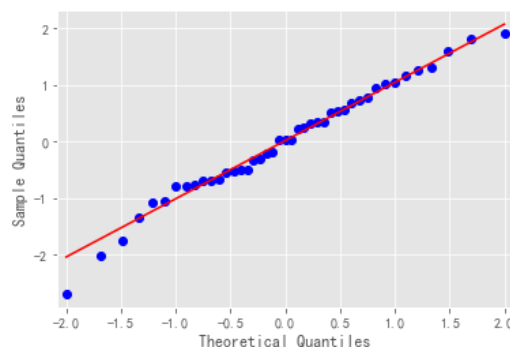


Figure 10: The Result D-W Test

Then we can observe the obvious changes in the data from the figure, and the corresponding changes in the reputation can also be clearly observed because the reputation and the score mentioned above are linked. From the figure shown below, it can be seen that reputation has a greater possibility of increasing.

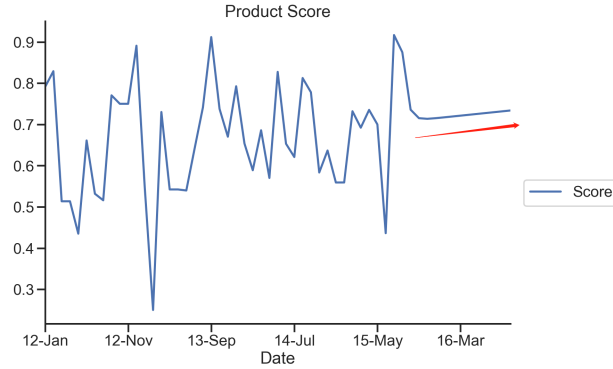


Figure 11: The Trend Generated by ARIMA

4.3 Non-linear Programming Model

As we explained in section 4.2, reputation affects the sales of online products. It means that reputation and its change over time conceal information about the quality of the product sales. Therefore, we map the success or failure of the product to the space of reputation. The base vector of this vector space contains star ratings, review score, helpfulness, etc., and their values range from zero to one.

We use a non-linear programming model to explore the combinations of dimensions that best indicate a potentially successful or failing product.

The model can be described as,

$$\begin{cases} \max & Total\ Score_{product(i)} = f(star\ rating, review\ score, helpfulness, auth, \alpha) \\ \min & Total\ Score_{product(i)} = f(star\ rating, review\ score, helpfulness, auth, \alpha) \end{cases}$$

$$\begin{cases} 0 \leq star\ rating \leq 1 \\ 0 \leq sentiment \leq 1 \\ 0 \leq helpfulness \leq 1 \\ auth \in A, A = \{\omega_1, \omega_2, \omega_3\} \\ \sum_{j=1}^2 \alpha_j = 1 \end{cases} \quad (8)$$

The objective function is iterated through the particle swarm optimization algorithm. The pseudo-code of the algorithm is as follows.

Algorithm 1: PSO in NLP

Input: $f(\mathcal{P})$, $\mathcal{P} \in \mathbb{D}^{star\ rating \times review\ score \times helpfulness \times auth \times \alpha}$, it_max , c_1 , c_2 , v
Output: $Combination(star\ rating, review\ score, helpfulness, auth, \alpha)$

- 1 **Initialize** $p = Particle_Initialization()$
- 2 **for** $i = 1$ to it_max **do**
- 3 **for each** particle p in \mathcal{P} **do**
- 4 $f_p = f(p)$
- 5 **if** f_p is better than $f(pBest)$ **then**
- 6 $pBest = p$
- 7 $gBest = best\ pin\ \mathcal{P}$
- 8 **for each** particle p in \mathcal{P} **do**
- 9 $v = v + c_1 \times rand \times (pBest - p) + c_2 \times rand \times (gBest - p)$ $p = p + v$

After several experiments, we get a pair of combinations. {star rate : 1, review sentiment : 0.75, helpfulness : 1} will best indicate the potentially successful product. While {star rate : 0, review sentiment : 0.5, helpfulness: 0.75} will best indicate the potentially failing product.

4.4 Customer Decision Model

To investigate the causal relationship between star ratings and reviews, we seek inspiration from social psychology. Hovland Persuasion Model is one of the most classic models in this field. The basic model of this approach can be described as “who said what to whom”: the source of the communication, the nature of the communication and the nature of the audience [5].

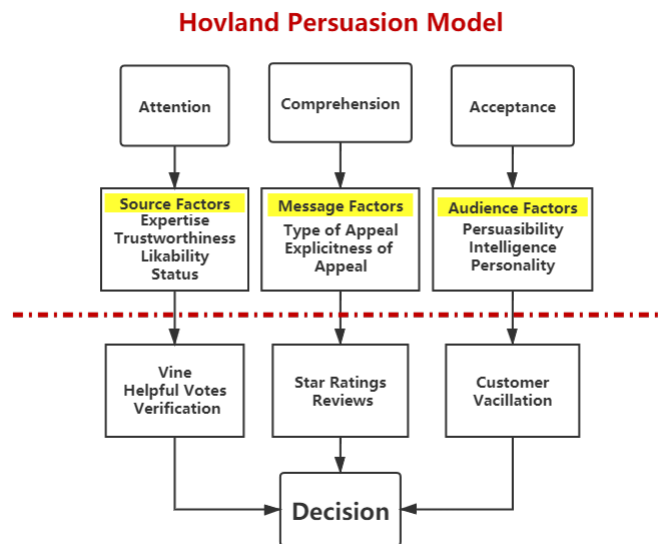


Figure 12: Hovland Persuasion Model Into Practice

4.5 Correlation Between Emotional Descriptors and Star Rating

To analyze specific quality descriptors, we classify all into eight categories using the NRC Emotion Lexicon, including anger, fear, trust, etc. The Sentiment and Emotion Lexicons is a

collection of lexicons that was entirely created by the experts of the National Research Council of Canada [6].

We then use the microwave dataset and find the sentiment distribution of all reviews. As the picture shows, the distribution of review headlines and review bodies is similar across all reviews. Words with emotions like anticipation, trust, and joy, make up the majority.

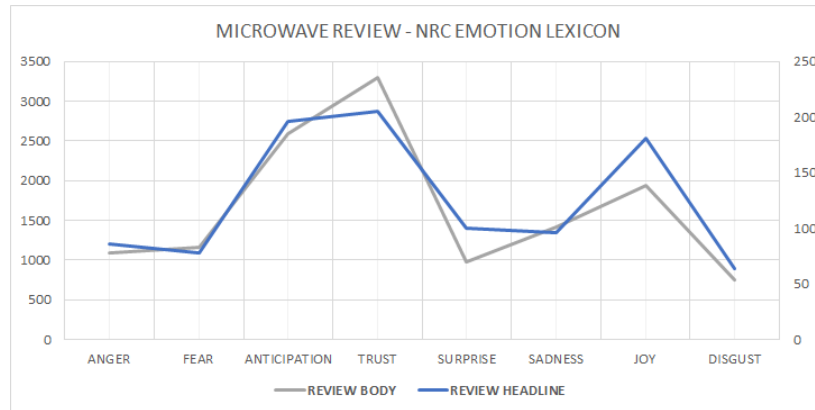


Figure 13: Microwave Review - NRC Emotion Lexicon

Generally speaking, in the bad reviews, customers can describe the expectations of the product, but there will not be too many negative descriptions in the good reviews. Therefore, negative descriptions were the focus of our attention.

We match the emotional intensity of eight categories with star ratings to see if they were somehow related. In the eight tables below, the left and right four show the relationship between the intensity of positive emotions and negative emotions and star ratings, respectively. The bars of each colour represent the frequency of that emotional intensity in each star rating level.

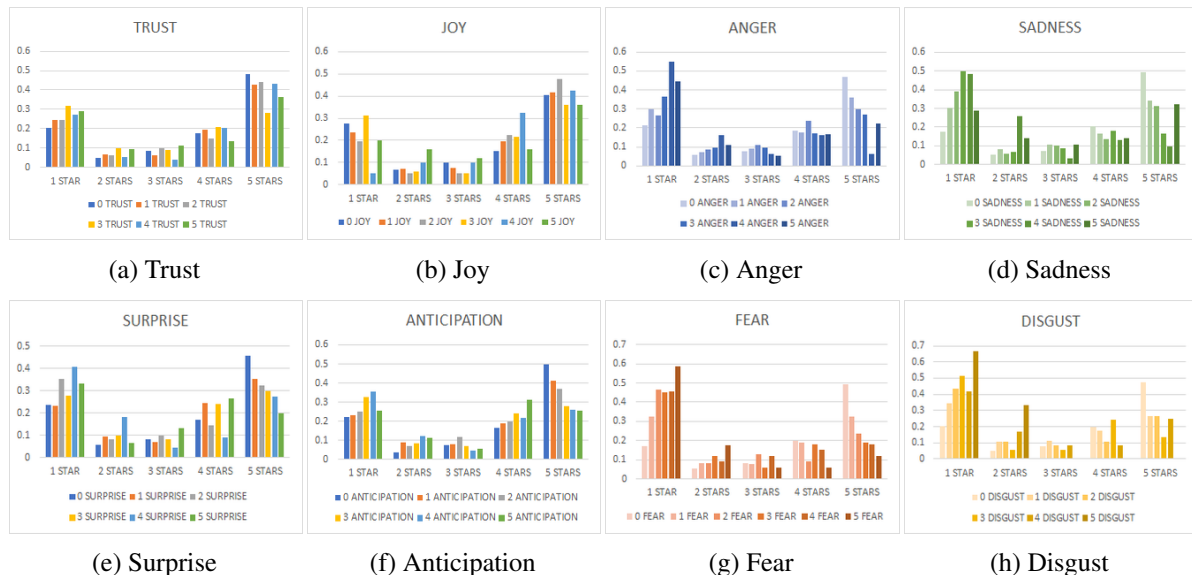


Figure 14: Emotional Intensity with Each Star Rating Level

In a specific star rating, the darker the color of the histogram is, the stronger the emotion of that star rating is.

In negative emotions, taken 'anger' as an example, reviews with angry words are more found in low stars ratings. Conversely, reviews that did not or rarely contain angry words were more observed in high star ratings. The same goes for the other three negative emotions.

In positive emotions, the distribution ratio of emotional intensity in star ratings is more uniform, especially in low star rating cases.

5 Sensitivity Analysis

An essential part of our model is the product measure in section 4.1. In our model, the weight (ω_i) is artificially judged based on experience. Changes in those weights may bring different results. We performed the sensitivity analysis by setting various combinations of weights to see the robustness of our model.

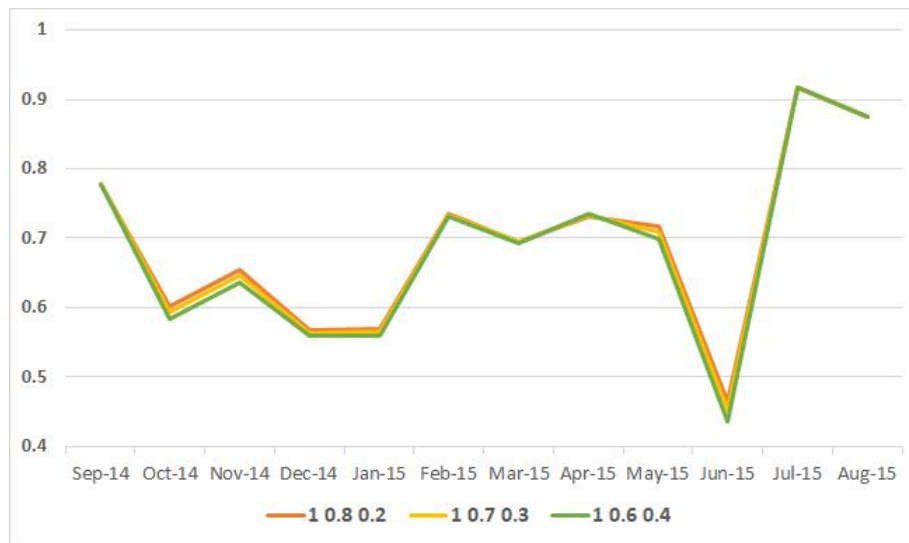


Figure 15: Product.127343313 Total Score

Each of these lines is the product score calculated in a set of given omegas. The green line is the weighted score line we originally designed. As is shown in the figure, when omega changes, the product score is only slightly affected.

6 Strengths and Weaknesses

6.1 Strengths

- We propose a new approach to combine star rating, review and helpfulness rating so that we can quantitatively calculate the overall score of a product and finally integrated them into an evaluation standard.
- We perform detailed and fully preprocessing of the review, and we use two different methods to analyze text-based review, including sentiment analysis using SentiWordNet and emotion analysis using NRC Emotion Lexicon.
- We have applied the Hovland persuasion model to our model, achieving the combination of the theory of social psychology and real-life context.
- The model is universal and can be used in other similar scenarios.
- We make detailed analysis when determining measures and giving suggestions.

6.2 Weaknesses

- In natural language processing, polysemous words such as "right" cannot be classified accurately.
- Misspelling of reviews affects the accuracy of our model.

References

- [1] J.D. Kelleher, B.M. Namee, and A. D’Arcy. Fundamentals of Machine Learning for Predictive Data Analytics: Algorithms, Worked Examples, and Case Studies. The MIT Press. MIT Press, 2015.
- [2] Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. SentiWordNet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. In Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC’10), Valletta, Malta, May 2010. European Language Resources Association (ELRA).
- [3] Nikolaos Korfiatis, Elena García-Bariocanal, and Salvador Sánchez-Alonso. Evaluating content quality and helpfulness of online product reviews: The interplay of review helpfulness vs. review content. Electronic Commerce Research & Applications, 11(3):205—217.
- [4] amazon. amazon. Website, 2020. <https://www.amazon.com/gp/vine/help>.
- [5] Timothy D. Wilson Aronson, Elliot and Robin M. Akert. Social Psychology. Upper Saddle River, N.J: Pearson Education, 2010.
- [6] National Research Council Canada. The sentiment and emotion lexicons. Website, 2016. <http://sentiment.nrc.ca/lexicons-for-research/>.

Appendices

Here are programmes we used in our model.

Analyze the sentiment of the review

```

1 import nltk
2 import re
3 import tqdm
4 import pandas as pd
5 from sklearn import preprocessing
6
7 class SentimentAnalysis(object):
8     def __init__(self, filename='SentiWordNet.txt'):
9         self.swn_pos = {'a': {}, 'v': {}, 'r': {}, 'n': {}}
10        self.swn_all = {}
11        self.build_swn(filename)
12    def geometric_weighted(self, score_list):
13        weighted_sum = 0
14        num = 1
15        for el in score_list:
16            weighted_sum += (el * (1 / float(2 ** num)))
17            num += 1
18        return weighted_sum
19    def build_swn(self, filename):
20        records = [line.split('\t') for line in open(filename)]
21        for rec in records:
22            words = rec[4].split()
23            pos = rec[0]
24            for word_num in words:
25                word = word_num.split('#')[0]
26                try:
27                    sense_num = int(word_num.split('#')[1])
28                except:
29                    continue
30                if word not in self.swn_pos[pos]:
31                    self.swn_pos[pos][word] = {}
32                self.swn_pos[pos][word][sense_num] = float(
33                    rec[2]) - float(rec[3])
34                if word not in self.swn_all:
35                    self.swn_all[word] = {}
36                self.swn_all[word][sense_num] = float(rec[2]) - float(rec
37                    [3])
38        for pos in self.swn_pos.keys():
39            for word in self.swn_pos[pos].keys():
40                newlist = [self.swn_pos[pos][word][k] for k in sorted(
41                    self.swn_pos[pos][word].keys())]
42                self.swn_pos[pos][word] = self.geometric_weighted(newlist)
43        for word in self.swn_all.keys():
44            newlist = [self.swn_all[word][k] for k in sorted(
45                self.swn_all[word].keys())]
46            self.swn_all[word] = self.geometric_weighted(newlist)
47    def score_word(self, word, pos):
48        try:
49            return self.swn_pos[pos][word]
50        except KeyError:
51            try:
52                return self.swn_all[word]
53            except KeyError:
54                return 0

```

```

54     def score(self, sentence):
55         impt = {'JJ', 'JJR', 'JJS', 'CC'}
56         non_base = {'VBD', 'VBG', 'VBN', 'VBP', 'VBZ', 'NNS', 'NNPS'}
57         word_database = pd.read_csv('data/word_database.csv', encoding="utf
            -8")
58         extend_word = word_database['extend_word']
59         limited_word = word_database['limited_word']
60         negations = {'not', 'n\'t', 'less', 'no', 'never', 'nothing', '
            nowhere', 'hardly', 'barely', 'scarcely',
61                     'nobody', 'none'}
62         stopwords = nltk.corpus.stopwords.words('english')
63         wnl = nltk.WordNetLemmatizer()
64         scores = []
65         tokens = [w.lower() for w in nltk.tokenize.word_tokenize(sentence)]
66         tagged = nltk.pos_tag(tokens)
67         new_token = []
68         new_tag = []
69         for word, pos in tagged:
70             if (pos not in impt and word not in negations and word.lower()
                not in extend_word) or word in limited_word:
71                 continue
72             else:
73                 new_tag.append((word, pos))
74                 new_token.append(word)
75         if len(tagged) == 1:
76             return self.score_word(tagged[0][0].lower(), self.pos_short(
                tagged[0][1])), 1
77         index = 0
78         for el in tagged:
79             pos = el[1]
80             try:
81                 word = re.match('(\w+)', el[0]).group(0).lower()
82                 start = index - 5
83                 if start < 0:
84                     start = 0
85                 neighborhood = tokens[start:index]
86                 if ((pos in impt) and (word not in stopwords) or word in
                    extend_word) and word not in limited_word:
87                     if pos in non_base:
88                         word = wnl.lemmatize(word, self.pos_short(pos))
89                         score = self.score_word(word, self.pos_short(pos))
90                         if word == 'if':
91                             score -= 0.3
92                         if word == 'supposed':
93                             score -= 0.2
94                         if len(negations.intersection(set(neighborhood))) > 0:
95                             r = negations.intersection(set(neighborhood))
96                             a = list(r)
97                             i_n = tokens.index(a[0])
98                             if tagged[i_n + 1][0] in impt or tagged[i_n + 1][0]
                                in extend_word:
99                                 score = -score
100                             scores.append(score)
101             except AttributeError:
102                 pass
103             index += 1
104         if len(scores) > 0:
105             return sum(scores) / float(len(scores)), len(scores)
106         else:
107             return 0, 0

```

```

108
109
110 def decontracted(phrase):
111     phrase = re.sub(r"won't", "will_not", phrase)
112     phrase = re.sub(r"can't", "can_not", phrase)
113     phrase = re.sub(r"n't", "_not", phrase)
114     phrase = re.sub(r"\ 're", "_are", phrase)
115     phrase = re.sub(r"\ 's", "_is", phrase)
116     phrase = re.sub(r"\ 'd", "_would", phrase)
117     phrase = re.sub(r"\ 'll", "_will", phrase)
118     phrase = re.sub(r"\ 't", "_not", phrase)
119     phrase = re.sub(r"\ 've", "_have", phrase)
120     phrase = re.sub(r"\ 'm", "_am", phrase)
121     return phrase
122 if __name__ == '__main__':
123     s = SentimentAnalysis(filename='SentiWordNet.txt')
124     data = pd.read_csv('data/microwave.tsv', sep='\t', header=0, encoding="
        utf-8")
125     review = data['review_body']
126     review_title = data['review_headline']
127     review_id = data['review_id']
128     star = data['star_rating']
129     star_dict = {}
130     title_reviews = {}
131     preprocessed_reviews = {}
132     score = {}
133     for index in tqdm.tqdm(range(0, data.shape[0])):
134         star_dict[review_id[index]] = star[index]
135         title_reviews[review_id[index]] = review_title[index]
136         try:
137             if '<br_>' in review[index]:
138                 tmp = review[index].copy()
139                 review[index] = tmp.replace('<br_>', '')
140                 without_short = decontracted(review[index])
141                 without_short1 = decontracted(review_title[index])
142                 without_number = re.sub("\S*d\S*", "", without_short).strip()
143                 without_number1 = re.sub("\S*d\S*", "", without_short1).strip()
144                 score1, len1 = s.score(without_number)
145                 score2, len2 = s.score(without_number1 + "_machine")
146                 if len1 == 0 and len2 == 0:
147                     continue
148                 if len2 == 0:
149                     preprocessed_reviews[review_id[index]] = score1
150                 elif len1 == 0:
151                     preprocessed_reviews[review_id[index]] = score2
152                 else:
153                     preprocessed_reviews[review_id[index]] = 0.1 * score1 + 0.9
154                     * score2
155             except:
156                 continue
157         del preprocessed_reviews[min(preprocessed_reviews, key=
            preprocessed_reviews.get)]
158         del preprocessed_reviews[max(preprocessed_reviews, key=
            preprocessed_reviews.get)]
159         min_value = preprocessed_reviews[min(preprocessed_reviews, key=
            preprocessed_reviews.get)]
160         max_value = preprocessed_reviews[max(preprocessed_reviews, key=
            preprocessed_reviews.get)]
161         ndict = {}

```

```

161     nstar = {}
162     s = pd.Series(preprocessed_reviews)
163     ndf = pd.DataFrame(s, columns=['score'])
164     std_scaler = preprocessing.StandardScaler()
165     std_label_data = std_scaler.fit_transform(ndf)
166     min_max_scaler = preprocessing.MinMaxScaler()
167     min_max_label_data = min_max_scaler.fit_transform(std_label_data)
168     index = 0
169     for key in preprocessed_reviews.keys():
170         ndict[key] = min_max_label_data[index][0]
171         if 'FiveStars' in title_reviews[key]:
172             ndict[key] = 1
173         elif 'FourStars' in title_reviews[key]:
174             ndict[key] = 0.75
175         elif 'ThreeStars' in title_reviews[key]:
176             ndict[key] = 0.5
177         elif 'TwoStars' in title_reviews[key]:
178             ndict[key] = 0.25
179         elif 'OneStar' in title_reviews[key]:
180             ndict[key] = 0
181         elif 'zero_star' in title_reviews[key]:
182             ndict[key] = 0
183         index += 1
184         nstar[key] = (star_dict[key] - 1) / 4
185     key1 = list(ndict.keys())
186     value2 = list(ndict.values())
187     value3 = list(nstar.values())
188     dataframe = pd.DataFrame({'review_id': key1, 'review_score': value2, '
189                             star_rate': value3, })
189     dataframe.to_csv("nor_score2.csv", index=False, sep=',')

```

Calculate the score of a single review

```

1  import numpy as np
2  import pandas as pd
3  import tqdm
4  import json
5
6  class NpEncoder(json.JSONEncoder):
7      def default(self, obj):
8          if isinstance(obj, np.integer):
9              return int(obj)
10         elif isinstance(obj, np.floating):
11             return float(obj)
12         elif isinstance(obj, np.ndarray):
13             return obj.tolist()
14         else:
15             return super(NpEncoder, self).default(obj)
16
17  if __name__ == '__main__':
18      # dict store origin data
19      data_tuple = {}
20      # read data
21      origin_data = pd.read_csv('data/hair_dryer.tsv', sep='\t', header=0,
22                               encoding="utf-8")
23      part_data = pd.read_csv('data/a1.csv', encoding="utf-8")
24      # prepare data
25      # get data from source data
26      review_id_origin = origin_data['review_id']
27      review_id_part = part_data['review_id']

```

```

27     # data from origin data
28     helpful = origin_data['helpful_votes']
29     total = origin_data['total_votes']
30     product_id = origin_data['product_parent']
31     date = origin_data['review_date']
32     vine = origin_data['vine']
33     verified = origin_data['verified_purchase']
34     # data from part data
35     review_score = part_data['review']
36     star_score = part_data['star']
37     review_id = list(review_id_part.values)
38     for index in tqdm.tqdm(range(0, len(origin_data.values))):
39         if review_id_origin[index] in review_id:
40             date_split = date[index].split('/')
41             index_part = review_id.index(review_id_origin[index])
42             if date_split[-1] not in data_tuple.keys():
43                 data_tuple[date_split[-1]] = {date_split[0]: {
44                     review_id_origin[index]: (
45                         str(product_id[index]), star_score[index_part],
46                         review_score[index_part], helpful[index],
47                         total[index], vine[index], verified[index], date[index]
48                     )}}
49             else:
50                 exist_data = data_tuple[date_split[-1]]
51                 if date_split[0] not in exist_data.keys():
52                     exist_data[date_split[0]] = {review_id_origin[index]: (
53                         product_id[index], star_score[index_part],
54                         review_score[index_part], helpful[index],
55                         total[index], vine[index], verified[index], date[
56                             index])}
57                 data_tuple[date_split[-1]] = exist_data
58             else:
59                 exist_data[date_split[0]][review_id_origin[index]] = (
60                     product_id[index], star_score[index_part],
61                     review_score[index_part], helpful[index],
62                     total[index],
63                     vine[index], verified[index], date[index])
64                 data_tuple[date_split[-1]] = exist_data
65 data_fre = {}
66 for key, value in data_tuple.items():
67     month = {}
68     for s_key, s_value in value.items():
69         tmp_star_dic = {}
70         tmp_review_dic = {}
71         for t_key, t_value in s_value.items():
72             if str(t_value[1]) not in tmp_star_dic.keys():
73                 tmp_star_dic[str(t_value[1])] = 1
74             else:
75                 tmp_star_dic[str(t_value[1])] += 1
76             if str(t_value[2]) not in tmp_review_dic.keys():
77                 tmp_review_dic[str(t_value[2])] = 1
78             else:
79                 tmp_review_dic[str(t_value[2])] += 1
80             month[s_key] = {'star': tmp_star_dic, 'review': tmp_review_dic}
81 data_pro = {}
82 for key, value in data_fre.items():
83     month = {}
84     for s_key, s_value in value.items():

```



```

81         tmp_star_dic = {}
82         tmp_review_dic = {}
83         sum_star = 0
84         sum_review = 0
85         for t_key, t_value in s_value['star'].items():
86             sum_star += t_value
87         for t_key, t_value in s_value['star'].items():
88             tmp_star_dic[t_key] = t_value / sum_star
89
90         for t_key, t_value in s_value['review'].items():
91             sum_review += t_value
92         for t_key, t_value in s_value['review'].items():
93             tmp_review_dic[t_key] = t_value / sum_review
94         month[s_key] = {'star': tmp_star_dic, 'review': tmp_review_dic}
95     data_pro[key] = month
96     e = {}
97     for key, value in data_pro.items():
98         month = {}
99         for s_key, s_value in value.items():
100             e_star = 0
101             e_review = 0
102             for t_key, t_value in s_value['star'].items():
103                 e_star += t_value * np.log(t_value)
104             e_star = (-1.0 / np.log(5)) * e_star
105             for t_key, t_value in s_value['review'].items():
106                 e_review += t_value * np.log(t_value)
107             e_review = (-1.0 / np.log(5)) * e_review
108             month[s_key] = {'star': e_star, 'review': e_review}
109         e[key] = month
110     w = {}
111     for key, value in e.items():
112         month = {}
113         for s_key, s_value in value.items():
114             e_star = s_value['star']
115             e_review = s_value['review']
116             w1 = (1 - e_star) / (2 - (e_star + e_review))
117             w2 = (1 - e_review) / (2 - (e_star + e_review))
118             month[s_key] = {'star': w1, 'review': w2}
119         w[key] = month
120     # model part 2
121     helpful_dic = {}
122     for key, value in w.items():
123         month = {}
124         for s_key, s_value in value.items():
125             review_item = {}
126             product_total_set = {}
127             reviews = data_tuple[key][s_key]
128             for t_key, t_value in reviews.items():
129                 if len(product_total_set) == 0:
130                     tmp_total = [t_value[4]]
131                     product_total_set[t_value[0]] = tmp_total
132                 else:
133                     if t_value[0] not in product_total_set.keys():
134                         tmp_total = [t_value[4]]
135                         product_total_set[t_value[0]] = tmp_total
136                     else:
137                         product_total_set[t_value[0]].append(t_value[4])
138             for t_key, t_value in reviews.items():
139                 if t_value[4] == 0:
140                     hvr = 0.5

```

```

141         else:
142             hvr = t_value[3] / t_value[4]
143             total_list = product_total_set[t_value[0]]
144             max_value = max(total_list)
145             min_value = min(total_list)
146             if (max_value - min_value) == 0:
147                 if max_value == 0:
148                     scale = 0
149                 else:
150                     scale = 1
151             else:
152                 scale = (t_value[4] - min_value) / (max_value -
153                     min_value)
154             review_item[t_key] = (str(t_value[0]), t_value[1], t_value
155                 [2], hvr, scale, t_value[5], t_value[6])
156             month[s_key] = review_item
157             helpful_dic[key] = month
158         # model part 3
159         right = [[0.4, 0.6], [1.0, 0]]
160         scores = {}
161         for key, value in helpful_dic.items():
162             month = {}
163             for s_key, s_value in value.items():
164                 score_item = {}
165                 product_score_set = {}
166                 product_k_set = {}
167                 alpha = w[key][s_key]
168                 for t_key, t_value in s_value.items():
169                     star_review_score = alpha['star'] * t_value[1] + alpha['
170                         review'] * t_value[2]
171                     hvr_k = t_value[3] * t_value[4]
172                     if t_value[5] == 'Y':
173                         i = 1
174                     else:
175                         i = 0
176                     if t_value[6] == 'Y':
177                         j = 1
178                     else:
179                         j = 0
180                     authority_k = right[i][j]
181                     score = authority_k * star_review_score + hvr_k *
182                         star_review_score
183                     k = authority_k + hvr_k
184                     if len(product_score_set) == 0:
185                         product_score_set[t_value[0]] = score
186                         product_k_set[t_value[0]] = k
187                     else:
188                         if t_value[0] not in product_score_set.keys():
189                             product_score_set[t_value[0]] = score
190                             product_k_set[t_value[0]] = k
191                         else:
192                             product_score_set[t_value[0]] += score
193                             product_k_set[t_value[0]] += k
194                     for t_key, t_value in product_score_set.items():
195                         score_item[t_key] = t_value / (product_k_set[t_key])
196                     month[s_key] = score_item
197             scores[key] = month
198         measure = {}
199         for key, value in scores.items():
200             for s_key, s_value in value.items():

```

```

197         for t_key, t_value in s_value.items():
198             if len(measure) == 0:
199                 month = {s_key: t_value}
200                 year = {key: month}
201                 measure[t_key] = year
202             else:
203                 if t_key not in measure.keys():
204                     month = {s_key: t_value}
205                     year = {key: month}
206                     measure[t_key] = year
207                 else:
208                     year = measure[t_key]
209                     if key not in year.keys():
210                         year[key] = {s_key: t_value}
211                         measure[t_key] = year
212                     else:
213                         year[key][s_key] = t_value
214                         measure[t_key] = year
215         with open('result/product_score_per_month_1.json', 'w') as f:
216             json.dump(measure, f, cls=NpEncoder)

```

NRC Emotion Lexicon: categorize the emotion of a word

```

1  import os
2  import sys
3  import json
4  import nltk
5  import numpy as np
6  import pandas as pd
7  from tqdm import tqdm
8  from textblob import TextBlob
9  from nltk.tokenize import sent_tokenize
10 from nltk.tokenize import TreebankWordTokenizer
11
12 data = pd.read_csv('original.csv')
13 data = pd.DataFrame(data)
14 data = data[["review_headline", "review_body"]]
15 nrc_lex = pd.read_csv("NRC-Emotion-Lexicon-Wordlevel-v0.92.txt", sep='\t')
16 def get_emotion(data, name):
17     emotions = []
18     for review in tqdm(data[name]):
19         total = [0, 0, 0, 0, 0, 0, 0, 0]
20         anger = 0
21         fear = 0
22         anticipation = 0
23         trust = 0
24         surprise = 0
25         sadness = 0
26         joy = 0
27         disgust = 0
28         if not review is np.nan:
29             lyrics_text = review
30             token_lyrics = sent_tokenize(lyrics_text)
31             for sentence in token_lyrics:
32                 lyric_words = TreebankWordTokenizer().tokenize(sentence)
33                 for word in lyric_words:
34                     anger_list = nrc_lex[nrc_lex['word'] == word][nrc_lex['emotion'] == 'anger'].index.tolist()
35                     if len(anger_list) == 1:
36                         anger += int(nrc_lex.iloc[int(anger_list[0])][

```

```

        association'])
37     fear_list = nrc_lex[nrc_lex['word'] == word][nrc_lex['
        emotion'] == 'fear'].index.tolist()
38     if len(fear_list) == 1:
39         fear += int(nrc_lex.iloc[int(fear_list[0])][ '
        association'])
40     anticipation_list = nrc_lex[nrc_lex['word'] == word][
41         nrc_lex['emotion'] == 'anticipation'].index.tolist
        ()
42     if len(anticipation_list) == 1:
43         anticipation += int(nrc_lex.iloc[int(
        anticipation_list[0])][ 'association'])
44     trust_list = nrc_lex[nrc_lex['word'] == word][nrc_lex['
        emotion'] == 'trust'].index.tolist()
45     if len(trust_list) == 1:
46         trust += int(nrc_lex.iloc[int(trust_list[0])][ '
        association'])
47     surprise_list = nrc_lex[nrc_lex['word'] == word][
        nrc_lex['emotion'] == 'surprise'].index.tolist()
48     if len(surprise_list) == 1:
49         surprise += int(nrc_lex.iloc[int(surprise_list[0])
        ][ 'association'])
50     sadness_list = nrc_lex[nrc_lex['word'] == word][nrc_lex
        ['emotion'] == 'sadness'].index.tolist()
51     if len(sadness_list) == 1:
52         sadness += int(nrc_lex.iloc[int(sadness_list[0])][ '
        association'])
53     joy_list = nrc_lex[nrc_lex['word'] == word][nrc_lex['
        emotion'] == 'joy'].index.tolist()
54     if len(joy_list) == 1:
55         joy += int(nrc_lex.iloc[int(joy_list[0])][ '
        association'])
56     disgust_list = nrc_lex[nrc_lex['word'] == word][nrc_lex
        ['emotion'] == 'disgust'].index.tolist()
57     if len(disgust_list) == 1:
58         disgust += int(nrc_lex.iloc[int(disgust_list[0])][ '
        association'])
59     total = [anger, fear, anticipation, trust, surprise, sadness,
        joy, disgust]
60     emotions.append(total)
61     else:
62         emotions.append(total)
63     return emotions
64 def gather(data, emotions):
65     result = {"id": [], "anger": [], "fear": [], "anticipation": [], "trust
        ": [], "surprise": [], "sadness": [],
66         "joy": [], "disgust": []}
67     i = 0
68     data = pd.DataFrame(data[["review_id"]])
69     print(data)
70     for row in data.values:
71         result["id"].append(row[0])
72         result["anger"].append(emotions[i][0])
73         result["fear"].append(emotions[i][1])
74         result["anticipation"].append(emotions[i][2])
75         result["trust"].append(emotions[i][3])
76         result["surprise"].append(emotions[i][4])
77         result["sadness"].append(emotions[i][5])
78         result["joy"].append(emotions[i][6])
79         result["disgust"].append(emotions[i][7])

```

```
80         i += 1
81     result = pd.DataFrame(result)
82     return result
83 if __name__ == '__main__':
84     data_tmp = pd.read_csv('original.csv')
85     data_tmp = pd.DataFrame(data_tmp)
86     result_body = gather(data_tmp, emotion_body)
87     result_body.to_csv('result_body.csv', index=None, sep=',')
88     result_title = gather(data_tmp, emotion_title)
89     result_title.to_csv('result_title.csv', index=None, sep=',')
90     emotion_body = get_emotion(data, 'review_body')
91     emotion_title = get_emotion(data, 'review_headline')
92     data_tmp = pd.read_csv('original.csv')
93     data_tmp = pd.DataFrame(data_tmp)
94     result_body = gather(data_tmp, emotion_body)
95     result_body.to_csv('result_body.csv', index=None, sep=',')
96     result_title = gather(data_tmp, emotion_title)
97     result_title.to_csv('result_title.csv', index=None, sep=',')
```
