

## Beat The Market: Comprehensive Exploration of Amazon Reviews and Ratings

Recently the Sunshine Company plans to launch three new products in the online marketplace, and our team is required to provide some insights of the given data, as well as giving strategic suggestions for improving the future sales and reputation of the products. Specific tasks are separated into two big questions.

As for question 1, we first apply data cleaning by removing unnecessary information, **tokenizing all texts and lemmatizing and stemming all words**. Then we apply **LDA topic model** to give an intuitive description of what reviews care about. Next we visualize the relationship of review amount, review length, star rating and the helpful votes by time and cross analysis. The results show that reviews with larger helpfulness ratings tend to accompany with high ratings and long review lengths. Moreover, **in the early stage the averaged number of helpful votes per review far outnumber that of the latter stage. In the mean time, the ratings, review lengths and other indexes dramatically fluctuate**. These all indicate the importance of keeping a good brand image in the cold start stage.

As to question 2(a), we propose three metrics for Sunshine to track: 1. **weighted rating ratio** which represents the occurring ratio of each rating weighted by number of helpful votes; 2. **weighted sentimental score for reviews**, where we apply **logistic regression** to calculate scores for each notional word and their weighted sums by helpfulness; 3. **preference vector**, where we sort out seven attributes for each product based on LDA's results , set up dictionaries which consist of related terms for each attribute and estimate people's preference ratio on these attributes based on weighted word frequency statistics with time decay. As a Result, Sunshine can allocate different efforts on improving different product features.

In question 2(b), we consider that the **reputation** of a product is related to the averaged star ratings, the authoritativeness of its reviews and the sales volume, among which we assume the sales volume is proportional to the review number in a fixed-length time window. Therefore the reputation over that fixed-length time window can be seen as the joint contribution of features concerning rate and reviews during that period, and after calculation, results show that for hair dryer and pacifier, their reputation scores increase in the early time and tend to be stable later, whereas those of microwave keep growing the whole time.

With regard to 2(c), we continue to use the time-varying **reputation** in 2(b) as a comprehensive indicator of both ratings and reviews and apply a **nested two-layer LSTM model** to predict its value for the review sequence, for considering that **this index takes nearly every informative given features into accounts (sales volume included)**.

In terms of 2(d), we consider the **ripple effect** of reviews. After analyzing the trend of the monthly averaged number of different ratings over time, we conclude that **reviews with rate 5 tend to incite more reviews**. Besides that, we unexpectedly observe that the reviews amount of rate 1 is significantly correlated with the length of reviews, and after **Granger Cause Test**, we find that **short reviews tend to lag within two months after a large ratio of low star ratings**.

Regarding 2(e), we sort out a dictionary for **affective words exclusively** and assign a score for each of them. Then we compute the new sentimental scores of all reviews and rank reviews into five ranks. After comparing the **confusion matrix** of star rating and review ranks, we find a **mapping asymmetry** that some reviews with strongly affective words had mild ratings, and reviews with more mildly affective words has extreme ratings. We explain this phenomena by visualizing that for those reviews, greater chances are that **both positive words and negative affective words occur, or strongly affective words are replaced by words that describe product attributes**.

Last but not least, we summarize the pros and cons of our solutions and present our insights to Sunshine's market director for the purpose of helping Sunshine take a lead in online market.

**Key Word:** logistic regression; LSTM; LDA topic model; Granger Cause Test; sentiment analysis



关注校苑数模公众号

获取更多资料

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Problem Background . . . . .	2
1.2	Clarification and Restatement . . . . .	2
1.3	Our Work . . . . .	2
<b>2</b>	<b>Problem 1: Data Preprocessing and Mining</b>	<b>4</b>
2.1	Data Cleaning . . . . .	4
2.2	Text Mining by LDA Topic Model[1] . . . . .	4
2.3	Overall Data Characteristics . . . . .	5
<b>3</b>	<b>Problem 2 (a): Ratings and Reviews Based Data Measures</b>	<b>8</b>
3.1	Weighted Rating Ratio . . . . .	8
3.2	Avg and Std of Weighted Sentimental Scores for Reviews[2][3][4] . . . . .	8
3.3	Users' Preference Vector . . . . .	9
<b>4</b>	<b>Problem 2 (b): Reputation Metric</b>	<b>11</b>
4.1	Reputation Metric . . . . .	11
4.2	Analysis on Reputation Variation Patterns . . . . .	11
<b>5</b>	<b>Problem 2 (c): Nested Two-Layer LSTM</b>	<b>12</b>
5.1	The Structure of the Nested Two-layer LSTM Model . . . . .	12
5.2	Analysis on Potential Success of Products . . . . .	14
<b>6</b>	<b>Problem 2 (d): Causal Effectiveness Between Reviews</b>	<b>15</b>
6.1	Ripple Effects of Extreme Ratings . . . . .	15
6.2	Causal Inference for Ratio of Low Rating and Review Length . . . . .	16
<b>7</b>	<b>Problem 2 (e)[5]: Correlation between Affective Words and Star Ratings</b>	<b>18</b>
7.1	Analysis for Alignment of Rate and Review Scored by Certain Affective Words . . . . .	18
7.2	Micro Observation of Asymmetry between Rate and Review[6] . . . . .	18
<b>8</b>	<b>Strengths and Weaknesses</b>	<b>20</b>
8.1	Strengths . . . . .	20
8.2	Weaknesses . . . . .	21
<b>9</b>	<b>Conclusion</b>	<b>21</b>
<b>10</b>	<b>The Letter to the Marketing Director of Sunshine Company</b>	<b>22</b>
<b>Appendices</b>		<b>26</b>
<b>Appendix A</b>	<b>LDA Topic Model for Microwave and Pacifier</b>	<b>26</b>
<b>Appendix B</b>	<b>Product Contrast of Microwave and Pacifier Over Time</b>	<b>27</b>
<b>Appendix C</b>	<b>Code</b>	<b>27</b>
C.1	Data Preprocess and Overall Analysis . . . . .	27
C.2	LDA Analysis . . . . .	33
C.3	Sentimental Score . . . . .	37
C.4	LSTM . . . . .	38
C.5	Reputation Model . . . . .	40



关注校苑数模公众号

获取更多资料

## 1 Introduction

### 1.1 Problem Background

Nowadays with more and more merchants joining the e-commerce camp, competition in Internet marketing has become increasingly intense. Meanwhile, as more and more customers participate in online reviews and interactions, comprehensive analysis of reviews and ratings plays an increasingly important role in understanding customers' pain points and specifying future strategies.

Specifically speaking, in Amazon's design, customers are allowed to choose a number from 1 to 5 to express their satisfaction level with the product, write any text-based messages as reviews and freely vote for other reviews they consider helpful. These are the main data sources for related companies to gain insights into the markets in which they participate, the timing of that participation, and the potential success of product design feature choices.

### 1.2 Clarification and Restatement

In this problem we are given three review data sets on three new products by Sunshine Company: a microwave oven, a baby pacifier and a hair dryer, and asked to provide the Company with an online sales strategy and important design patterns to enhance the attractiveness of the product. The data sets include product types, star ratings, review titles and bodies, helpfulness votes, certification information and date. We should only use these data to solve the following problems:

- Analyze the three product data sets, qualitatively or quantitatively measure and describe these data to obtain some quantifiable indicators, patterns or relationships to guide product sales plan of Sunshine.
- Based on the above analysis, the following special needs should be further addressed:
  - Design a metric based on the rating and review content to represent the information value of the review, thus providing Sunshine Company with a way to find the most valuable review.
  - Find time-based measures and patterns which indicate whether the market reputation of specific product is increasing or decreasing.
  - Identify a combination of text-based measures and ratings-based measures that best indicate the potential success or failure of a product.
  - Find correlations between star ratings and review properties such as quantity, quality, etc.
  - Find correlations between star ratings and specific review words.
- Write a letter to the Marketing Director of Sunshine Company by summarizing the analysis and results, and providing rational recommends.

### 1.3 Our Work

Our work follows the following workflow, as shown in the Fig. 1. The level-1 data is the 15 product properties directly provided in the problem. The level-2 data is based on the level-1 data, including product class, helpfulness rating and emotion measure of the review content. Based on this, level-3 data is further processed into a score that can be used to measure a review's information value. Then level-4 data, which indicates a reasonable measure of a product's prospects, is aggregated for a reputation index in problem 2(b)(c).



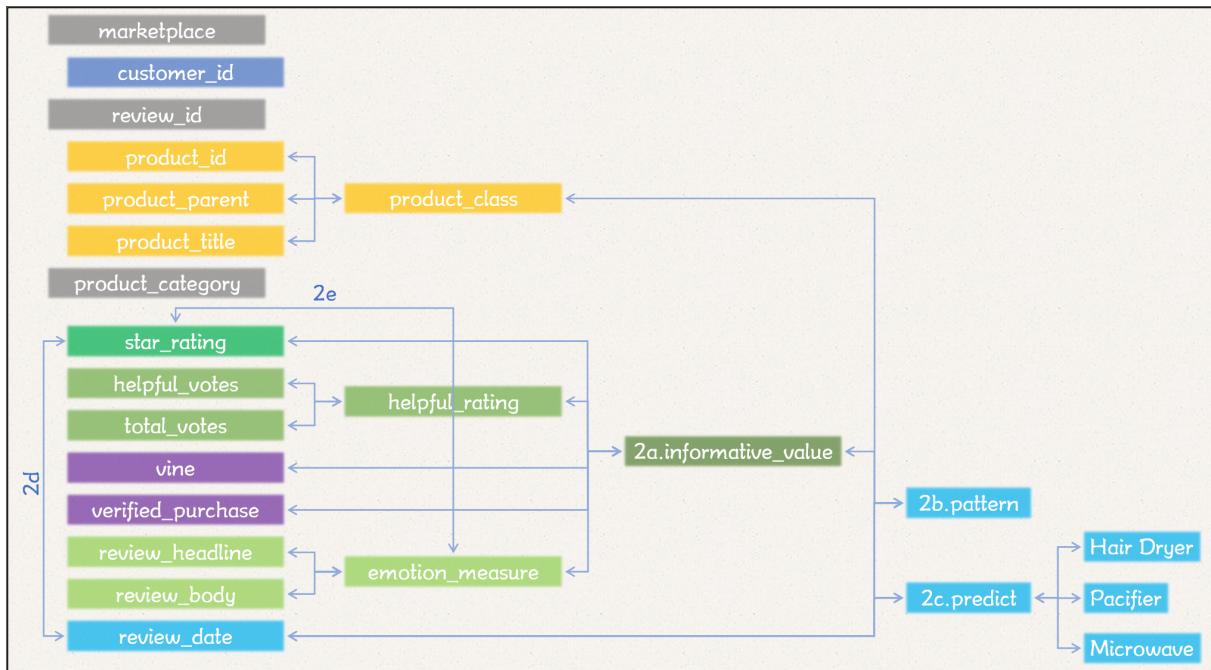


Figure 1: Workflow

- In the problem 1, we simply process the level-2 data, and then analyze the relationship between the review amount, star rating and date, the ratio of various star rating levels, review amount of customers, as well as the relationship between the review amount, review length, star rating and the helpfulness ratings.
- In problem 2(a), we propose three metrics: 1. weighted rating ratio, 2. weighted sentimental score, 3. preference vector. They depict the attributes of a review from the perspectives of star rating, emotion and what reviewers care about. Further, we can measure the value of reviews from these three dimensions in needs.
- In problem 2(b), we consider the reputation of a product in a fixed-length time window as a joint contribution of total review number, averaged rate, helpful votes, review readability and averaged sentimental scores of reviews over that period. And after calculation, we observe its changing values for each product over time.
- In problem 2(c), we continue to use the time-varying reputation in 2(b) as a comprehensive index of both ratings and reviews and LSTM is applied to predict their values in the next 1000 reviews. Since we assume the sales amount is proportional to the total review number, the reputation is associated with the sales and takes nearly every informative given features into accounts (sales amount included). Also because the training and testing losses of LSTM are too small to be valued, we are confident in the power of this index.
- In problem 2(d), we consider the ripple effect of reviews. After analyzing the trend of the monthly average number of different star ratings over time, we conclude that reviews with a score of 5 often lead to more reviews. Furthermore, we unexpectedly find that the reviews amount of star ratings 1 is significantly correlated with the length of reviews, and after **Granger Cause Test**, we find that short reviews tend to lag within two months after a large ratio of low star ratings.
- In problem 2(e), we extract the emotion words , score and rank the reviews based on the words. After comparing the confusion matrix of rating and review ranks, we find some mapping asymmetricity that some reviews with strongly affective words have mild ratings, and reviews with more mildly affective



words have extreme ratings. We explain this phenomena by showing that for those reviews, greater chances are that both positive words and negative affective words occur, or strongly affective words are replaced by words that described product properties.

## 2 Problem 1: Data Preprocessing and Mining

In this section, we preprocess the data set and analyze the relationship between review amount, star rating, helpfulness rating, customer and date, and present some findings.

### 2.1 Data Cleaning

#### 1. General Data Cleaning:

The given data set presents 15 properties for the three products: hair dryer, microwave and pacifier. In all of these properties, marketplace which is either "US" or "us", review id which is unique for all items and product category which is exclusive for these three products, are useless for our subsequent analysis and we simply delete the three properties. For vine and verified purchase, we replace all n and N with N, and all y and Y with Y. Then we clean all given texts by replacing some special Chinese characters with equivalent ASCII characters as much as possible, and eliminating the remaining special garbled codes and symbols, for further text-based emotion analysis and keyword extraction. In the end, we eliminate 214 unqualified items and leave 11424, 1606 and 18784 items for hair dryer, microwave and pacifier, respectively.

#### 2. Text Data Pre-processing[7]:

- (a) Tokenization: Split texts into sentences and sentences into words. Lowercase words and remove punctuation.
- (b) Words that have fewer than 3 characters are removed. All stopwords are removed.
- (c) Words are lemmatized: words in third person are changed to first person and verbs in past and future tenses are changed into present.
- (d) Words are stemmed: words are reduced to their root form.

### 2.2 Text Mining by LDA Topic Model[1]

Since the review texts are rich and complex, it is of great importance to dig out the topic words so as to clearly grasp what consumers really care about. And here we accomplish this task by applying the Latent Dirichlet Allocation (LDA) model, a powerful NLP method for topic word extraction.

By simply running the standard python codes of LDA, we get results that will be useful in the subsequent estimation of users' preference vectors. We show results of LDA for hair dryer in Fig. 2, and the rest results can be seen in the Appendix. A.



关注校苑数模公众号

获取更多资料

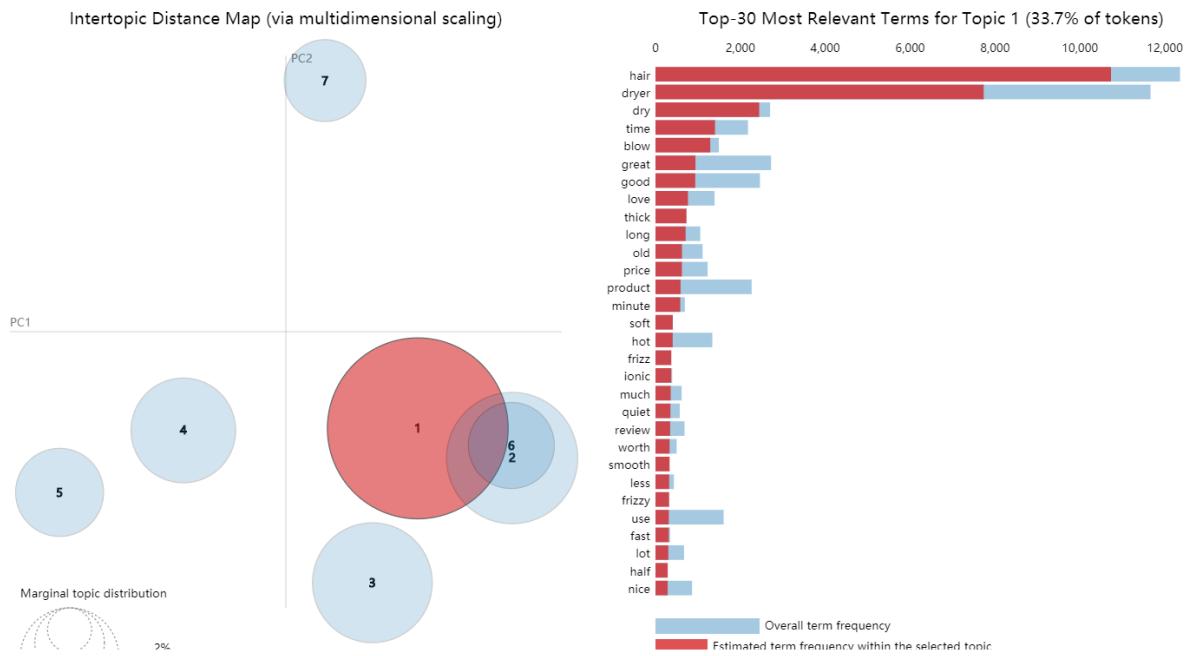


Figure 2: LDA topic model for hair dryer

### 2.3 Overall Data Characteristics

First, we analyze **the changes of some main features over time**[8]: We visualize the review amount of three products each month in Fig. 3(a) and observe that the reviews amount has increased exponentially over time. Besides that, the averaged rating score per month is further calculated and displayed, where reviews which occur less than 2 times in a month are not included in the statistics, as shown in Fig. 3(b):

1. The averaged review rating of all three products is generally between 3.5 and 4.5.
2. The rating scores tend to converge over time.
3. The reputation of microwave seems to be continuously increasing in recent years.

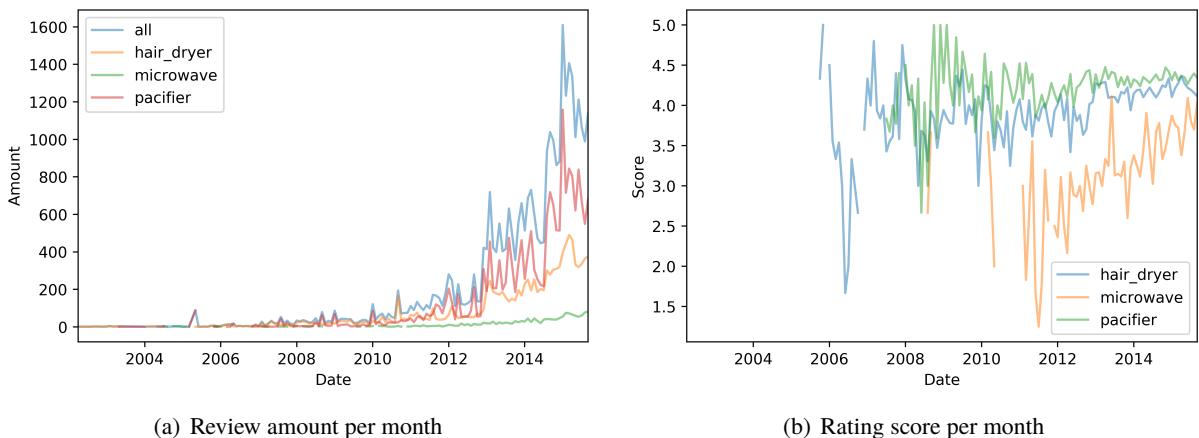


Figure 3: Review status over time

In order to understand the **distribution of review ratings**, we plot the proportion of five ratings for three products, shown in Fig. 4(a). Obviously, most of the ratings are 5 and the next most are 4, where microwave



is an exception for owning too many 1 star ratings. This may be partly because microwave ovens are not easy to be accepted in the early stage.

Next we count the number of reviews of each customer and take the logarithm of the Y-axis to avoid a large difference in the order of magnitude. As can be seen in Fig. 4(b), **most customers only review once, and only a few customers review more than 5 times**.

Fig. 4(c) shows that most of the review votes are useful, where we only consider reviews with more than 10 votes, which represent the review quality appropriately.

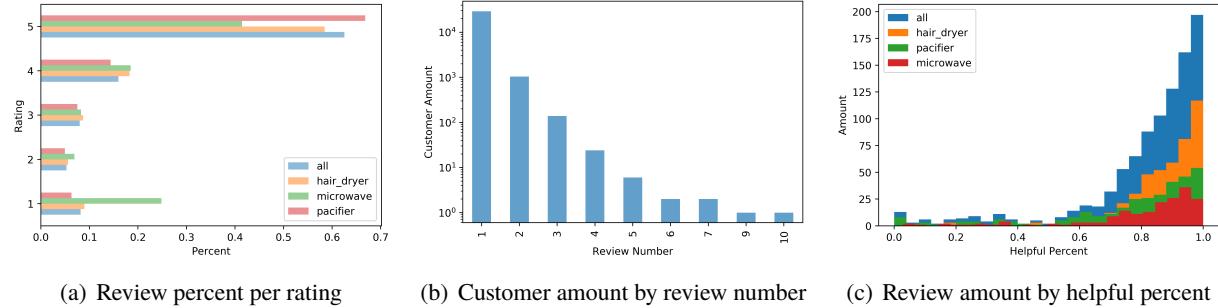


Figure 4: Review amount from different angles

To further understand **the characteristics of reviews with high helpfulness ratings**, we display the averaged review length and averaged star rating in Fig. 5. Besides, the time-varying averaged helpful votes of each review are shown in Fig. 6.

- In Fig. 5, reviews with higher helpfulness ratings tend to have more words and higher star ratings, possibly for the reason that **too short reviews are easily ignored, few customers will write long reviews with nonsense, or they will lose their interests in a product after seeing low ratings and quickly leave without giving a certificate for reviews they've just seen**, etc.
- In Fig. 6, helpful votes of a single review in the early stage far outnumber the those in the recent period. This may be because in the early stage there are too few comments for people to refer to, thus people naturally attach more importance of each comment in the cold start period. Moreover, reviews with more helpful votes may be displayed in the front of the comment zone, which enlarges their probabilities to get more votes. As stated in [9] and [10], the above phenomenon are called early birds bias and winner circle bias, respectively.

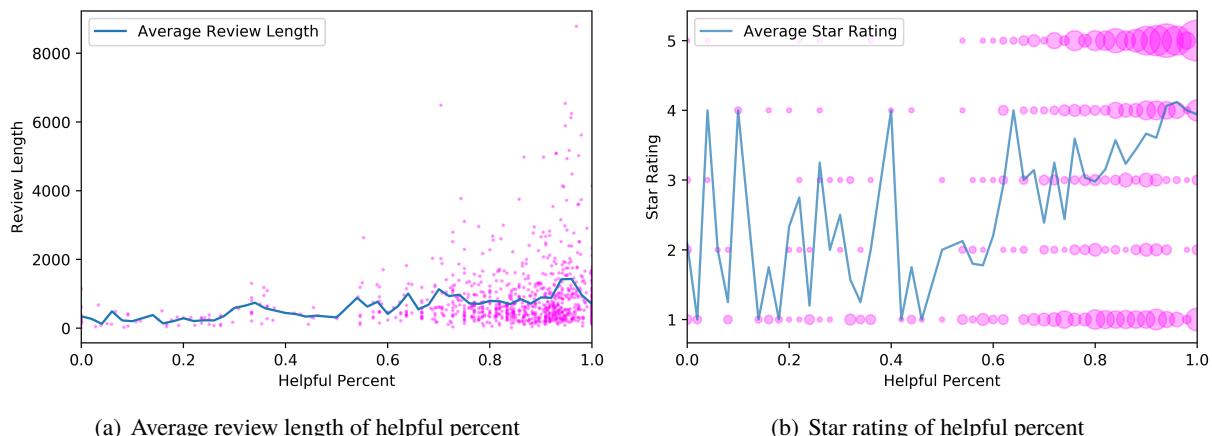


Figure 5: Review status by helpful percent



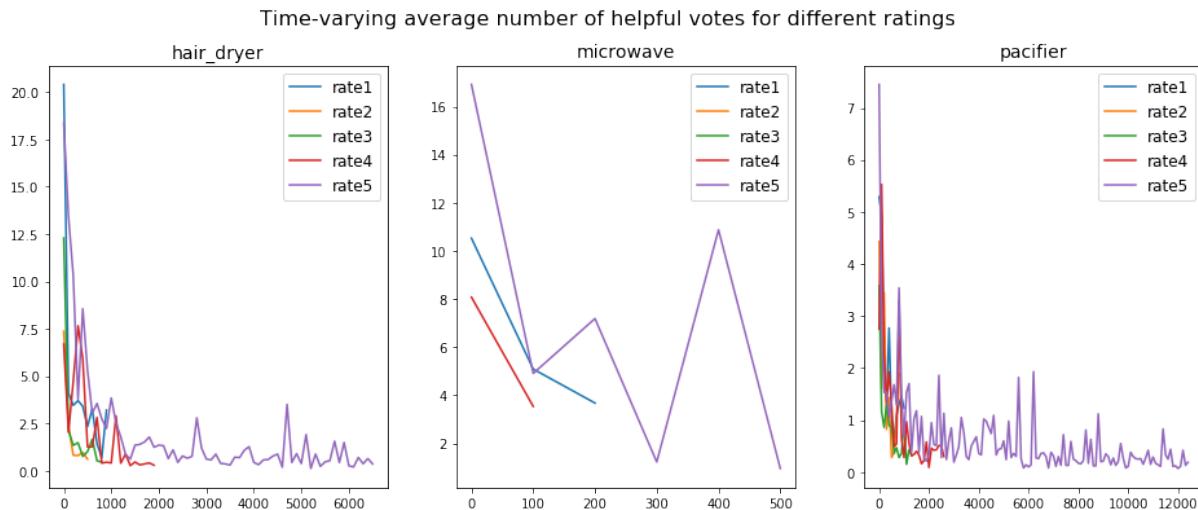


Figure 6: Time-varying average number of helpful votes for different ratings

In addition, we conduct a detailed analysis of various commodities under the three products to find out the relationship between the review amount and the averaged star rating and the averaged length of reviews on the same commodity. The results in hair dryer data set are shown in Fig. 7, and the rest are shown in the Appendix B. We can find it seems that longer reviews are often followed by improved ratings, and longer periods of high ratings will lead to an increase of reviews, which partly reflects the increase on sales.

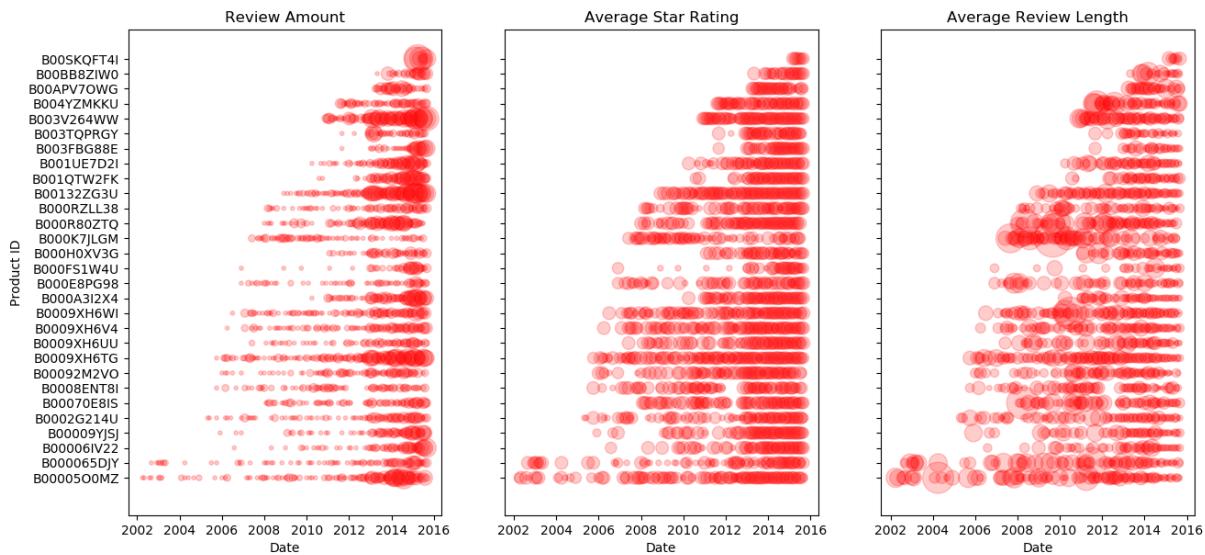


Figure 7: Product contrast of hair dryer over time

**Based on our previous analysis, Sunshine Company should note that:**

1. The number of reviews has increased rapidly in recent years, which reflects the importance of Sunshine's frequent interaction with reviewers.
2. The star rating will experience large fluctuations in the early stage of the product launch, during which time the brand image should be maintained cautiously.



3. Since reviews with high star ratings and more words tend to get more helpful votes, Sunshine had better closely watch out this type of reviews and discover customers' pain points based on them[11].

### 3 Problem 2 (a): Ratings and Reviews Based Data Measures

We design three types of data measures based on ratings and reviews that are informative for Sunshine Company to track, i.e. weighted rating ratio, averages and standard deviations of weighted sentimental scores for reviews and users preference vectors.

#### 3.1 Weighted Rating Ratio

1. For  $i \in [5]$ , denote the timestamps when a new rating  $i$  is created as  $A_i$ .
2. For the timestamp  $j \in A_i$ , denote the number of helpful votes at time  $j$  as  $HN_j$ .
3. Assign a weight of  $HN_j + 1$  to time  $j$ , and we obtain that the weighted ratio for rating  $i$  can be formulated as

$$wrr_i = \frac{\sum_{j \in A_i} (HN_j + 1)}{\sum_{i \in [5]} \sum_{j \in A_i} (HN_j + 1)}$$

$wrr_i$  represents **the occurring ratio of rating  $i$  weighted by helpfulness** among all five ratings.

According to the formula above and the given data, we obtain the  $wrr_i$  values as follows Fig. 8:

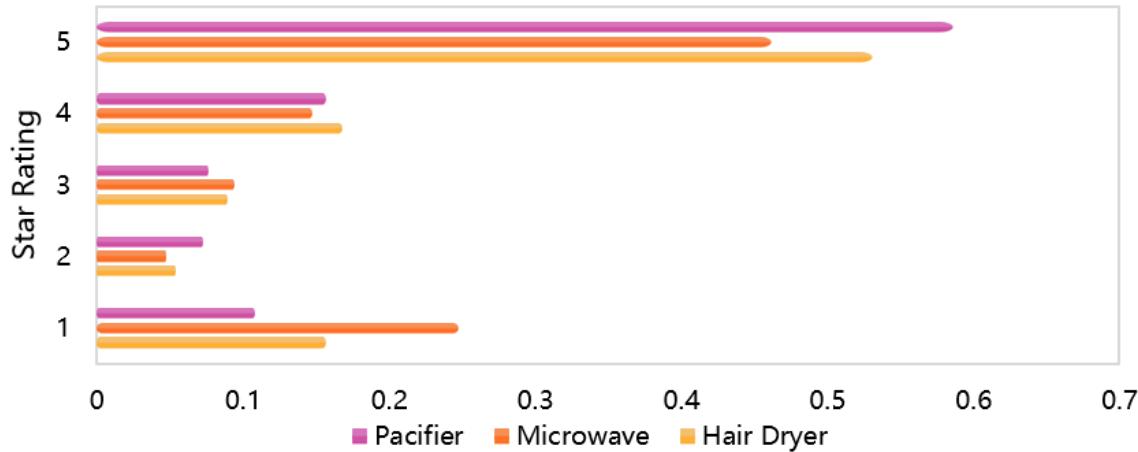


Figure 8: Weighted rating ratio of different rating for different product

#### 3.2 Avg and Std of Weighted Sentimental Scores for Reviews[2][3][4]

1. Calculate the normalized sentimental scores of notional words appearing in the review texts by logistic regression Algorithm 1:

**Note:** The reason for which we adopt logistic regression[4] is that normally the sentimental scores of negative words are negative, those of positive words are positive and those of neutral words should be zero. The sigmoid function in logistic regression happens to symmetrically map  $\mathbb{R}$  into  $[0, 1]$  with zero as the center of x axis. As a result,  $S$  can be naturally solved by this method.



**Algorithm 1:** Logistic Regression

**Input:** notional words set  $\mathcal{W} = \{w_1, \dots, w_n\}$ , review set  $\{re_t\}_{t \in T}$  and rating sets  $\{ra_t\}_{t \in T}$ .

**Output:** sentimental scores  $\mathcal{S} = \{s_1, \dots, s_n\}$

Map  $ra_t$  to  $y_t$  by  $\{1 : 0, 2 : 0.25, 3 : 0.5, 4 : 0.75, 5 : 1\}$ .

$$\mathcal{S} = \operatorname{argmin}_S \sum_{t \in [T]} (y_t - \frac{1}{1+e^{\sum_{s \in S} s 1_s \{re_t\}}})^2$$

where  $1_s \{re_t\} = 1$  if the notional word  $s$  occurred in the  $t$ -th review, otherwise 0.

2. For  $t \in [t]$ , denote the sentimental score of the  $t$ -th review as  $ss_t$ , and we obtain that

$$ss_t = \sum_{s \in S} s 1_s \{re_t\}$$

Then the weighted averaged reviews' sentimental score for rating  $i$  can be formulated

$$\text{avg}_i = \frac{\sum_{j \in A_i} (HN_j + 1) \times ss_j}{\sum_{i \in [5]} \sum_{j \in A_i} (HN_j + 1)}$$

and the weighted standard deviation of reviews' sentimental score for rating  $i$  can be formulated

$$\text{std}_i = \sqrt{\frac{\sum_{j \in A_i} (HN_j + 1) \times (ss_j - \text{avg}_i)^2}{\sum_{i \in [5]} \sum_{j \in A_i} (HN_j + 1)}}$$

According to the formula above and the given data, we obtain the  $\text{avg}_i, \text{avg}_i$  values as follows Fig. 9:

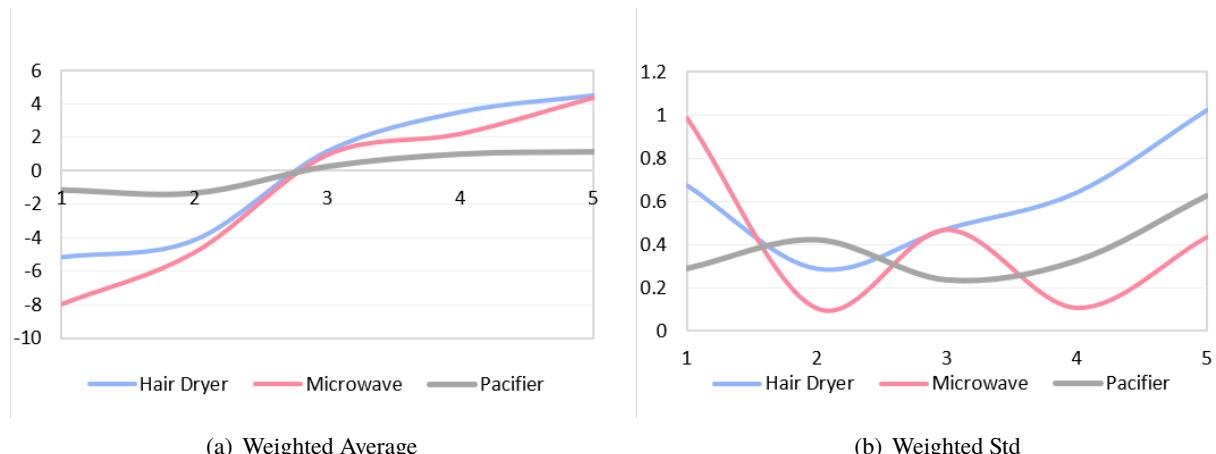


Figure 9: Weighted sentimental score index of reviews for different ratings

### 3.3 Users' Preference Vector

Since requirement measurements are fundamental to product positioning and strategic marketing development, we attach great importance to computing the preference vectors of consumers by analyzing the review data.

Owing to the topic analysis by the Latent Dirichlet Allocation (LDA) model and detailed observations of product reviews, we identify eight main product attributes for hair dryer, microwave and pacifier respectively. Lexicons of 46, 51 and 45 words are constructed to identify the attribute terms in the product reviews. The classified attributes and relative number of terms are described in Table 1.



Table 1: Classified attributes and relative number of terms

Hair dryer		Microwave		Pacifier	
Attribute	Num	Attribute	Num	Attribute	Num
Lifespan	5	Lifespan	5	Lifespan	5
Weight	3	Weight	3	Size and shape	9
After-sales service	9	After-sales service	9	Packing	8
Heat	3	Heat	3	Softness	3
Materials&Structure	19	Materials&Structure	24	Materials&Structure	13
Price	3	Cost	3	Cost	3
Safety	4	Safety	4	Safety	4

### Next is to estimate users' preference ratio for each attribute.

Taking the hair dryer product as an example, we denote users' preference degree and preference ratio for attribute  $i$  as  $PD_i$  and  $PR_i$ . Then it holds  $PR_i = \frac{PD_i}{\sum_{j \in [8]} PD_j}$ .

To calculate  $PD_i$ , we must clarify that

1. The older the review is, the less importance it is in  $PD_i$ .
2. Reviews with neutral ratings are attached little importance in  $PD_i$ .
3. The more helpful votes the review has, the more weight it takes up in  $PD_i$ .

Based on the considerations above, we define  $PD_i = \sum_{t \in T} (|ra_t - 3| + 0.5) \times 1_{it} \times (HN_t + 1)$ , hence

$$PR_i = \frac{\sum_{t \in T} \gamma^{1-t} (|ra_t - 3| + 0.5) \times 1_{it} \times (HN_t + 1)}{\sum_{i \in [8]} \sum_{t \in T} \gamma^{1-t} (|ra_t - 3| + 0.5) \times 1_{it} \times (HN_t + 1)}$$

where  $\gamma$  is a discounted factor of time, and we set it to be 0.999. According to the formula above and the given data, different users' preference vectors of three products are displayed in Table 2:

Table 2: Classified attributes and relative weights in reviews

Hair dryer		Microwave		Pacifier	
Attribute	PR	Attribute	PR	Attribute	PR
Lifespan	0.136	Lifespan	0.161	Lifespan	0.064
Weight	0.035	Weight	0.016	Size and shape	0.285
After-sales service	0.143	After-sales service	0.191	Packing	0.157
Heat	0.152	Heat	0.094	Softness	0.076
Materials&Structure	0.250	Materials&Structure	0.324	Materials&Structure	0.291
Price	0.049	Cost	0.048	Cost	0.020
Safety	0.235	Safety	0.166	Safety	0.107

Having computed the preference vectors, the Sunshine Company can readjust its strategies accordingly. Taking hair dryer as an example, the preference ratio of safety, materials quality, heat, lifespan and after sales service take up nearly 90 percentage, therefore the Sunshine Company should pay more attention to these attributes and assign proper funds to improve these aspects by their preference ratios respectively.



关注校苑数模公众号

获取更多资料

## 4 Problem 2 (b): Reputation Metric

### 4.1 Reputation Metric

The reputation of a product is related to many factors (e.g. the average star rating it receives, the authoritativeness of its reviews [12][13] and the sales volume). As we assume that the sales volume over a period of time is in proportion to the number of reviews the product receives. The reputation score a product earns over a period of time can be seen as the joint contribution of all reviews it receives, which indicating both the quality and sales volume.

Here we define a function  $\mathcal{R}(r_i)$ , which maps an instance of review  $r_i$  ( $r_i$  contains any information we might use from the review  $i$ ) into a scalar indicating the contribution of the review to the reputation for the product. Taking plenty of factors into account, the contribution score  $\mathcal{R}$  can be formulated as:

$$\mathcal{R}(r_i) = \frac{1}{t_i + 1} (wrr_i \times \alpha_{vine} \times \alpha_{verified} \times rating_i) \quad (1)$$

The scaling parameter  $t$  is the time interval (unit: day) between  $r_i$  and the last review. The larger the time interval, the less reviews the product receives, indicating less sales volume during that period of time.  $wrr_i$  is the weighted rate radio calculated in problem 2(a), which contains the information of helpful votes. The ratio  $\alpha_{vine}$  and  $\alpha_{verified}$  is formulated as:

$$\begin{cases} \alpha_{vine} = 0.2 * vine + 0.8 \\ \alpha_{verified} = 0.2 * verified + 0.8 \end{cases} \quad (2)$$

where we set the parameter  $vine$  and  $verified$  to 1 if the review has "vine" or "verified purchase" tags. Otherwise we set them to 0.

Next, a sliding window with size  $m$  is applied to the review table. A successive sequence of reviews, representing a period of time, constructs the average reputation score the product earns during that time:

$$rep_i^p = \frac{1}{m} \sum_{j=0}^{m-1} \mathcal{R}(r_{i+j}) \quad (3)$$

### 4.2 Analysis on Reputation Variation Patterns

By applying the sliding window strategy above to formulate the average reputation score, we can get a pattern of reputation variation with smooth changes because the difference between adjacent reputation scores is scaled by the window size  $m$ :

$$rep_{i+1}^p - rep_i^p = \frac{1}{m} (\mathcal{R}(i+m) - \mathcal{R}(i)) \quad (4)$$

In this section, we set the window size  $m$  equal to 200 and 500. As shown in Fig. 10 and Fig. 11, the larger the window size  $m$  is, more violently the reputation curves fluctuate.



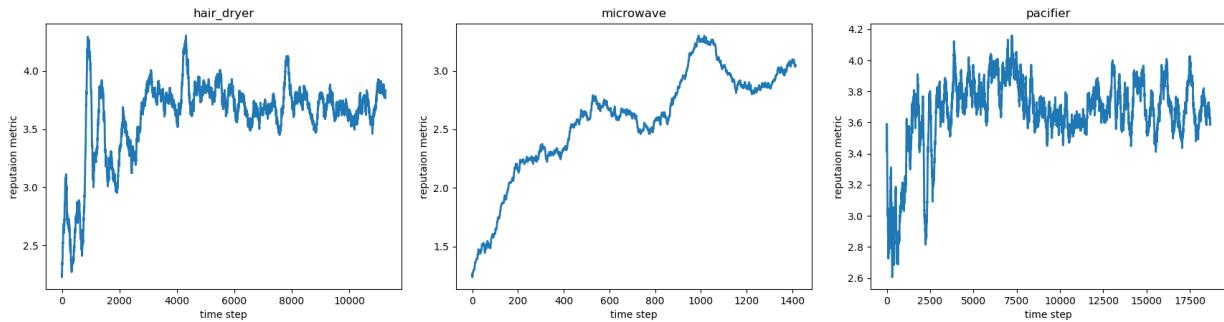


Figure 10: Reputation score of three products over time with window size  $m = 200$

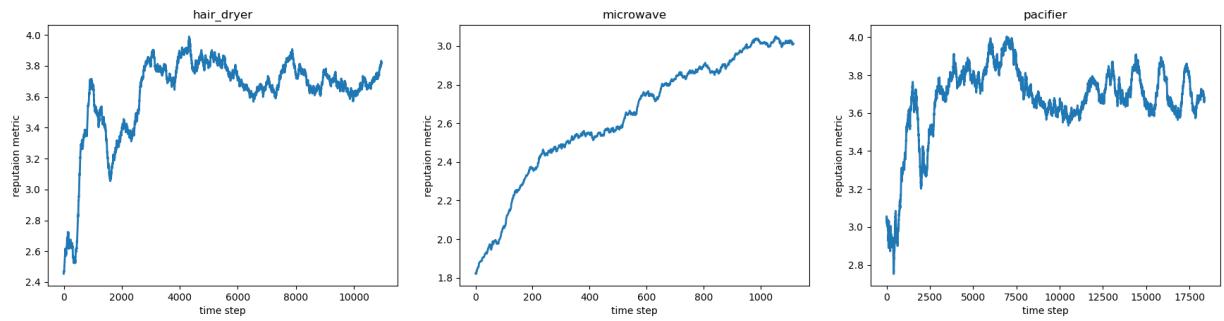


Figure 11: Reputation score of three products over time with window size  $m = 500$

Although different window size  $m$  results in different volatility of the reputation curve, they all generally have the same variation patterns for each product. For hair dryer and pacifier, their reputation scores increase in the early time and tend to be stable later. For microwave, its reputation score keeps growing the whole time. These trends also roughly correspond to the pattern of star ratings shown in Fig. 3, which is a lateral support for the correctness of our reputation metric.

## 5 Problem 2 (c): Nested Two-Layer LSTM

### 5.1 The Structure of the Nested Two-layer LSTM Model

For problem 2(c), we continue using the reputation score  $rep_i^P$  as the metric representing the success of a product. A potentially successful or failing product means that the reputation score keeps increasing or decreasing in the future and reach a threshold value. To predict the reputation score, we use Long and Short-Term Memory (LSTM)[14][15] to model the variation pattern of the review sequence.

LSTM model is a variant of recurrent neural network (RNN). It is well-suited to classifying, processing and making predictions based on time series data. A common LSTM unit is composed of a cell, an input gate, an output gate and a forget gate. The cell remembers values over arbitrary time intervals and the three gates regulate the flow of information into and out of the cell. The structure of our nested two-layer LSTM model is shown in Fig. 12. The model can be divided into three separate modules: input module, LSTM module and output module.



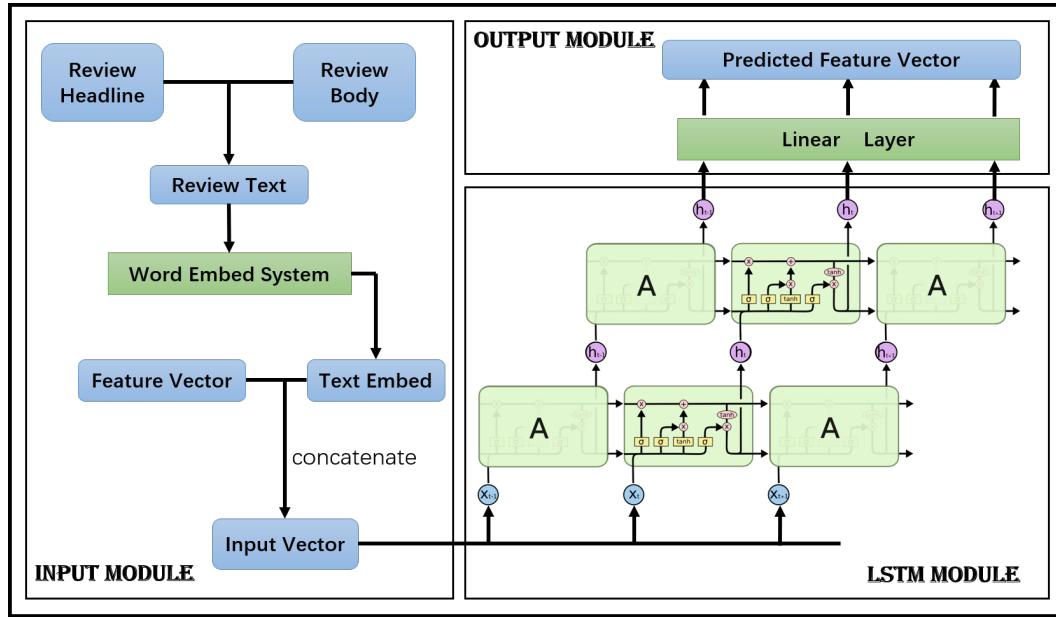


Figure 12: The structure of nested two-layer LSTM model

### 1. Input Module

For each review instance at a time stamp, we first feed the review text (combining both review headline and review body) into the embedding system to generate the sentence embedding  $v_s$  which is represented as a 200-dimension vector. Then we generate the input vector  $v_{input}$  for the 2-layer LSTM model with 200-dimension hidden size by concatenating sentence embedding and other useful features, which is represented as a 6-dimension feature vector  $v_f$ . More information about the feature vector  $v_f$  is shown in Table 3.

Table 3: Feature Definition of Vector  $v_f$ 

Dimension	Definition	Domain
$x_1$	Star rating	{1, 2, 3, 4, 5}
$x_2$	Helpful rating	[0, 1]
$x_3$	This review is vine voice or not	{0, 1}
$x_4$	This review is verified purchase or not	{0, 1}
$x_5$	The readability of this review	[0, 1]
$x_6$	Time interval (Unit: day) between this review and the last one	$\mathbb{N}$

All the features above are already defined in the previous sections. Therefore, the input vector for the LSTM can be formulated as:

$$\begin{cases} v_s = \text{WordEmbed}(review_{headline} + review_{body}) \in \mathbb{R}^{200} \\ v_f = [x_1, x_2, x_3, x_4, x_5, x_6] \in \mathbb{R}^6 \\ x = v_{input} = [v_s, v_f] \in \mathbb{R}^{206} \end{cases} \quad (5)$$

### 2. LSTM Module

The workflow of one-layer LSTM unit is shown below:



$$\begin{cases} i = \sigma(W_{ii}x + b_{ii} + W_{hi}h + b_{hi}) \\ f = \sigma(W_{if}x + b_{if} + W_{hf}h + b_{hf}) \\ g = \tanh(W_{ig}x + b_{ig} + W_{ho}h + b_{ho}) \\ o = \sigma(W_{io}x + b_{io} + W_{ho}h + b_{ho}) \\ c' = f * c + i * g \\ h' = o * \tanh(c') \end{cases} \quad (6)$$

The notations for equations above is shown in Table 4:

Table 4: Notations for one-layer LSTM

Symbol	Definition
$x \in \mathbb{R}^{206}$	The concatenated input vector for the LSTM
$h \in \mathbb{R}^{200}$	Hidden state, containing encoded information for the sequence flow
$c \in \mathbb{R}^{200}$	Cell state, tracking dependencies between the elements in the input sequence
$i \in \mathbb{R}^{200}$	Input gate, controlling the extent to which a new value flows into the cell
$f \in \mathbb{R}^{200}$	Forget gate, controlling the extent to which a value remains in the cell
$g \in \mathbb{R}^{200}$	Gathered input value from input $x$ and current hidden state
$o \in \mathbb{R}^{200}$	Output gate, controlling the extent to which the cell is used to compute outputs
$W$	The weight matrix for transitions
$b$	The bias for transitions
$\sigma$	The sigmoid function
$\tanh$	The hyperbolic tangent function

As for our nested two-layer LSTM, a LSTM unit is overlaid with another LSTM unit. We feed the hidden state of the bottom LSTM into the upper LSTM again and use the hidden state from the upper state as the input for the output module.

### 3. Output Module

In the output module, we simply feed the hidden state of the upper LSTM into a linear layer to make predictions on the feature vector, which has the same structure as  $v_{input}$  in the input module.

$$v_{pred} = Wv_h + b \quad (7)$$

## 5.2 Analysis on Potential Success of Products

Before we show the analysis results, there are several training details to be stated:

### 1. Loss function

During the training, we use mean square error (MSE) as the loss function:

$$MSE(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (8)$$

where  $y$  is the true target vector and  $\hat{y}$  is the predicted feature vector and  $N$  is the number of dimension.

### 2. Normalization

Due to different domains of each dimension in the input vector  $v_{input}$ , a normalization technique is needed. Before we feed the input vector into the LSTM module, we apply the following per-feature normalization operation:



$$x_i = \frac{x_i - \mu_i}{\sigma_i} \quad (9)$$

where  $\mu_i$  and  $\sigma_i$  are the mean and standard deviation of the  $i^{th}$  feature.

### 3. Training data

To make our model robust, we generate training data by sampling from the original review sequence with different lengths and different starting points. Our model is trained on these generated review sequences for 30 epochs (20 batches each epoch) with learning rate  $\alpha = 0.8$  and LBFGS optimizer.

After finishing the training period, we feed the original whole review sequence into the well-trained LSTM model to predict the next  $n$  review items ( $n = 0.4 \times \text{original}$ ). With the predicted review information, we can compute the future reputation score for each product. Therefore, a product with a reputation score that keeps increasing or stable in the future can be identified as a potentially successful product. Otherwise, it can be a potentially failing product.

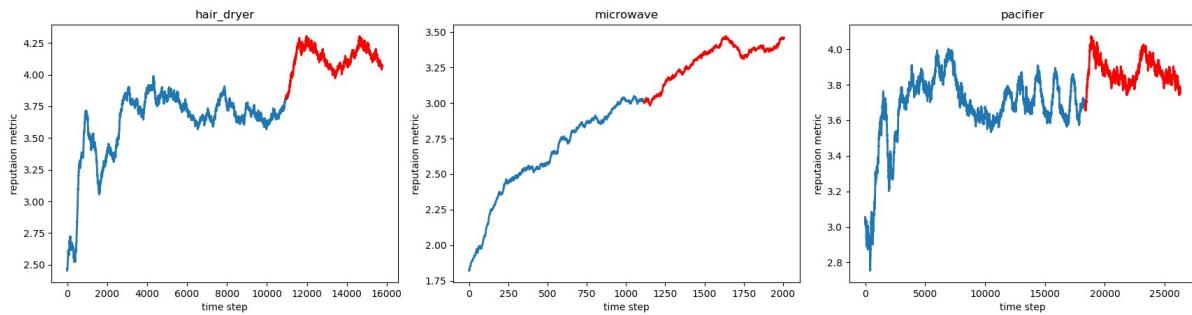


Figure 13: The future reputation score for each product with window size  $m = 500$ . The blue line represents the historical reputation scores. The red line represents our predicted values.

As shown in Fig. 13, hair dryer, which keeps a stable flow of reviews, will meet a rising trend in the future and then get back to the stable state again. Product microwave will keep its growth trend until it meets the upper bound. Therefore, these two products are potentially successful due to the increasing reputation scores in the future. As for product pacifier, it just keeps the stable state in the future.

## 6 Problem 2 (d): Causal Effectiveness Between Reviews

### 6.1 Ripple Effects of Extreme Ratings

To figure out whether specific star ratings incite more reviews, we calculate the monthly review numbers for each rating, seen in Fig. 14.



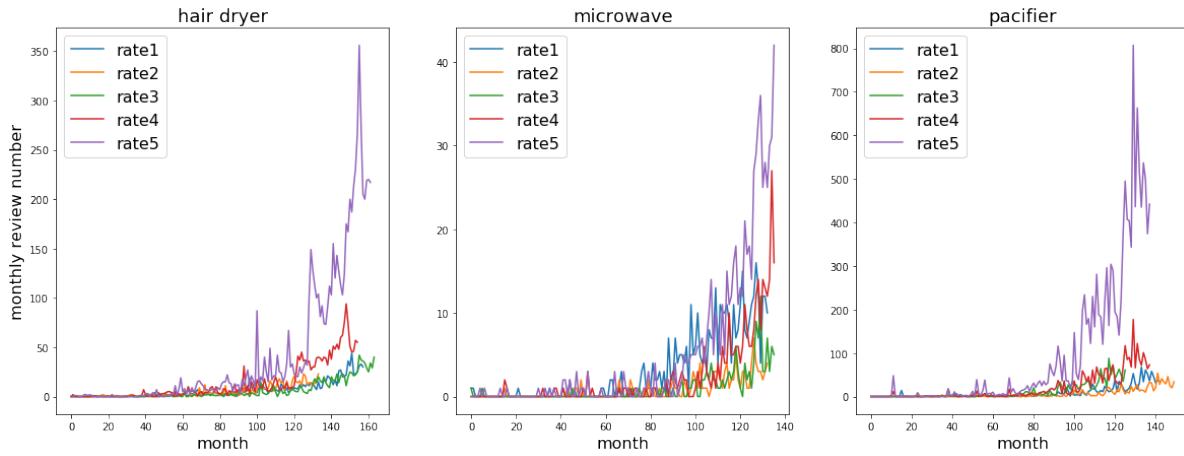


Figure 14: Monthly review number per rating

As can be seen in Fig. 14, at the beginning, reviews of all ratings are few. Even if the number of reviews with rate 5 usher in a few small peaks, the numbers of reviews with other ratings seem not to be affected. The real inflection point is when review number with rate 5 surges, the numbers of reviews corresponding to other ratings also increase. After that, the total number of comments rises almost exponentially. Besides, the peaks of other-rated reviews are almost always behind those of rate 5. In summary, we speculate that **a certain amount of rate 5 may incite more reviews**.

## 6.2 Causal Inference for Ratio of Low Rating and Review Length

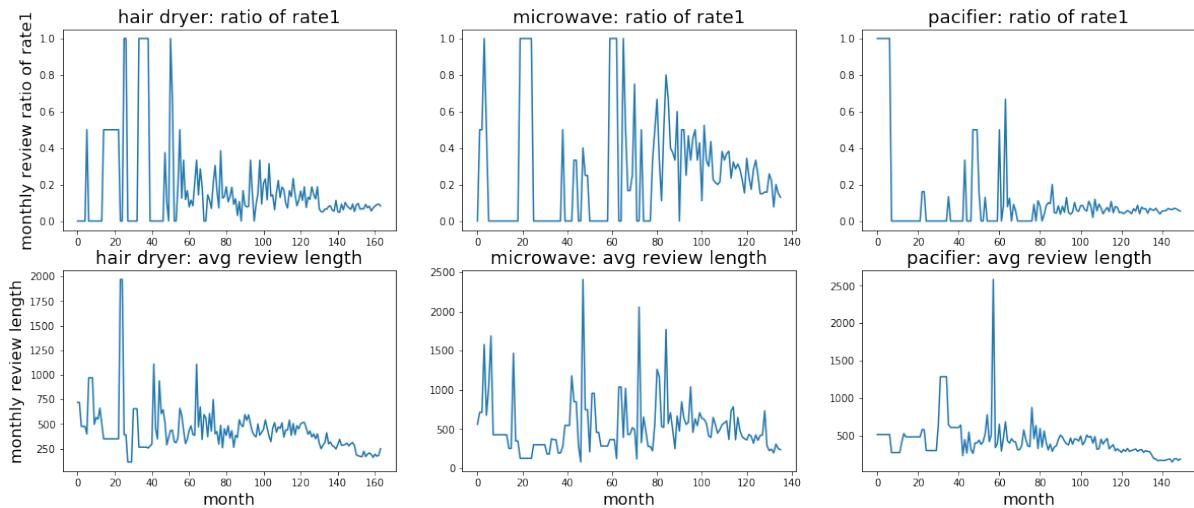


Figure 15: Monthly ratio of reviews with rate 1 and monthly review length

During the process of data mining, we unexpectedly find that the ratio trend of reviews with rate 1 is surprisingly similar to the trend of all reviews' average length. To further dig out which is the cause and which is the lag variable, we apply the Granger Cause Test for causal inference.

### Introduction and running steps of Granger Cause Test:

1. Granger Cause Test is a well-known method in Statistics to determine if a particular variable comes after another in the time series and what the lags are. It is a bottom up procedure which assumes the



data-generating processes in any time series are independent variables; then the data sets are analyzed to see if they are correlated.

## 2. Running the Test[16]:

The null hypothesis for the test is that lagged x-values do not explain the variation in y. Therefore, we run the **F-Test** to examine this hypothesis.

- State the null hypothesis and alternate hypothesis. For example,  $y(t)$  does not Granger-cause  $x(t)$ .
- Choose the lags. This mostly depends on how much data you have available. One way to choose lags i and j is to run a model order test (i.e. use a model order selection method). It might be easier just to pick several values and run the Granger test several times to see if the results are the same for different lag levels. The results should not be sensitive to lags.
- Find the f-value. Two equations can be used to find if  $\beta_j = 0$  for all lags j:

$$y(t) = \sum_{i=1}^T \alpha_i y(t-i) + c_1 + v_1(t)$$

$$y(t) = \sum_{i=1}^T \alpha_i y(t-i) + \sum_{j=1}^T \beta_j x(t-j) + c_2 + v_2(t)$$

Two equations for Granger causality: Restricted (top) and unrestricted (bottom).

Similarly, these equations test to see if  $y(t)$  Granger-causes  $x(t)$ :

$$x(t) = \sum_{i=1}^T \alpha_i x(t-i) + c_1 + u_1(t)$$

$$x(t) = \sum_{i=1}^T \alpha_i x(t-i) + \sum_{j=1}^T \beta_j y(t-j) + c_2 + u_2(t)$$

- Calculate the f-statistic using the following equation:

$$F = \frac{((SS'_E - SS_E)/m)}{MS_E} \sim F(m, df_E)$$

- If the p-value for this test is less than the designed value of  $\alpha$ , then we reject the null hypothesis and conclude that x causes y (at least in the Granger causality sense). Otherwise, change the lags and reperform F-test.

Table 5: Granger Causality: Number of lags (no zero) 2

ssr based F test	$F=5.6862, p=0.0041, df\_denom=157, df\_num=2$
ssr based ch12 test	$ch12=11.7346, p=0.0028, df=2$
likelihood ratio test	$ch12=11.3290, p=0.0035, df=2$

From the results in Table 5, we can claim with at least 95 percent confidence level that the averaged review length comes after the the ratio of reviews with rate 1 and the lags are 2 data points, that is, combining with Fig. 15, **short reviews often come after large ratio of extremely low ratings with lags within 2 month.**



关注校苑数模公众号

获取更多资料

## 7 Problem 2 (e)[5]: Correlation between Affective Words and Star Ratings

Nowadays, a new career named water army unconsciously bloom. People who engage in this occupation act as Internet ghostwriters paid to post online comments with particular content. In order to help achieve or obstruct the online shopkeepers' sales plan while not being monitored by anomaly detection machines, they may disturb the normal alignment of ratings and review texts. Moreover, **since different people have different rating standards, their mappings from reviews to ratings are various.** Therefore, chances are that: **for some people, specific quality descriptors of their text-based reviews such as enthusiastic and disappointed are strongly associated with rating levels, for others the relations are relatively looser to different extents.** As a result, analyzing the alignment of rates and review texts is of great importance.

### 7.1 Analysis for Alignment of Rate and Review Scored by Certain Affective Words

Recall that in section 3.2. we calculate the sentimental scores for each review text, which represent different people's attitude to products. In this section, we

1. update all sentimental scores by **just considering the combinations of affective words such as wonderful and terrific.**
2. After that, we further divide the reviews into 5 ranks by 4 thresholds of sentimental scores. The i-th thresholds, denoted as  $thre_i$ , are set to be  $\frac{avg_i + avg_{i+1}}{2}$ , where  $avg_i$  is the weighted averaged reviews sentimental score for rating i defined in section 3.2.

So far, we've classify all ratings and reviews into 5 groups respectively. Next, we count the number of items where the rating is rank i and the review is rank j for  $i, j \in [5]$ , generate the confusion matrix and compute the recall, precision and macro F1 value for each product[17], as shown in Fig. 16.

Assume rate 1,2 and review rank 1,2 belong to negative attitudes, rate 4,5 and review rank 4,5 belong to positive attitudes. From the figures above, it's easily seen that positive reviews are **seldom mapped into negative ratings and so are negative reviews.** Nevertheless, **crossover mappings between rate 1 and review rank 2 or rate 3 and review rank 4 and etc. seem to take up a certain percentage** that we cannot ignore.

Also, from the bottom right part of Fig. 16(d), **the recall, precision and macro F1 value of rate and review alignment for all three products are between 0.6 and 0.8, which are not very satisfying.** Among these results, **the mapping asymmetry for microwave is the most significant.**

### 7.2 Micro Observation of Asymmetry between Rate and Review[6]

From subsection 7.1, we find that there is some asymmetry of rate-review mapping when it comes to the sentimental analysis. Next, we aim to detailedly dig out the reason of this asymmetry, i.e. why reviews with strongly affective words has mild ratings, and why reviews with more mildly affective words has extreme ratings.

1. **Both positive words and negative affective words occur in the same reviews.**

By scanning the review texts, we find that some consumers detailedly illustrate their attitude towards products from high expectation or satisfaction to disappointment. As a result, we particularly analyze the reviews whose ranks deviate from their ratings, the results are shown as follows Table 6.

2. **Strongly affective words are replaced by words that describe properties.**



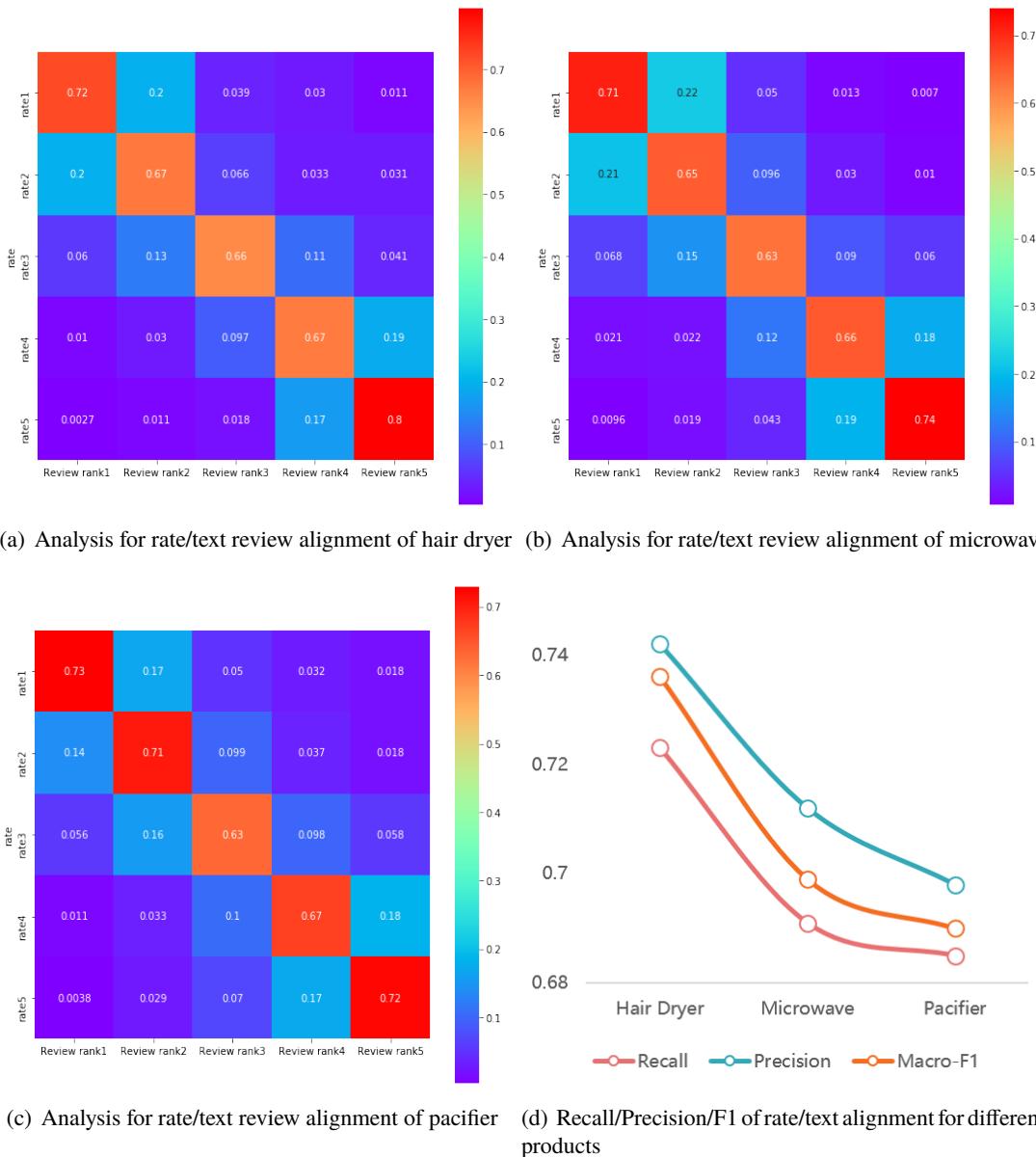


Figure 16: Analysis for alignment of rate and review scored by certain affective words

Table 6: Ratio for both occurrence of positive and negative affective words:  
reviews with and without rating deviation, respectively.

	Hair Dryer	Microwave	Pacifier
With rate/text deviation	0.12	0.078	0.098
Without rate/text deviation	0.191	0.163	0.152

People may naturally equal mildly affective words to neutral ratings, while through text data mining, we observe that there are a proportion of extreme ratings whose mildly affective words outnumber strongly affective words.

To dig out the reason of this phenomena, we select reviews with rate 1 and rate 5, calculate the averaged length and number of property words such as heat and magnetrons for extreme reviews, i.e. reviews with rate 1 and rate 5. After that, we compute the same statistics for extreme reviews with more neutral affective words than strongly affective words. The results are displayed in Fig. 17.



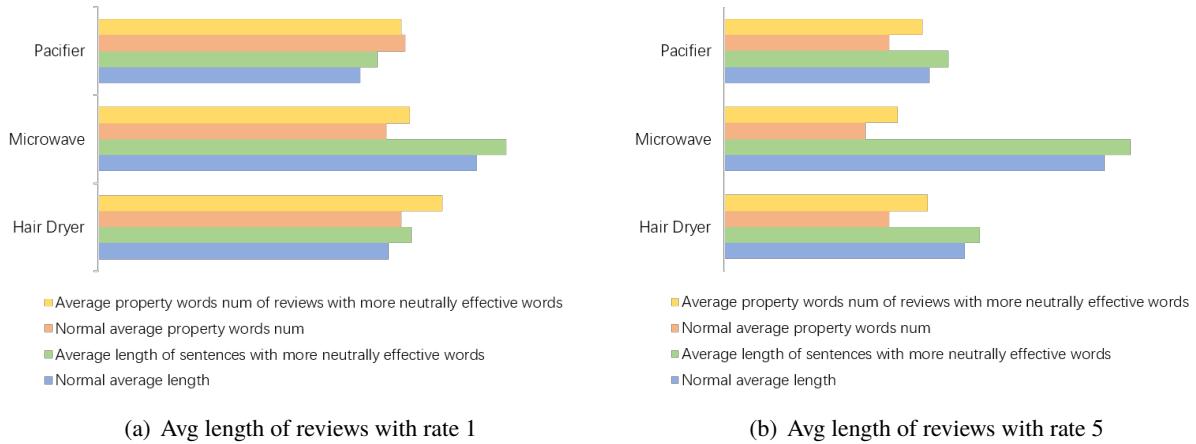


Figure 17: Comparisons between normal extreme reviews and extreme reviews with more mildly affective words.

From the figures above, it is easily seen that extreme reviews with more neutral words tend to consist of more property words and have longer lengths. With further observation of detailed examples, we find that this is largely because some reviewers are **objective enough to give factual descriptions** about products instead of piling up their strongly affective attitude towards them. **For these group of customers, the extents of their affective words in reviews are not closely related to ratings.**

## 8 Strengths and Weaknesses

### 8.1 Strengths

1. We thoroughly apply data cleaning by removing extremely short words in reviews and columns with similar meanings, tokenizing all texts and lemmatizing and stemming all words. Besides, we also apply LDA topic model to give a intuitive description of what reviews care about[18].
2. We vividly and diversely visualize our cross-dimensional analysis to explore the correlation between data in all directions. Some hidden features such as review readability and time gap of adjacent reviews are also taken into account[19].
3. We comprehensively apply multiple models to go deep into the patterns of reviews and ratings:
  - (a) We apply LDA topic model to  $7 \times 3$  dictionaries for different attributes of each product and apply weighted average with time decay to generate users' preference vectors.
  - (b) We apply sentiment analysis and logistic regression to estimate sentimental score of each review. Furthermore, we divide reviews into 5 rank , visualize the confusion matrix between rate and review rank and discover the mapping asymmetry of rate and reviews with regard to people's attitude.
  - (c) We use LSTM to compute the time-varying coefficients of joint contribution for the reputation index, which is powerful in globally dealing with moderate scale time series.
  - (d) Granger Cause Test is carried out to perform the causal inference for monthly ratios of low ratings and review lengths.



## 8.2 Weaknesses

1. We simply assume the sales column to be proportional to number of reviews, which may be too simplistic to be realistic.
2. Owing to lacking the price data and other necessary information, we fail to figure out the thorough reason for which all the three products spend a long time to go through the cold start period.
3. Due to the time limit, the terms in our dictionaries for different attributes of each products may be not enough to give an accurate profile of consumers' reviews.

## 9 Conclusion

Nowadays, online ratings and reviews for products play a more and more important role of influencing users' purchasing decisions. An in-depth analysis of this information is of great significance to companies' marketing strategies.

In this project, we make a detailed data cleaning of the given dataset. After that, we vividly and diversely visualize our cross-dimensional analysis to explore the correlation between data in all directions. We design three sub-index and a comprehensive index to indicator the products' reputation and prospects. Logistic regression, LSTM, LDA topic model, Granger Cause Test, sentiment analysis and confusion matrix method are applied to assist us in digging out the important and unexpected insights for each product, which allow us to make informed suggestions for the Sunshine Company.



关注校苑数模公众号

获取更多资料

## 10 The Letter to the Marketing Director of Sunshine Company

Dear director,

Considering today's increasingly fierce competition in e-commerce, accurately grasping customer needs and specifying appropriate marketing strategies are of vital importance to improve corporate profits and product visibility. As response to your company's requirement, we are here pretty glad to have the opportunity to introduce our research and suggestions to you, with the hope that it may give you some insights of the future strategies.

### 1. Pay attention to and properly guide the early reviews.

Throughout the whole records of consumers' ratings and reviews for hair dryer, microwave and pacifiers, we find the averaged number of helpful votes for a single review in the first 5 stage is above 3.2 times than that in recent period. Moreover, the descriptions in some reviews with rate 2 look even worse than that with rate 1, and so are descriptions in reviews with rate 4,5. Since helpful votes and the fitness of review description associate tightly with people's belief and the reputation propagation of products, we suggest Sunshine to properly guide and closely track the trend of early reviews, e.g. make guidelines of which rate corresponds to which extent of attitude, so as to smoothly go through the cold start stage.

### 2. Keep a good image of brand for all three products.

Based on our statistics, the number of reviews has increased rapidly in recent years, which reflects the great prospects of the online sales market. Whereas at the same time, it should be noted that the star rating and averaged review length will experience large fluctuations in the early stage of the product launch. This phenomena may give customers a bad or fuzzy brand impression. Moreover as the data shows, the early review volume, an indirect reflection of brand popularity, has been hovering at a lower level. As a result, if Sunshine can perform well, e.g. keep a good image of brand, in the early stage, great chances are that it will take a lead in such a highly competitive market, which may exert continuously positive effects on your subsequent development.

### 3. Apply our reputation index to monitor product dynamics.

Since the review and rating data is rich and complex, refining a comprehensive indicator that summarizes product reputation and sales is necessary for Sunshine to quickly adjust its promotion strategy. In our project, we design a reputation index which takes into account the star rating, time gap of adjacent reviews, review helpfulness rate and readability, information of vine and verification. The time-varying coefficients of joint contribution are calculated by LSTM, whose losses on testing dataset are small enough to be emphasized. Note that the time gap of neighboring reviews is a reflection of product popularity and we assume it to be proportional to the sales volume, hence our reputation metric can be a great indicator of not only people's attitude towards products but Sunshine's profits.

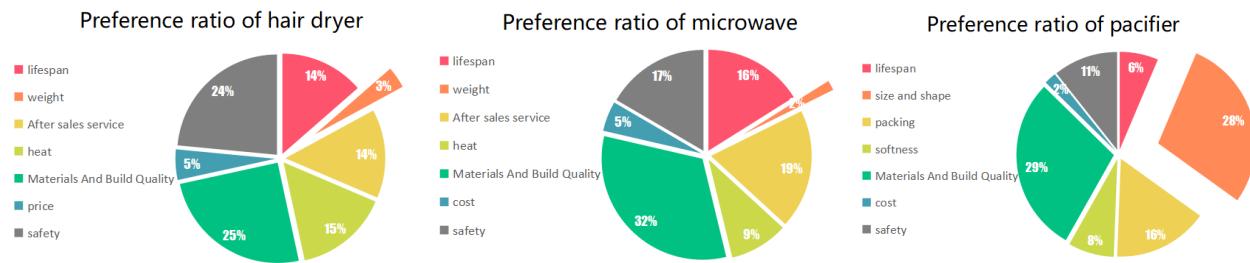
### 4. Attach different importance to products' properties according to people's preferences.

Through the text mining of LDA topic model, we extract 7 topics for each product, which represent the product properties that customers care about most. By an integrative consideration of reviews' keywords, ratings, helpful votes and the time decay of ratings, we estimate customers' preference vectors in 8 topics for each products, i.e.



关注校苑数模公众号

获取更多资料



As can be seen from the figures above, except for the tight control of products' materials quality, we suggest that your company ensure the security, heat and after sales service of hair dryer and microwave, as well as the size, shape and packing of pacifier.

We are really appreciated for this opportunity to assist you in building up an online marketing strategy, and we are convinced that our proposal can be utilized in improvement of your competence for the three products. Please feel free to contact us for further information on the project.

Sincerely yours

MCM 2020 Team



## References

- [1] Topic modeling using latent dirichlet allocation(lda) and gibbs sampling explained! [Online]. Available: <https://medium.com/analytics-vidhya/topic-modeling-using-lda-and-gibbs-sampling-explained-49d49b3d1045>
- [2] A. Bhatt, A. Patel, H. Chheda, and K. Gawande, “Amazon review classification and sentiment analysis,” *International Journal of Computer Science and Information Technologies*, vol. 6, no. 6, pp. 5107–5110, 2015.
- [3] T. U. Haque, N. N. Saber, and F. M. Shah, “Sentiment analysis on large scale amazon product reviews,” in *2018 IEEE International Conference on Innovative Research and Development (ICIRD)*. IEEE, 2018, pp. 1–6.
- [4] D. G. Kleinbaum, K. Dietz, M. Gail, M. Klein, and M. Klein, *Logistic regression*. Springer, 2002.
- [5] S. Dhanasobhon, P.-Y. Chen, M. Smith, and P.-y. Chen, “An analysis of the differential impact of reviews and reviewers at amazon. com,” *ICIS 2007 Proceedings*, p. 94, 2007.
- [6] P.-Y. Chen, S. Dhanasobhon, and M. D. Smith, “All reviews are not created equal: The disaggregate impact of reviews and reviewers at amazon. com,” *Com (May 2008)*, 2008.
- [7] Topic modeling and latent dirichlet allocation (lda) in python. [Online]. Available: <https://towardsdatascience.com/topic-modeling-and-latent-dirichlet-allocation-in-python-9bf156893c24>
- [8] J. Leino and K.-J. Räihä, “Case amazon: ratings and reviews as part of recommendations,” in *Proceedings of the 2007 ACM conference on Recommender systems*, 2007, pp. 137–140.
- [9] M. Li, L. Huang, C.-H. Tan, and K.-K. Wei, “Helpfulness of online product reviews as seen by consumers: Source and content features,” *International Journal of Electronic Commerce*, vol. 17, no. 4, pp. 101–136, 2013.
- [10] J. Liu, Y. Cao, C.-Y. Lin, Y. Huang, and M. Zhou, “Low-quality product review detection in opinion summarization,” in *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, 2007, pp. 334–342.
- [11] T. Wong, “Exploratory data analysis of amazon. com book reviews,” 2009.
- [12] S.-M. Kim, P. Pantel, T. Chklovski, and M. Pennacchiotti, “Automatically assessing review helpfulness,” in *Proceedings of the 2006 Conference on empirical methods in natural language processing*, 2006, pp. 423–430.
- [13] N. Korfiatis, E. GarcíA-Bariocanal, and S. SáNchez-Alonso, “Evaluating content quality and helpfulness of online product reviews: The interplay of review helpfulness vs. review content,” *Electronic Commerce Research and Applications*, vol. 11, no. 3, pp. 205–217, 2012.
- [14] F. A. Gers, J. Schmidhuber, and F. Cummins, “Learning to forget: Continual prediction with lstm,” 1999.
- [15] K. Greff, R. K. Srivastava, J. Koutník, B. R. Steunebrink, and J. Schmidhuber, “Lstm: A search space odyssey,” *IEEE transactions on neural networks and learning systems*, vol. 28, no. 10, pp. 2222–2232, 2016.
- [16] Granger causality: Definition, running the test. [Online]. Available: <https://www.statisticshowto.datasciencecentral.com/granger-causality/>



- [17] J. T. Townsend, "Theoretical analysis of an alphabetic confusion matrix," *Perception & Psychophysics*, vol. 9, no. 1, pp. 40–50, 1971.
- [18] S. Huang, J.-T. Sun, J. Wu, M. Wang, and Z. Chen, "Product review search," Sep. 4 2008, uS Patent App. 12/024,930.
- [19] J. C. De Albornoz, L. Plaza, P. Gervás, and A. Díaz, "A joint model of feature mining and sentiment analysis for product review rating," in *European conference on information retrieval*. Springer, 2011, pp. 55–66.



关注校苑数模公众号

获取更多资料

# Appendices

## Appendix A LDA Topic Model for Microwave and Pacifier

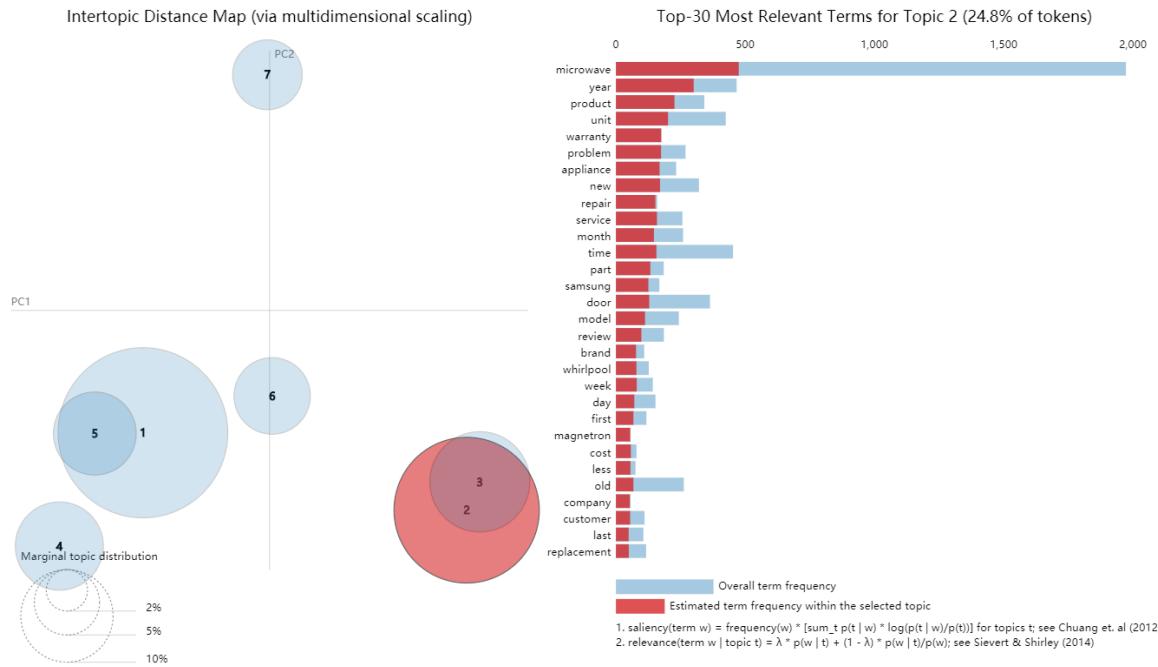


Figure 18: LDA topic model for microwave

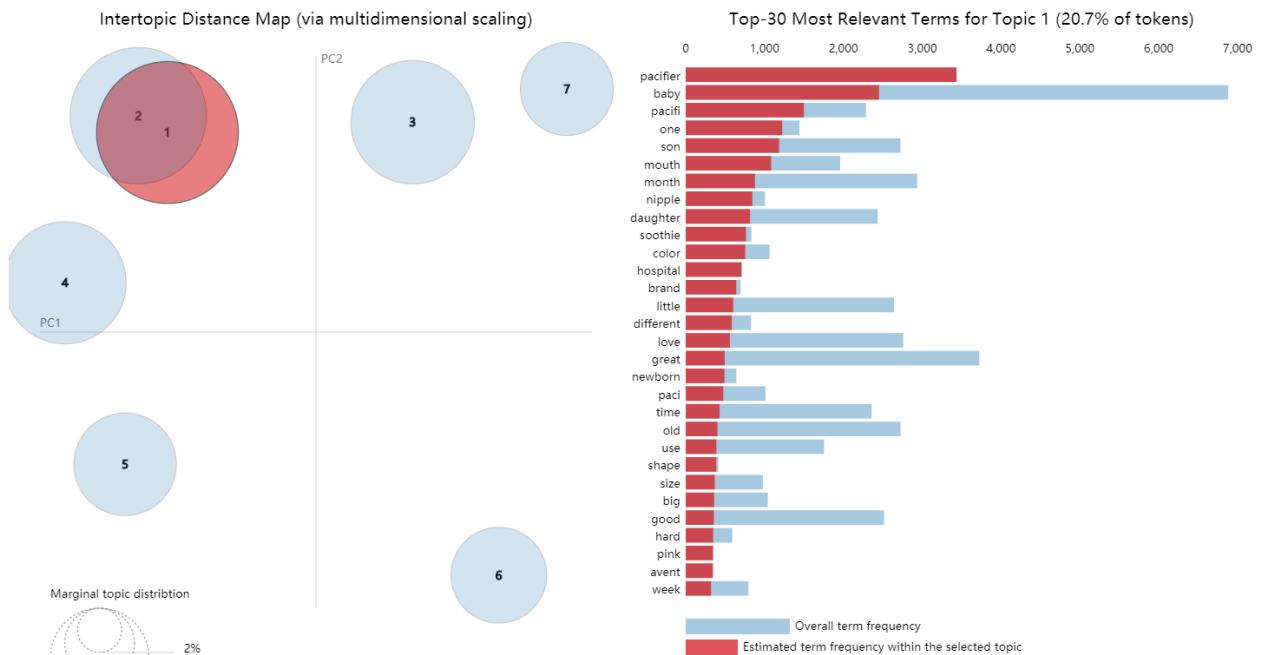


Figure 19: LDA topic model for pacifier



## Appendix B Product Contrast of Microwave and Pacifier Over Time

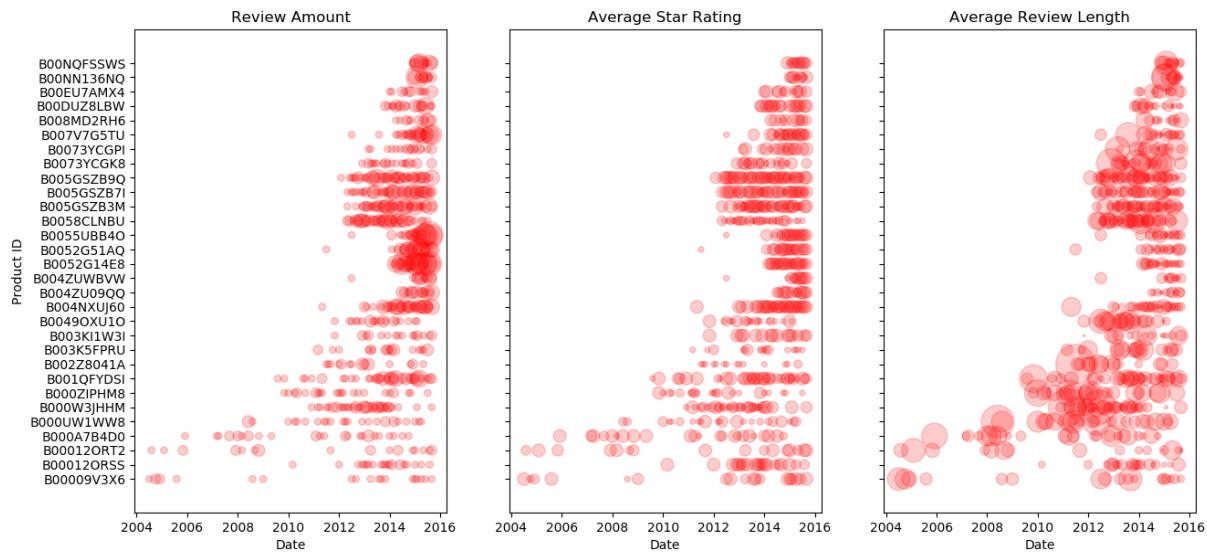


Figure 20: Product contrast of microwave over time

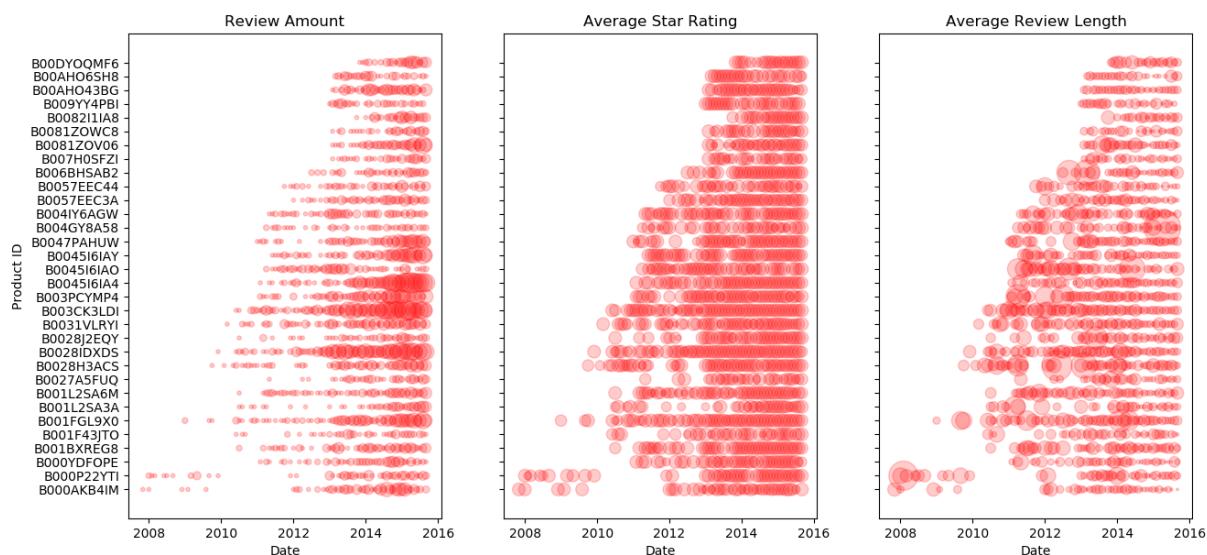


Figure 21: Product contrast of pacifier over time

## Appendix C Code

### C.1 Data Preprocess and Overall Analysis

```

1 ##########
2 ### 0. Clean Data
3 #####
4
5
6 import numpy as np
7 import pandas as pd

```



关注校苑数模公众号  
获取更多资料

```

8
9 hair_dryer = pd.read_csv('data/hair_dryer.tsv', sep = '\t')
10 microwave = pd.read_csv('data/microwave.tsv', sep = '\t')
11 pacifier = pd.read_csv('data/pacifier.tsv', sep = '\t')
12
13
14 def washer(dataset):
15     new_dataset = dataset.drop(['marketplace', 'product_category'], axis=1)
16     new_dataset.loc[new_dataset['vine'] == 'n', 'vine'] = 'N'
17     new_dataset.loc[new_dataset['vine'] == 'y', 'vine'] = 'Y'
18     new_dataset.loc[new_dataset['verified_purchase'] == 'n', 'verified_purchase'] = 'N'
19     new_dataset.loc[new_dataset['verified_purchase'] == 'y', 'verified_purchase'] = 'Y'
20
21     transfer = [('', '\''), ('', '\"'), ('', '''), ('', ''"), ('', ','), ('', '.'), ('',
22     '!'), ('', '...'), ('', '-'), ('', '-')]
23     accept = '[^a-zA-Z0-9_!?,.\\"+-=;()\\[\\]<>*#~&$^@%//|\\\\\\\\]'
24     for pair in transfer:
25         new_dataset['product_title'] = new_dataset['product_title'].str.replace(pair[0], pair[1])
26         new_dataset['review_headline'] = new_dataset['review_headline'].str.replace(
27             pair[0], pair[1])
28         new_dataset['review_body'] = new_dataset['review_body'].str.replace(pair[0],
29             pair[1])
30         trash = new_dataset[new_dataset['product_title'].str.contains(accept, na=True)
31                             + new_dataset['review_headline'].str.contains(accept, na=
32                             True)
33                             + new_dataset['review_body'].str.contains(accept, na=True)]
34
35     new_dataset = new_dataset[-(new_dataset['product_title'].str.contains(accept, na=
36     True)
37                             + new_dataset['review_headline'].str.contains(accept, na=
38     True)
39                             + new_dataset['review_body'].str.contains(accept, na=True))]
40
41
42 new_hair_dryer, trash1 = washer(hair_dryer)
43 new_microwave, trash2 = washer(microwave)
44 new_pacifier, trash3 = washer(pacifier)
45 trash = pd.concat([trash1, trash2, trash3])
46
47
48 del hair_dryer, microwave, pacifier, trash1, trash2, trash3
49
50
51 ##########
52 ### 1. Prepare: Import, Time Process, Review Process
53 #####
54
55
56 import matplotlib.pyplot as plt
57
58 def data_process(dataset):
59     dataset['review_date'] = pd.to_datetime(dataset['review_date'])
60     dataset['year'] = dataset['review_date'].dt.year
61     dataset['month'] = dataset['review_date'].dt.month
62     dataset['day'] = dataset['review_date'].dt.day

```



关注校苑数模公众号

获取更多资料

```
63     dataset['helpful_rate'] = dataset['helpful_votes'] / dataset['total_votes']
64     dataset['review_length'] = dataset['review_headline'].str.len() + dataset['
65     review_body'].str.len()
66
66 data_process(new_hair_dryer)
67 data_process(new_microwave)
68 data_process(new_pacifier)
69 all = pd.concat([new_hair_dryer, new_microwave, new_pacifier])
70
71 #####
72 ###### 2.1. Plot 1
73 #####
74 #####
75
76
77 import calendar
78
79 legend = ['all', 'hair_dryer', 'microwave', 'pacifier']
80
81 review_num_to_date = pd.concat([all.groupby(['year', 'month'])['star_rating'].count(),
82                                 new_hair_dryer.groupby(['year', 'month'])['star_rating'
83                                     ].count(),
84                                 new_microwave.groupby(['year', 'month'])['star_rating'
85                                     ].count(),
86                                 new_pacifier.groupby(['year', 'month'])['star_rating'
87                                     ].count()], axis=1)
87 review_num_to_date.columns = legend
88 review_num_to_date = review_num_to_date.reset_index()
89 review_num_to_date['date'] = review_num_to_date['year'] + review_num_to_date['month'] /
90                             12
91 review_num_to_date = review_num_to_date.drop(['year', 'month'], axis=1).set_index('date
92 ')
93 review_num_to_date.plot(legend=True, alpha=0.5)
94 # plt.title('Review Amount per Month')
95 plt.xlabel('Date')
96 plt.ylabel('Amount')
97 plt.savefig("figure/Review Amount per Month.png", dpi=500, bbox_inches = 'tight')
98 plt.show()
99
100
101 def func(dataset):
102     series = dataset.groupby(['year', 'month'])['star_rating'].agg(['count', 'mean'])
103     series[series['count']<3] = None
104     series = series.reset_index()
105     series['date'] = series['year'] + series['month'] / 12
106     series = series.drop(['year', 'month', 'count'], axis=1).set_index('date')
107     return series
108 ser_all, ser_hair, ser_micro, ser_paci = func(all), func(new_hair_dryer), func(
109     new_microwave), func(new_pacifier)
110 review_num_to_date = pd.concat([ser_hair, ser_micro, ser_paci], axis=1)
111 review_num_to_date.columns = legend[1:]
112 review_num_to_date.plot(legend=True, alpha=0.5)
113 # plt.title('Rating Score per Month')
114 plt.xlabel('Date')
115 plt.ylabel('Score')
116 plt.xticks(np.arange(0, 168, 30))
117 plt.savefig("figure/Rating Score per Month.png", dpi=500, bbox_inches = 'tight')
118 plt.show()
119
120
121 del legend, review_num_to_date, ser_all, ser_hair, ser_micro, ser_paci
```



关注校苑数模公众号

获取更多资料

```
119 #####  
120 ### 2.2. Plot 2  
121 #####  
122  
123  
124 review_percent_to_rating = pd.concat([all.groupby('star_rating')['star_rating'].count()  
    / all.shape[0],  
    new_hair_dryer.groupby('star_rating')['star_rating'].  
    count() / new_hair_dryer.shape[0],  
    new_microwave.groupby('star_rating')['star_rating'].  
    count() / new_microwave.shape[0],  
    new_pacifier.groupby('star_rating')['star_rating'].  
    count() / new_pacifier.shape[0]], axis=1)  
125 review_percent_to_rating.columns = ['all', 'hair_dryer', 'microwave', 'pacifier']  
126 review_percent_to_rating.plot.barh(legend=True, alpha = 0.5)  
127 # plt.title('Review Percent per Rating')  
128 plt.xlabel('Percent')  
129 plt.ylabel('Rating')  
130 plt.savefig("figure/Review Percent per Rating.png",dpi=500,bbox_inches = 'tight')  
131 plt.show()  
132  
133  
134 del review_percent_to_rating  
135  
136  
137#####  
138### 2.3. Plot 3  
139#####  
140#####  
141#####  
142#####  
143#####  
144#####  
145 for can in [all, new_hair_dryer, new_pacifier, new_microwave]:  
    can[can['total_votes'] > 10]['helpful_rate'].plot.hist(bins=25)  
146 # plt.title('Review Amount of Helpful Percent')  
147 plt.xlabel('Helpful Percent')  
148 plt.ylabel('Amount')  
149 plt.legend(['all', 'hair_dryer', 'pacifier', 'microwave'])  
150 plt.savefig("figure/Review Amount of Helpful Percent.png",dpi=500,bbox_inches = 'tight')  
151 )  
152 plt.show()  
153  
154  
155 del can  
156  
157  
158#####  
159### 2.4. Plot 4  
160#####  
161#####  
162#####  
163 all.groupby('customer_id')['customer_id'].agg({'num': 'count'}).reset_index().groupby('  
    num').count().plot(kind='bar', logy=True, legend=False, alpha=0.7)  
164 # plt.title('Customer Amount for Review Number')  
165 plt.xlabel('Review Number')  
166 plt.ylabel('Customer Amount')  
167 plt.savefig("figure/Customer Amount for Review Number.png",dpi=500,bbox_inches = 'tight')  
168 )  
169 plt.show()  
170  
171#####  
172### 2.5. Plot 5  
173#####  
174#####
```



关注校苑数模公众号

获取更多资料

```
175 def trans(k):
176     if k>=0:
177         return int(k*50)/50
178     return None
179
180 review_length = all[all['total_votes']>10][['helpful_rate', 'review_length']]
181 review_length['helpful_level'] = pd.Series([trans(k) for k in review_length['
    helpful_rate']], index=review_length.index)
182 review_length.groupby('helpful_level')['review_length'].agg({'average_length': 'mean'
    }).plot()
183 plt.scatter(review_length['helpful_rate'], review_length['review_length'], s=(2,), c='#
    ff00ff', alpha=0.3)
184 # plt.title('Average Review Length of Helpful Percent')
185 plt.xlabel('Helpful Percent')
186 plt.ylabel('Review Length')
187 plt.legend(['Average Review Length'])
188 plt.savefig("figure/Average Review Length of Helpful Percent.png", dpi=500, bbox_inches =
    'tight')
189 plt.show()
190
191 star_rating = all[all['total_votes']>10][['helpful_rate', 'star_rating']].groupby(
    'star_rating')['helpful_rate'].mean()
192
193
194 del review_length
195
196
197 #####
198 ##### 2.6. Plot 6
199 #####
200 #####
201
202 star_rating = all[all['total_votes']>10][['helpful_rate', 'star_rating']]
203 star_rating['helpful_level'] = pd.Series([trans(k) for k in star_rating['helpful_rate'
    ]], index=star_rating.index)
204 star_rating.groupby('helpful_level')['star_rating'].agg({'Average Star Rating': 'mean'
    }).plot(alpha=0.7)
205 helpful_star_frequency = star_rating.groupby(['helpful_level', 'star_rating']).count().r
    eset_index()
206 helpful_star_frequency.columns = ['helpful_level', 'star_rating', 'frequency']
207 plt.scatter(helpful_star_frequency['helpful_level'], helpful_star_frequency['
    star_rating'], helpful_star_frequency['frequency']*10,
208             c='#ff00ff', alpha=0.3)
209
210 # plt.title('Star Rating of Helpful Percent')
211 plt.xlabel('Helpful Percent')
212 plt.ylabel('Star Rating')
213 plt.legend(['Average Star Rating'])
214 plt.savefig("figure/Star Rating of Helpful Percent.png", dpi=500, bbox_inches = 'tight')
215 plt.show()
216
217
218 del star_rating, helpful_star_frequency
219
220
221 #####
222 ##### 2.7. Plot 7
223 #####
224
225
226 def product_date_process(dataset, threshold):
227     product = dataset.groupby('product_id')['star_rating'].agg({'count': 'count'})
```



关注校苑数模公众号

获取更多资料

```
229 product = product[product['count']>threshold].reset_index()['product_id']
230
231 product_date = dataset.groupby(['product_id', 'year', 'month'])['star_rating', 'review_length'].agg({'star_rating':['count', 'mean'], 'review_length': 'mean'})
232 product_date.columns = ['count', 'avg_star', 'avg_length']
233 product_date = product_date.reset_index()
234 product_date['date'] = product_date['year'] + product_date['month'] / 12
235 product_date = product_date.drop(['year', 'month'], axis=1).set_index(['product_id', 'date']).loc[product]
236
237 return product_date
238
239
240 product_date_h = product_date_process(new_hair_dryer, 400)
241 p, (ax1, ax2, ax3) = plt.subplots(nrows=3, ncols=1, sharex=True)
242 product_date_h['count'].unstack().T.sort_index().plot(ax=ax1, figsize=(5,7))
243 product_date_h['avg_star'].unstack().T.sort_index().plot(ax=ax2, legend=False)
244 product_date_h['avg_length'].unstack().T.sort_index().plot(ax=ax3, legend=False)
245 ax1.set_ylabel('Review Amount')
246 ax2.set_ylabel('Average Star Rating')
247 ax3.set_ylabel('Average Review Length')
248 plt.savefig("figure/Product Infomation for Date of Hair Dryer.png", dpi=500, bbox_inches = 'tight')
249 plt.show()
250
251
252 product_date_m = product_date_process(new_microwave, 78)
253 p, (ax1, ax2, ax3) = plt.subplots(nrows=3, ncols=1, sharex=True)
254 product_date_m['count'].unstack().T.sort_index().plot(ax=ax1, figsize=(5,7))
255 product_date_m['avg_star'].unstack().T.sort_index().plot(ax=ax2, legend=False)
256 product_date_m['avg_length'].unstack().T.sort_index().plot(ax=ax3, legend=False)
257 ax1.set_ylabel('Review Amount')
258 ax2.set_ylabel('Average Star Rating')
259 ax3.set_ylabel('Average Review Length')
260 plt.xticks(np.arange(2012, 2016))
261 plt.savefig("figure/Product Infomation for Date of Microwave.png", dpi=500, bbox_inches = 'tight')
262 plt.show()
263
264
265 product_date_p = product_date_process(new_pacifier, 250)
266 p, (ax1, ax2, ax3) = plt.subplots(nrows=3, ncols=1, sharex=True)
267 product_date_p['count'].unstack().T.sort_index().plot(ax=ax1, figsize=(5,7))
268 product_date_p['avg_star'].unstack().T.sort_index().plot(ax=ax2, legend=False)
269 product_date_p['avg_length'].unstack().T.sort_index().plot(ax=ax3, legend=False)
270 ax1.set_ylabel('Review Amount')
271 ax2.set_ylabel('Average Star Rating')
272 ax3.set_ylabel('Average Review Length')
273 plt.savefig("figure/Product Infomation for Date of Pacifier.png", dpi=500, bbox_inches = 'tight')
274 plt.show()
275
276
277 product_date_h = product_date_process(new_hair_dryer, 100).reset_index()
278 p, (ax1, ax2, ax3) = plt.subplots(nrows=1, ncols=3, sharey=True, figsize=(15,7))
279 ax1.scatter(product_date_h['date'], product_date_h['product_id'], product_date_h['count'] * 10, c='#ff0000', alpha=0.2)
280 ax2.scatter(product_date_h['date'], product_date_h['product_id'], product_date_h['avg_star'] * 20, c='#ff0000', alpha=0.2)
281 ax3.scatter(product_date_h['date'], product_date_h['product_id'], product_date_h['avg_length'] / 5, c='#ff0000', alpha=0.2)
282 ax1.set_title('Review Amount')
283 ax2.set_title('Average Star Rating')
284 ax3.set_title('Average Review Length')
```



关注校苑数模公众号

获取更多资料

```

284 ax1.set_xlabel('Date')
285 ax2.set_xlabel('Date')
286 ax3.set_xlabel('Date')
287 ax1.set_ylabel('Product ID')
288 plt.savefig("figure/Product Contrast for Date of Hair Dryer.png",dpi=100,bbox_inches =
289   'tight')
290 plt.show()
291
292 product_date_m = product_date_process(new_microwave, 15).reset_index()
293 p, (ax1, ax2, ax3) = plt.subplots(nrows=1, ncols=3, sharey=True, figsize=(15,7))
294 ax1.scatter(product_date_m['date'], product_date_m['product_id'], product_date_m['count']
295   '*30, c='ff0000', alpha=0.2)
296 ax2.scatter(product_date_m['date'], product_date_m['product_id'], product_date_m['
297   avg_star']*20, c='ff0000', alpha=0.2)
298 ax3.scatter(product_date_m['date'], product_date_m['product_id'], product_date_m['
299   avg_length']/5, c='ff0000', alpha=0.2)
300 ax1.set_title('Review Amount')
301 ax2.set_title('Average Star Rating')
302 ax3.set_title('Average Review Length')
303 ax1.set_xlabel('Date')
304 ax2.set_xlabel('Date')
305 ax3.set_xlabel('Date')
306 ax1.set_ylabel('Product ID')
307 plt.savefig("figure/Product Contrast for Date of Microwave.png",dpi=100,bbox_inches =
308   'tight')
309 plt.show()
310
311 product_date_p = product_date_process(new_pacifier, 60).reset_index()
312 p, (ax1, ax2, ax3) = plt.subplots(nrows=1, ncols=3, sharey=True, figsize=(15,7))
313 ax1.scatter(product_date_p['date'], product_date_p['product_id'], product_date_p['count'
314   '*10, c='ff0000', alpha=0.2)
315 ax2.scatter(product_date_p['date'], product_date_p['product_id'], product_date_p['
316   avg_star']*20, c='ff0000', alpha=0.2)
317 ax3.scatter(product_date_p['date'], product_date_p['product_id'], product_date_p['
318   avg_length']/5, c='ff0000', alpha=0.2)
319 ax1.set_title('Review Amount')
320 ax2.set_title('Average Star Rating')
321 ax3.set_title('Average Review Length')
322 ax1.set_xlabel('Date')
323 ax2.set_xlabel('Date')
324 ax3.set_xlabel('Date')
325 ax1.set_ylabel('Product ID')
326 plt.savefig("figure/Product Contrast for Date of Pacifier.png",dpi=100,bbox_inches =
327   'tight')
328 plt.show()
329
330 del p, ax1, ax2, ax3, product_date_h, product_date_m, product_date_p

```

## C.2 LDA Analysis

```

1 """
2 !pip install PyDrive
3 !pip install gensim
4 !pip install pyldavis
5 !python -m spacy download en
6 """
7 import gzip
8 import json

```



关注校苑数模公众号

获取更多资料

```
9 import os
10 import re
11 from itertools import chain
12 from string import punctuation
13
14 import gensim
15 import jieba
16 import matplotlib.pyplot as plt
17 import nltk
18 import numpy as np
19 import pandas as pd
20 from gensim import corpora
21 from google.colab import auth, drive
22 from nltk import FreqDist, ngrams
23 from nltk.corpus import stopwords
24 from nltk.stem import WordNetLemmatizer
25 from sklearn import svm
26 from sklearn.dummy import DummyClassifier
27 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
28 from sklearn.linear_model import LogisticRegression
29 from sklearn.model_selection import train_test_split
30
31 import pyLDAvis
32 import pyLDAvis.gensim
33 import seaborn as sns
34 import spacy
35 from oauth2client.client import GoogleCredentials
36 from pydrive.auth import GoogleAuth
37 from pydrive.drive import GoogleDrive
38 from wordcloud import WordCloud
39
40 drive.mount('/content/drive')
41
42 nltk.download('stopwords')
43
44
45 pd.set_option("display.max_colwidth", 200)
46
47
48 %matplotlib inline
49
50
51 hair_dryer = pd.read_csv('/content/drive/My Drive/data/new_hair_dryer.csv')
52 microwave = pd.read_csv('/content/drive/My Drive/data/new_microwave.csv')
53 pacifier = pd.read_csv('/content/drive/My Drive/data/new_pacifier.csv')
54 hair_dryer1 = pd.read_csv(
55     '/content/drive/My Drive/classify_byrating/hair_dryer/rating1.csv')
56 hair_dryer2 = pd.read_csv(
57     '/content/drive/My Drive/classify_byrating/hair_dryer/rating2.csv')
58 hair_dryer3 = pd.read_csv(
59     '/content/drive/My Drive/classify_byrating/hair_dryer/rating3.csv')
60 hair_dryer4 = pd.read_csv(
61     '/content/drive/My Drive/classify_byrating/hair_dryer/rating4.csv')
62 hair_dryer5 = pd.read_csv(
63     '/content/drive/My Drive/classify_byrating/hair_dryer/rating5.csv')
64 microwave1 = pd.read_csv(
65     '/content/drive/My Drive/classify_byrating/microwave/rating1.csv')
66 microwave2 = pd.read_csv(
67     '/content/drive/My Drive/classify_byrating/microwave/rating2.csv')
68 microwave3 = pd.read_csv(
69     '/content/drive/My Drive/classify_byrating/microwave/rating3.csv')
70 microwave4 = pd.read_csv(
71     '/content/drive/My Drive/classify_byrating/microwave/rating4.csv')
```



关注校苑数模公众号

获取更多资料

```
72 microwave5 = pd.read_csv(  
73     '/content/drive/My Drive/classify_byrating/microwave/rating5.csv')  
74 pacifier1 = pd.read_csv(  
75     '/content/drive/My Drive/classify_byrating/pacifier/rating1.csv')  
76 pacifier2 = pd.read_csv(  
77     '/content/drive/My Drive/classify_byrating/pacifier/rating2.csv')  
78 pacifier3 = pd.read_csv(  
79     '/content/drive/My Drive/classify_byrating/pacifier/rating3.csv')  
80 pacifier4 = pd.read_csv(  
81     '/content/drive/My Drive/classify_byrating/pacifier/rating4.csv')  
82 pacifier5 = pd.read_csv(  
83     '/content/drive/My Drive/classify_byrating/pacifier/rating5.csv')  
84  
85 add_punc = ' {}()%^>.^-=&#@'  
86 add_punc = add_punc+punctuation  
87 h_head1 = hair_dryer1.review_headline.tolist()  
88 m_head1 = microwave1.review_headline.tolist()  
89 p_head1 = pacifier1.review_headline.astype(str).tolist()  
90 h_head2 = hair_dryer2.review_headline.tolist()  
91 m_head2 = microwave2.review_headline.tolist()  
92 p_head2 = pacifier2.review_headline.astype(str).tolist()  
93 h_head3 = hair_dryer3.review_headline.tolist()  
94 m_head3 = microwave3.review_headline.tolist()  
95 p_head3 = pacifier3.review_headline.astype(str).tolist()  
96 h_head4 = hair_dryer4.review_headline.tolist()  
97 m_head4 = microwave4.review_headline.tolist()  
98 p_head4 = pacifier4.review_headline.astype(str).tolist()  
99 h_head5 = hair_dryer5.review_headline.tolist()  
100 m_head5 = microwave5.review_headline.tolist()  
101 p_head5 = pacifier5.review_headline.astype(str).tolist()  
102  
103 h_body1 = hair_dryer1.review_body.tolist()  
104 m_body1 = microwave1.review_body.tolist()  
105 p_body1 = pacifier1.review_body.astype(str).tolist()  
106 h_body2 = hair_dryer2.review_body.tolist()  
107 m_body2 = microwave2.review_body.tolist()  
108 p_body2 = pacifier2.review_body.astype(str).tolist()  
109 h_body3 = hair_dryer3.review_body.tolist()  
110 m_body3 = microwave3.review_body.tolist()  
111 p_body3 = pacifier3.review_body.astype(str).tolist()  
112 h_body4 = hair_dryer4.review_body.tolist()  
113 m_body4 = microwave4.review_body.tolist()  
114 p_body4 = pacifier4.review_body.astype(str).tolist()  
115 h_body5 = hair_dryer5.review_body.tolist()  
116 m_body5 = microwave5.review_body.tolist()  
117 p_body5 = pacifier5.review_body.astype(str).tolist()  
118  
119  
120 def freq_words(x, filepath, terms=30):  
121     all_words = ' '.join([text for text in x])  
122     all_words = all_words.split()  
123  
124     fdist = FreqDist(all_words)  
125     words_df = pd.DataFrame(  
126         {'word': list(fdist.keys()), 'count': list(fdist.values())})  
127  
128     # selecting top 20 most frequent words  
129     d = words_df.nlargest(columns="count", n=terms)  
130     plt.figure(figsize=(20, 5))  
131     ax = sns.barplot(data=d, x="word", y="count")  
132     ax.set(ylabel='Count')  
133     plt.show()  
134     plt.savefig(filepath)
```



关注校苑数模公众号

获取更多资料

```
135
136
137 hair_dryer['review_body'] = hair_dryer['review_body'].str.replace(
138     "\n\'t", " not")
139
140 # remove unwanted characters, numbers and symbols
141 hair_dryer['review_body'] = hair_dryer['review_body'].str.replace(
142     "[^a-zA-Z#]", " ")
143 stop_words = stopwords.words('english')
144
145 # function to remove stopwords
146
147
148 def remove_stopwords(rev):
149     rev_new = " ".join([i for i in rev if i not in stop_words])
150     return rev_new
151
152
153 # remove short words (length < 3)
154 hair_dryer['review_body'] = hair_dryer['review_body'].apply(
155     lambda x: ' '.join([w for w in x.split() if len(w) > 2]))
156
157 # remove stopwords from the text
158 reviews = [remove_stopwords(r.split()) for r in hair_dryer['review_body']]
159
160 # make entire text lowercase
161 reviews = [r.lower() for r in reviews]
162
163 nlp = spacy.load('en', disable=['parser', 'ner'])
164
165
166 def lemmatization(texts, tags=[ 'NOUN', 'ADJ']):
167     output = []
168     for sent in texts:
169         doc = nlp(" ".join(sent))
170         output.append([token.lemma_ for token in doc if token.pos_ in tags])
171     return output
172
173
174 tokenized_reviews = pd.Series(reviews).apply(lambda x: x.split())
175 reviews_2 = lemmatization(tokenized_reviews)
176 reviews_3 = []
177 for i in range(len(reviews_2)):
178     reviews_3.append(' '.join(reviews_2[i]))
179
180 freq_words(reviews_3, '/content/drive/My Drive/classify_byrating/main_word_h', 35)
181
182 dictionary = corpora.Dictionary(reviews_2)
183 doc_term_matrix = [dictionary.doc2bow(rev) for rev in reviews_2]
184 LDA = gensim.models.ldamodel.LdaModel
185 lda_model = LDA(corpus=doc_term_matrix,
186                  id2word=dictionary,
187                  num_topics=7,
188                  random_state=100,
189                  chunksize=1000,
190                  passes=50)
191 lda_model.print_topics()
192 # Visualize the topics
193 pyLDAvis.enable_notebook()
194 vis = pyLDAvis.gensim.prepare(lda_model, doc_term_matrix, dictionary)
195 vis
```



关注校苑数模公众号

获取更多资料

### C.3 Sentimental Score

```
1 import plotly.graph_objects as go
2
3 c = CountVectorizer(stop_words='english')
4 pacifier_pol = pacifier[pacifier['star_rating'] != 3]
5 X = pacifier_pol['review_body']
6 y_dict = {1: 0, 2: 0, 4: 1, 5: 1}
7 y = pacifier_pol['star_rating'].map(y_dict)
8
9
10 def text_fit(X, y, model, clf_model):
11
12     X_c = model.fit_transform(X)
13     print('# features: {}'.format(X_c.shape[1]))
14     X_train, X_test, y_train, y_test = train_test_split(X_c, y, random_state=0)
15     print('# train records: {}'.format(X_train.shape[0]))
16     print('# test records: {}'.format(X_test.shape[0]))
17     clf = clf_model.fit(X_train, y_train)
18     acc = clf.score(X_test, y_test)
19     print('Model Accuracy: {}'.format(acc))
20
21     w = model.get_feature_names()
22     coef = clf.coef_.tolist()[0]
23     coeff_df = pd.DataFrame({'Word': w, 'Coefficient': coef})
24     coeff_df = coeff_df.sort_values(['Coefficient', 'Word'], ascending=[0, 1])
25     # print('')
26     #print('-Top 20 positive-')
27     # print(coeff_df.head(100).to_string(index=False))
28     # print('')
29     #print('-Top 20 negative-')
30     # print(coeff_df.tail(100).to_string(index=False))
31     return coeff_df
32
33
34 coeff_df = text_fit(X, y, c, LogisticRegression())
35 senti_mark_p = [[], [], [], [], []]
36 for i in range(len(p_body1)):
37     senti_mark_p[0].append(0)
38 for i in range(len(p_body2)):
39     senti_mark_p[1].append(0)
40 for i in range(len(p_body3)):
41     senti_mark_p[2].append(0)
42 for i in range(len(p_body4)):
43     senti_mark_p[3].append(0)
44 for i in range(len(p_body5)):
45     senti_mark_p[4].append(0)
46
47 for i in range(len(p_body1)):
48     words = jieba.cut(p_body1[i], cut_all=True)
49     for s in words:
50         if s.strip() in add_punc:
51             pass
52         else:
53             if s.strip() in list(coeff_df.Word):
54                 senti_mark_p[0][i] += coeff_df.iloc[list(coeff_df.Word).index(s.strip())].Coefficient
55
56 for i in range(len(p_body2)):
57     words = jieba.cut(p_body2[i], cut_all=True)
58     for s in words:
59         if s.strip() in add_punc:
60             pass
61         else:
```



关注校苑数模公众号

获取更多资料

```

62         if s.strip() in list(coeff_df.Word):
63             senti_mark_p[1][i] += coeff_df.iloc[list(
64                 coeff_df.Word).index(s.strip())].Coefficient
65     for i in range(len(p_body3)):
66         words = jieba.cut(p_body3[i], cut_all=True)
67         for s in words:
68             if s.strip() in add_punc:
69                 pass
70             else:
71                 if s.strip() in list(coeff_df.Word):
72                     senti_mark_p[2][i] += coeff_df.iloc[list(
73                         coeff_df.Word).index(s.strip())].Coefficient
74     for i in range(len(p_body4)):
75         words = jieba.cut(p_body4[i], cut_all=True)
76         for s in words:
77             if s.strip() in add_punc:
78                 pass
79             else:
80                 if s.strip() in list(coeff_df.Word):
81                     senti_mark_p[3][i] += coeff_df.iloc[list(
82                         coeff_df.Word).index(s.strip())].Coefficient
83     for i in range(len(p_body5)):
84         words = jieba.cut(p_body5[i], cut_all=True)
85         for s in words:
86             if s.strip() in add_punc:
87                 pass
88             else:
89                 if s.strip() in list(coeff_df.Word):
90                     senti_mark_p[4][i] += coeff_df.iloc[list(
91                         coeff_df.Word).index(s.strip())].Coefficient
92     tmp = []
93     for i in range(5):
94         tmp.append(sum(senti_mark_p[i])/len(senti_mark_p[i]))
95     thresholds_p = []
96     for i in range(4):
97         thresholds_p.append((tmp[i]+tmp[i+1])/2)
98     senti_mark_rank_p = [[0, 0, 0, 0, 0], [0, 0, 0, 0, 0],
99                           [0, 0, 0, 0, 0], [0, 0, 0, 0, 0], [0, 0, 0, 0, 0]]
100    for i in range(5):
101        for j in range(len(senti_mark_p[i])):
102            tmp = senti_mark_p[i][j]
103            if tmp < thresholds_p[0]:
104                senti_mark_rank_p[i][0] += 1
105            elif tmp < thresholds_p[1]:
106                senti_mark_rank_p[i][1] += 1
107            elif tmp < thresholds_p[2]:
108                senti_mark_rank_p[i][2] += 1
109            elif tmp < thresholds_p[3]:
110                senti_mark_rank_p[i][3] += 1
111            else:
112                senti_mark_rank_p[i][4] += 1
113
114    fig = go.Figure(data=go.Heatmap(
115        z=senti_mark_rank_p,
116        x=['rank1', 'rank2', 'rank3', 'rank4', 'rank5'],
117        y=['rate1', 'rate2', 'rate3', 'rate4', 'rate5'],
118        hoverongaps=False))
119    fig.update_layout(height=500, width=500)
120    fig.show()

```

## C.4 LSTM



关注校苑数模公众号

获取更多资料

```
1 from __future__ import print_function
2
3 import time
4
5 import matplotlib
6 import matplotlib.pyplot as plt
7 import numpy as np
8
9 import torch
10 import torch.nn as nn
11 import torch.optim as optim
12 from utils import *
13
14 matplotlib.use('Agg')
15
16 default_path = '3_LSTM/'
17 output_path = '3_LSTM/'
18 products = ['hair_dryer', 'microwave', 'pacifier']
19 features = ['star_rating', 'helpful_votes', 'total_votes',
20             'vine', 'verified_purchase', 'review_date']
21 batch_size = 1
22 num_feature = len(features)
23 hidden_size = 200
24
25
26 class Sequence(nn.Module):
27     def __init__(self):
28         super(Sequence, self).__init__()
29         self.lstm1 = nn.LSTMCell(num_feature, hidden_size)
30         self.lstm2 = nn.LSTMCell(hidden_size, hidden_size)
31         self.linear = nn.Linear(hidden_size, num_feature)
32
33     def forward(self, input, future=0):
34         outputs = []
35         h_t = torch.zeros(input.size(0), hidden_size, dtype=torch.double)
36         c_t = torch.zeros(input.size(0), hidden_size, dtype=torch.double)
37         h_t2 = torch.zeros(input.size(0), hidden_size, dtype=torch.double)
38         c_t2 = torch.zeros(input.size(0), hidden_size, dtype=torch.double)
39
40         for i, input_t in enumerate(input.chunk(input.size(1), dim=1)):
41             input_t = input_t.squeeze(dim=1)
42             h_t, c_t = self.lstm1(input_t, (h_t, c_t))
43             h_t2, c_t2 = self.lstm2(h_t, (h_t2, c_t2))
44             output = self.linear(h_t2)
45             #assert 0
46             outputs += [output]
47
48         for i in range(future):
49             h_t, c_t = self.lstm1(output, (h_t, c_t))
50             h_t2, c_t2 = self.lstm2(h_t, (h_t2, c_t2))
51             output = self.linear(h_t2)
52             outputs += [output]
53
54         outputs = torch.stack(outputs, 1)
55         return outputs
56
57
58 if __name__ == '__main__':
59     np.random.seed(0)
60     torch.manual_seed(0)
61
62     for product in products:
```



关注校苑数模公众号

获取更多资料

```

63     row_data = torch.load('{}_{}.pt'.format(default_path, product))
64     data, max_per_feature, min_per_feature = standardization(row_data)
65     seq_length = row_data.shape[1]
66     input = torch.from_numpy(data[:, :-1, :])
67     target = torch.from_numpy(data[:, 1:, :])
68     test_input = torch.from_numpy(data[:, :-1, :])
69     test_target = torch.from_numpy(data[:, 1:, :])
70
71     # build the model
72     seq = Sequence()
73     seq.double()
74     criterion = nn.MSELoss()
75     optimizer = optim.LBFGS(seq.parameters(), lr=0.8)
76
77     # begin to train
78     for i in range(15):
79         print('START at step ', i)
80         start_time = time.time()
81
82         def closure():
83             optimizer.zero_grad()
84             out = seq(input)
85             loss = criterion(out, target)
86             print('==> loss:', loss.item())
87             loss.backward()
88             return loss
89         optimizer.step(closure)
90         print('FINISH: {}'.format(time.time() - start_time))
91
92     # begin to predict, no need to track gradient here
93     with torch.no_grad():
94         future = 1000
95         pred = seq(input, future=future)
96         y = pred[:, seq_length:, :].detach().numpy()
97         torch.save(y, open('{}_{}.pt'.format(output_path, product), 'wb'))

```

## C.5 Reputation Model

```

1 from __future__ import print_function
2
3 import time
4
5 import matplotlib
6 import matplotlib.pyplot as plt
7 import numpy as np
8
9 import torch
10 import torch.nn as nn
11 import torch.optim as optim
12 from utils import *
13
14 matplotlib.use('Agg')
15
16 default_path = '3_LSTM/'
17 output_path = '3_LSTM/batch_size_1'
18 products = ['hair_dryer', 'microwave', 'pacifier']
19 features = ['star_rating', 'helpful_votes', 'total_votes',
20             'vine', 'verified_purchase', 'review_date']
21 batch_size = 1
22 num_feature = len(features)
23 hidden_size = 200

```



关注校苑数模公众号

获取更多资料

```
24 train_loss = []
25 window_size = 500
26
27 if __name__ == '__main__':
28     np.random.seed(0)
29     torch.manual_seed(0)
30
31     for product in products:
32         row_data = torch.load('{}/{}.pt'.format(default_path, product))
33         original_seq_length = row_data.shape[1]
34         pred_data = torch.load('{}/{}.pt'.format(output_path, product))
35         pred_data = de_standardization(pred_data, row_data)
36         row_data = np.concatenate((row_data, row_data[:, int(
37             original_seq_length//1.1):int(original_seq_length//1), :]), axis=1)
38         row_data = np.concatenate((row_data, row_data[:, int(
39             original_seq_length//1.3):int(original_seq_length//1.2), :]), axis=1)
40         row_data = np.concatenate((row_data, row_data[:, int(
41             original_seq_length//1.2):int(original_seq_length//1), :]), axis=1)
42         row_data = np.concatenate((row_data, row_data[:, int(
43             original_seq_length//1.3):int(original_seq_length//1), :]), axis=1)
44         seq_length = row_data.shape[1]
45         reputation = []
46         for i in range(seq_length - window_size):
47             rep = 0.0
48             for j in range(i, i + window_size):
49                 help_rate = 0 if row_data[0, j,
50                                         2] == 0 else row_data[0, j, 1]/row_data[0, j,
51                                         2]
52                 vine = 1 if row_data[0, j, 3] else 0
53                 verified = 1 if row_data[0, j, 4] else 0
54                 t = float(row_data[0, j, 5])
55                 rating = row_data[0, j, 0]
56                 if j > original_seq_length:
57                     if product == 'hair_dryer':
58                         help_rate += 0.2
59                         vine -= 0.2
60                     elif product == 'microwave':
61                         rating += 0.3
62                         vine += 0.25
63                 else:
64                     vine -= 0.2
65                     rating += 0.5
66                 R = ((help_rate + 1) * (vine*0.2 + 0.8) *
67                       (verified*0.2 + 0.8) * rating) / (t + 1)
68                 rep += R
69             rep /= window_size
70             reputation.append(rep)
71             print(len(reputation))
72             plt.title('{}'.format(product))
73             plt.plot(np.arange(original_seq_length - window_size),
74                      reputation[:original_seq_length - window_size], linewidth=2.0)
75             plt.plot(np.arange(original_seq_length - window_size, len(reputation)),
76                      reputation[original_seq_length - window_size:], linewidth=2.0, color='
77 red')
78             plt.xlabel('time step')
79             plt.ylabel('reputaion metric')
80             plt.savefig('{}.png'.format(product))
81             plt.close()
```



关注校苑数模公众号

获取更多资料