

Sentiment Analysis of Products' Reviews Containing English and Hindi Texts

Jyoti Prakash Singh¹(✉), Nripendra P. Rana², and Wassan Alkhowaiter³

¹ National Institute of Technology Patna, Patna, Bihar, India
jyotip.singh@gmail.com

² School of Management, Swansea University, Swansea SA2 8PP, UK
n.p.rana@swansea.ac.uk

³ Qassim University, Buraidah, Saudi Arabia
w.alkhowaiter@gmail.com

Abstract. The online shopping is increasing rapidly because of its convenience to buy from home and comparing products from their reviews written by other purchasers. When people buy a product, they express their emotions about that product in the form of review. In Indian context, it is found that the reviews contain Hindi text along with English. It is also found that most of the Hindi text contains opinionated words like *bahut achha*, *bakbas*, *pesa wasool* etc. We have tried to find out different Hindi texts appearing in product reviews written on Indian E-commerce portals. We have also developed a system which takes all those reviews containing Hindi as well as English texts and find out the sentiment expressed in that review for each attribute of the product as well as a final review of the product.

Keywords: Sentiment analysis · POS-Tagging · Review analysis · Product summarization

1 Introduction

The life style of society is changing with the penetration of Internet, and E-commerce in every corner of the world. Earlier, the advertisement and friends recommendations were a major source of information while buying a product. The number of recommendations was a limited one to compare similar products of different brands. Nowadays, as the e-commerce business has grown up, they are offering more products. The e-commerce websites also request their customers to write their experience about the product they brought in the form of a product review. These reviews offer significant information to buyers about the product they are planning to buy and also enable them to compare products of different brands. The reviews help consumers to choose the best products by comparing them based on other consumers' evaluation of the products. It also aids in the improvement of the product by informing the manufacturers about the advantages and defects of their products. The number of reviews about the

products grows with the growth of e-commerce businesses. It becomes very difficult for buyers and sellers to manually analyze a large number of reviews and get any meaningful information. This attracts a lot of researchers to automate the analysis of reviews and get valuable information hidden in the reviews [3, 5].

The reviews written by Indian buyers are mainly in English, but it contains some Hindi texts (written in English Scripts only) also as Hindi is a prevalent language in India. Some of the most widely used Hindi words like *bahut achha*, *bakbas*, *pesa wasool* are found in a number of reviews. Most of these words are opinionated and contain strong opinions in the form of *good* or *bad*. Most of the earlier work done in the area of finding polarity of opinions for product reviews neglect these texts as they are mainly developed for English texts only. As per the best of our knowledge, no work has yet been reported which consider the correction of these typos and includes the sentiment of the Hindi words along with English texts.

In this work, we have proposed a sentiment analysis system which works for reviews containing both English as well as Hindi opinionated texts. First of all we have gathered possible Hindi opinionated texts from reviews appearing on Indian popular E-commerce sites such as *amazon.in*, *flipkart.com*, *snapdeal.com*, *shopclues.com* and so on. These Hindi texts are preprocessed and their equivalent English words are found. The summarized review of the product is then calculated consulting sentiwordnet database.

The rest of the paper is organized as follows: The proposed system architecture and algorithm are discussed in Sect. 2. In Sect. 3, we present our results and finally in Sect. 4, we conclude the paper.

2 Proposed Work

Product reviews contents from popular Indian e-commerce sites like *amazon.in*, *flipkart.com*, *snapdeal.com*, *shopclues.com* are collected as our dataset. The dataset has a lot of typos in the form of joint words like *verygood* as well as abbreviations containing numerals such as *gr8* for *great*. The dataset also contains Hindi words like "*bakbas*", "*bekar*", "*achchha*" (written in English script only). Some sample typos gathered from various Indian e-commerce websites are given in Tables 2 and 3. Table 1 contains words of Hindi Texts typed in English along with their English equivalent text. Table 2 contains some popular abbreviations used online for review, chatting, etc. along with their correct form in English. Some joint words (missing space) are shown in Table 3.

One of the primary focuses of this work is to pre-process the product review available on Indian E-commerce sites so that the reviews contain only English text. Once reviews are converted to English text, Part of Speech (POS) Tagging to the text is done using wordnet [4] database. Once POS tagging is done, the adjective, noun, and adverb are extracted. Further, sentiwordnet [1, 2] database is used to assign numerical values to the adjectives contained in the review. The proposed system architecture is shown in Fig. 1.

Table 1. List of wrong words in Hindi

Wrong words	Corrected words
Ye	This
Achha	Good
G8t	Great
N8t	Night
H	Is
Som1	Some one

Table 2. List of wrong words in English

Wrong words	Corrected words
Gud	Good
Gooood	Good
Exclent	Excellent
Bd	Bad
Awesm	Awesome

A pseudo code for our proposed system is given below.

Proposed Algorithm:

Step 1: Tokenize based on space.

Step 2: Consult wordnet

If word is matched, then go to POS tagger.

else correct word and go to POS Tagger

Step 4: Noun, Adverb, and Adjective are stored in frequent feature database.

Step 5: Generate the product summary with the help of SentiWordNet Lexical databases.

The working of our scheme is traced with the aid of an example presented here. *Yeh achha camera h. eski pictre quality bahut achhi hai. The pics resolution is enough. Zoom is bakbas. focus is verygd bt not g8t.*

The document (Complete review) is broken down into several sentences based on [., [?], And [!] Mark. For example review, sentences are:

S1. Yeh achha camera h.

S2. eski pictre quality bahut achhi hai.

S3. The pics resolution is enough.

S4. Zoom is bakbas.

S5. focus is verygd bt not g8t.

Yeh is a Hindi word whose English equivalent is *This*. *Achha* is another Hindi word meaning *good* in English. The complete review is written in English after correcting and converting every word to English.

Table 3. List of joint words

Jointed words	Corrected words
Verygood	Very good
Verybad	Very bad
Bahutbura	Bahut bura

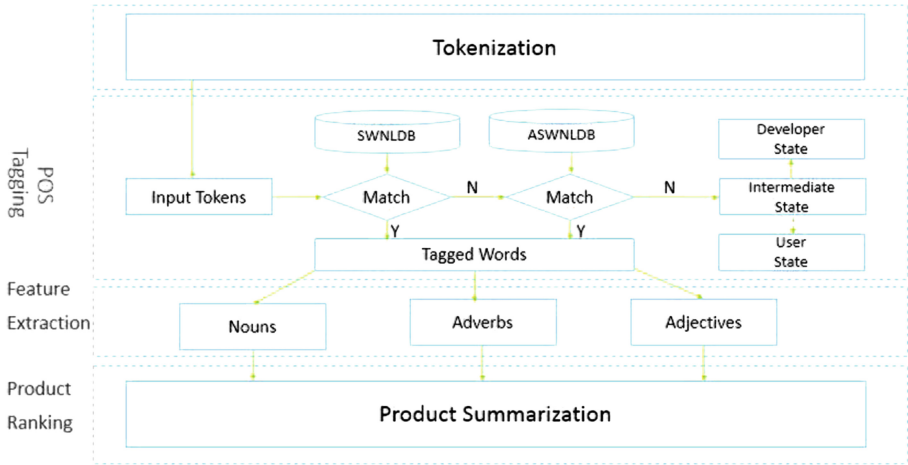


Fig. 1. Proposed system architecture

- s1. This is a good camera.
- s2. Its picture quality is very good.
- s3. The picture resolution is enough.
- s4. Zoom is bad.
- s5. Focus is very good but not great.

The POS tagging is applied as given below for just one sentence. We have used Penn Treebank tagset for Part of Speech Tagging

This===== [This_DT]
 Good===== [Good_JJ]
 Camera===== [Camera_NN]
 Is===== [Is_VBZ]
 . =====[_.UH]

Next step is to consult the sentiwordnet database to find the priority of the adjective to find the sentiment value of the sentence. The score of every adverb and adjective are given in the SentiWordNet lexical database. We have listed here some of them given in Table 4 which are going to be used in the above example.

Table 4. Score list of words

Word	Orientation of word	Score
Good	Positive	.75
Great	Positive	.875
Awesome	Positive	.875
Excellent	Positive	1
Well	Positive	.75
Average	Positive	.375
Enough	Neutral	.875
Bad	Negative	.65
Very	Nil	.5
Not	Negation	-1

Table 5. Sentence wise score

Sentence number	Score	Score type
S1	0.750	Positive
S2	1.250	Positive
S3	0.125	Positive
S4	0.625	Negative
S5	0.275	Positive

Where nil represents the neither positive nor negative orientation of words. And negation represents the multiplier factor which having -1 value.

For above example, sentence s1 has good as adjective whose sentiment score is $+0.75$. In this sentence there is no adverb or negation, so the sentiment score for sentence s1 is $+0.75$. For second sentence s2, adverb (*very*) is there with the adjective (*good*), so the sentiment score of s2 is 1.25, which is a sum of scores of good (0.75) and very (0.5). The sentiment score of each sentence is shown in Table 5.

The polarity of the review of a product is determined by finding the polarity of each feature of the product across all reviews and finding a weighted sum of all features.

3 Result

We have collected 1100 reviews from *flipkart.com* and *amazon.in* of three popular mobile brands in India at the time of writing this paper. The results show that for Android based smart-phone people are talking about features like *camera*, *battery*, *memory*, *processor*, *RAM*, *display*, *price*, *weight* and *phone*. Out of these features, battery, camera and display are found to be more prominent ones across

Table 6. Summary of the review of product **Samsung Galaxy S3 Neo**, **Asus Zenfone 2** and **Honor 4X**

Phone feature		Galaxy S3 Neo	Zenfone 2	Honor 4X
<i>Feature Name</i>	<i>Scores Type</i>	<i>Percentage(%)</i>	<i>Percentage(%)</i>	<i>Percentage(%)</i>
Camera	Positive	100	66.39	60.17
	Negative	0.0	33.64	39.82
Battery	Positive	61.76	90.82	66.09
	Negative	38.23	9.18	33.9
Memory	Positive	100	100	49.24
	Negative	0.0	0.0	50.76
Processor	Positive	100	100	49.25
	Negative	0.0	0.0	50.74
RAM	Positive	100	100	48.24
	Negative	0.0	0.0	51.75
Display	Positive	90.5	0.0	100
	Negative	9.5	100	0.0
Price	Positive	0.0	No opinion	No opinion
	Negative	100	No opinion	No opinion
Weight	Positive	100	0.0	No opinion
	Negative	0.0	100	No opinion
Phone	Positive	55.07	89.09	62.02
	Negative	44.93	10.91	37.98
Overall	Positive	68.86	75.14	62.60
	Negative	31.134	21.15	32.26

all phones. The results are shown in Table 6. *no opinion* in both positive and negative rows shows that no one has given any opinion about that feature of that product.

4 Conclusion

We have designed a sentiment analysis system which can take reviews written in Hindi as well as English texts and find the sentiment of customers for that product. We have taken a dictionary based approach to correct the wrong words and replace Hindi text with their English equivalent. We further want to extend this system with a machine learning algorithm to correct the wrong words and Hindi words. The final opinion score is a weighted average of all the features of the product under consideration. We are also working to identify the most prominent features of a product to calculate the final opinion score.

References

1. Baccianella, S., Esuli, A., Sebastiani, F.: Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. In: LREC, vol. 10, pp. 2200–2204 (2010)
2. Esuli, A., Sebastiani, F.: Sentiwordnet: a publicly available lexical resource for opinion mining. In: Proceedings of the 5th Conference on Language Resources and Evaluation (LREC 2006), pp. 417–422 (2006)
3. Hu, M., Liu, B.: Mining and summarizing customer reviews. Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. KDD 2004, pp. 168–177. ACM, New York (2004)
4. Miller, George A.: Wordnet: a lexical database for english. *Commun. ACM* **38**(11), 39–41 (1995)
5. Popescu, A.-M., Etzioni, O.: Extracting product features and opinions from reviews. Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing. HLT 2005, pp. 339–346. ACL, Stroudsburg (2005)