A Mini Project in

# Product Review Sentimental Analysis

*Submitted in partial fulfillment of
the requirements for the award of the degree of*

**Bachelor of Technology
In
Computer Science and Engineering**

By

| Reg. No. | Name Of Student |
|----------|-----------------|
| 20134044 | Shivam Mohan |
| 20134086 | Shreyash Hisariya |
| 20134022 | Prashant Agrawal |
| 20134166 | Ankit Kumar Sharma |
| 20134164 | Banda Prashanth Yadav |

Project Group – CS26

To the

**COMPUTER SCIENCE AND ENGINEERING DEPARTMENT**

**MOTILAL NEHRU NATIONAL INSTITUTE OF TECHNOLOGY**

**ALLAHABAD**

April 2016

# UNDERTAKING

I declare that the work presented in this project titled "**Product Review Analysis**", submitted to the Computer Science and Engineering Department, Motilal Nehru National Institute of Technology, Allahabad, for the award of the Bachelor of Technology degree in Computer Science & Engineering, is my original work. I have not plagiarized or submitted the same work for the award of any other degree. In case this undertaking is found incorrect, I accept that my degree may be unconditionally withdrawn.

April 2016.

# CERTIFICATE

Certified that the work contained in the project titled "**Product Review Analysis**", by *Shivam Mohan, Shreyash Hisariya, Prashant Agrawal, Ankit Kumar Sharma and Banda Prashanth Yadav* has been carried out under my supervision and that this work has not been submitted elsewhere for a degree.

_____

(Er. Rupesh Dewang)

**Computer Science and Engineering Department**

Motilal Nehru National Institute of Technology Allahabad

# Preface

On-line reviews on services and products are fast becoming an important factor while buying/using that product/service. India's e-Commerce market focuses on the various sub-segments of the e-Commerce market and highlights the factors driving growth across these segments and the challenges.

India's e-commerce market was worth about $3.9 billion in 2009, it went up to $12.6 billion in 2013. In 2013, the e-retail segment was worth US$2.3 billion. About 70% of India's e-commerce market is travel related. According to Google India, there were 35 million online shoppers in India in 2014 and is expected to cross 100 million mark by end of year 2016. India's retail market is estimated at $470 billion in 2011 and is expected to grow to $675 Bn by 2016 and $850 Bn by 2020, – estimated CAGR of 10%.

Buying decisions on these websites are greatly influenced by the product reviews given by previous purchasers. In this regard, many reviews are written by people who have not actually purchased/used the product. These reviews may misguide a purchaser's decision regarding a particular product.

# Acknowledgements

# Contents

# Chapter 1

# Introduction

Nowadays, most products and services are listed online and welcome online feedback from their customers. The feedback provided is mostly in the form of user-ratings and reviews (which gives in-depth analysis of the rating given by a user). Ratings are quantified values hence does not explain details of the product/service. Reviews on the other hand are useful data which a potential customer uses in order to make a choice.

As the number of reviews of a particular product may be large, it is a tedious task to analyze each and every review manually. So here we came up with an idea to calculate the percentage of positive and negative reviews programmatically, which enables the buyers to compare and choose between different products in the market.

## 1.1 Motivation

The motivation for doing this project was primarily an interest in undertaking a challenging project. In the current scenario reviews play a very important role.

88% Of Consumers Trust Online Reviews As Much As Personal Recommendations Each new review written about a product on Company's site increases the amount of unique content site offers on that product , it will be seen as having higher relevance and be useful in deciding future strategies.

Reviews serve as a strong source of "word-of-mouth" communication for next customer in his purchase decision.

### 1.1.1 Some Wonderful Minds

Bing Liu, professor of opinion mining at the University of Illinois at Chicago, is working on detecting fake or deceptive opinions on social media platforms like re-views, Facebook, Twitter, Weibo, and forum discussion sites, and also is a author of sentiment analysis and opinion mining. Nitin Jindal is a professor at University of Illinois at Chicago. It is interesting to know that Jindal and Liu together had worked on many research papers together on detecting review spammers.

# Chapter 2

# Hypothesis

## 2.1  Initial Approach

### 2.1.1 Web Crawler

A Web crawler is a computer program that browses the World Wide Web in a methodical, automated manner to gather or collect particular set of data.

We designed a web crawler which could collect reviews of mobiles from an E-Commerce website (Snap-deal) which could be treated as a dataset for further analysis of reviews. We performed some preprocessing and remove the unwanted data present.

Next step would be classifying the product overall whether the mobile product is favorable or non-favorable to buy.

## 2.2  Data Collection

The content format of the **Snapdeal** website is 48 products per page which comprises 10 reviews per product per page. In order to obtain all the reviews for each product of a particular category (mobiles), a python script was written that automatically traverses each page and extracts the reviews.

We have obtained the HTML Source page by using the *Requests* library and then relevant tags are extracted using the *BeautifulSoup4* library. We traversed the links from the tags and extracted the reviews of the products. This process is repeated for all the products. We even made this to work with authenticated proxy.

## 2.3   Preprocessing

### 2.3.1   Punctuation Removal

Punctuations do not contribute anything to the meaning and context of words when one word is processed at a time. We wrote a Python script that removed punctuation from the documents and gave us a single document.

### 2.3.2   Stop words Removal

Sometimes, some extremely common words which would appear to be of little value in helping select, documents matching the user need are excluded from the vocabulary entirely. These words are called **stop words**. To remove such words, we used NLTK's built-in library.

### 2.3.3   Lemmatizing

The goal of lemmatization is to reduce inflectional forms and sometimes derivationally related forms of a word to a common base form. For instance:

*am, are, is => be*
*car, cars, car's, cars' => car*

The result of this mapping of text will be something like:

*The boy's cars are different colors =>*
*The boy car be differ color.*

## 2.4   Sentiment Calculation

Sentiment Analysis is the process of determining whether a piece of writing is positive, negative or neutral. It's also known as opinion mining, deriving the opinion or attitude of a speaker. It is used to identify and extract subjective information in source materials. It aims to determine the overall attitude and contextual polarity of a document with respect to some topic.

We have calculated the overall polarity of a product considering all the collected reviews.

We have used two dictionaries, one for positive and other for negative words. These dictionaries contain words with their corresponding weightages, which helps in summing up the polarity of the review. By comparing each word from the reviews with these dictionaries, we determine the polarity of the products.

## 2.5   Probability Determination:

We have determined probability of a particular word, whether the word has been used in terms of positive or negative respect. We have performed this probability determination by using the Naïve Bayes theorem.

**Naïve Bayes classifiers** are a family of simple probabilistic classifiers based on applying Bayes' theorem with strong (naive) independence assumptions between the features.

Abstractly, naive Bayes is a conditional probability model: given a problem instance to be classified, represented by a vector $x = (x_1, \ldots, x_n)$ representing some n features (independent variables), it assigns to this instance probabilities.

$$p(C_k | x_1, \ldots, x_n)$$

for each of *K* possible outcomes or *classes*.

The problem with the above formulation is that if the number of features *n* is large or if a feature can take on a large number of values, then basing such a model on probability tables is infeasible. We therefore reformulate the model to make it more tractable. Using Bayes' theorem, the conditional probability can be decomposed as

$$p(C_k | \mathbf{x}) = \frac{p(C_k)\ p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

In plain English, using Bayesian probability terminology, the above equation can be written as

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$

# Chapter 3

## Result

## 3.1  Sample Product Reviews:

## 3.2   Sentiment Calculation

iPhone_5S_16GB_Space_Gray

```
1 Number of sentiment words : 2537
2
3 Most used positive words :
4 good  -  365
5 best  -  175
6 great  -  113
7 like  -  102
8 happy  -  84
9 nice  -  81
10 awesome  -  71
11 perfect  -  69
12 well  -  69
13 fast  -  68
14 work  -  67
15 love  -  66
16 genuine  -  63
17 recommend  -  39
18 excellent  -  33
19 thank  -  30
20 easy  -  29
21 worth  -  29
22 satisfy  -  28
23 satisfied  -  26
24
25 Most used negative words :
26 problem  -  33
27 issue  -  22
28 cheap  -  13
29 bad  -  10
30 lag  -  10
31 doubt  -  9
32 hang  -  8
33 worry  -  7
36 hassle  -  5
37 expensive  -  4
38 delay  -  4
39 hard  -  3
40 costly  -  3
41 complaint  -  3
42 fuss  -  3
43 fault  -  3
44 poor  -  3
45 loose  -  3
46
47 Percentage of positive words : 92.279877
48 Percentage of negative words : 7.720123
```

## 3.3  Naïve Bayes Theorem Implementation



```
shivam@shivam-HP-ProBook-4430s: ~/Desktop/python-naive-bayes-master
shivam@shivam-HP-ProBook-4430s:~/Desktop/python-naive-bayes-master$ python nb_3.
py
Count of 'buy' word in review file 130.0
Count of 'dont' word in review file 38.0
Probability of 'buy' word 0.00319685233001
Probability of 'dont' word 0.000934464527235
Probability of buy when review is positive  : 0.0824735680434
Probability of buy when review is negative  : 0.0483428387864
shivam@shivam-HP-ProBook-4430s:~/Desktop/python-naive-bayes-master$
```

# Chapter 4

# Conclusion

So far we have successfully:

- Extracted mobile reviews from *Snapdeal* website.
- Converted them to excel format for ease of storage and implementing database model (if need arises).
- Calculated sentiment value of all mobile product reviews (score in percentage).
- Applied Naïve Bayes theorem to calculate conditional probability of occurrence of word in positive or negative sentiment.

# Chapter 5

## Limitations

- Inability to detect Sarcasm. Oxford dictionary defines sarcasm as the use of irony to mock or convey contempt.

  *Example*: *The price of iPhone is very cheap.*

- Inability to process slangs.

- Inability to expand general abbreviations.

# Bibliography

- Myle Ott: Deceptive Opinion Spam Corpus v1.4. http://myleott.com/op spam/
- NLTK Documentation. http://www.nltk.org/
- BeautifulSoup4 Documentation: http://www.crummy.com/software/BeautifulSoup/bs4/doc
- Math Works Documentation.
- Nitin Jindal and Bing Liu. Opinion Spam and Analysis (2008)