

# ONLINE PRODUCTS RECOMMENDATION SYSTEM USING SVM CLASSIFICATION AND DETECT FRAUD REVIEWS

PRIYANGHA.K<sup>1</sup>, NAGAPPAN.V.K<sup>2</sup>

<sup>1</sup>PG Student, Dept of CSE, Bharathiyar College of Engineering and Technology, Karaikal.

<sup>2</sup>Assistant Professor, Dept. of CSE, Bharathiyar College of Engineering and Technology, Karaikal.

**Abstract**— User review is a crucial component of open mobile app market such as Google Play Store. These markets allow users to submit feedback for downloaded apps in the form of a star ratings and opinions in the form of text reviews. Users read these reviews in order to gain insight into the app before they buy or download it. The user opinion about the product also influences on the purchasing decisions of potential users; indeed play a key role in the generation of revenue for the developers. The product can contain large volumes of reviews and it is impossible for a user to skim through thousands of reviews to find the opinion of other users about the features he/she is interested in. Towards this end, we propose a methodology to automatically extract the features of an app from its corresponding reviews using machine learning technique. Moreover, our proposed methodology aid user to compare the features across multiple apps, using the sentiments, expressed in their associated reviews. The proposed methodology can be used to understand user's preference to a certain mobile app and can uncover the relational behind why users prefer an app over other. Ranking fraud in the mobile App market refers to fraudulent or deceptive activities which have a purpose of bumping up the Apps in the popularity list. Indeed, it becomes more and more frequent for App developers to use shady means, such as inflating their Apps' sales or posting phony App ratings, to commit ranking fraud. While the importance of preventing ranking fraud has been widely recognized, there is limited understanding and research in this area.

**Keywords-** Support Vector Machine, Unified Modelling Language, Iterative Self Organising Data, K-Nearest Neighbors.

## 1. INTRODCUTION

Opinions are at the center of almost all human activities and are an important reflection of our behavior. Our belief in reality, our perceptions and our choices depend on how others see and appreciate the world. For this reason, we often refer to others' opinions when we need to make a

decision. This does not apply only to individuals. It also applies to organizations. In the real world, companies and organizations always want to receive opinions and comments from consumers or the public about their products and services. Individual consumers want to know the views of current users before a product is

purchased or the opinions of others about political candidates before giving a vote on political elections. Getting public opinion and consumer perspectives has long been a major workload for marketing, public relations and political campaign companies. With the help of social networks (for example, criticism, forum discussions, blogs, microblogs, Twitter comments and posts on social networking sites) and as a result of increased power in decision-making of social media, individuals and organizations have become inevitable to take into account the content of these media. In recent years, industrial activities involving sentiment analysis is also developing rapidly large number of new initiatives have emerged in this area. Many large corporations have developed their own sentiment analysis systems to measure the quality of on-site services, thereby creating awareness in the business and social environment

To Identify the Fraud in Ranking of product review system. It can implement product reviews rating from 1 reviews, which aim to automatically identify important product aspects from online consumer reviews.

The important aspects are identified according to two observations:

(a) The important aspects of product are usually commented by a large number of consumers.

(b) Consumers' opinions on the important aspects greatly influence their overall opinions on the product. In particular, given consumer reviews of product, we first identify the product aspects by clustering and determine consumers' opinions on these aspects via a sentiment classifier.

We then implement Support vector classification to identify the opinion words by simultaneously considering the reviews collection and the influence of consumers' opinions given to each aspect on their overall opinions. The experimental results on popular product reviews demonstrate the effectiveness of our approach. We further apply the review ranking results to the application of sentiment classification, and improve the performance significantly.

## **2. RELATED WORK**

### **A Joint Segmentation and Classification Framework for Sentence Level Sentiment Classification**

A joint segmentation and classification framework for sentence-level sentiment classification. It is widely recognized that phrasal information is crucial for sentiment classification. However, existing sentiment classification algorithms typically split a sentence as a word sequence, which does not effectively handle the inconsistent sentiment polarity between a phrase and the words it contains, such as {"not bad," "bad"} and {"a great deal of," "great"}. We address this issue by developing a joint framework for sentence-level sentiment classification. It simultaneously generates useful segmentations and predicts sentence-level polarity based on the segmentation results. Specifically, we develop a candidate generation model to produce segmentation candidates of a sentence; a segmentation ranking model to score the usefulness of a segmentation candidate for sentiment classification; and a classification model for predicting the sentiment polarity of a segmentation. We train the joint framework directly from sentences annotated with only sentiment polarity, without using any syntactic or sentiment annotations in segmentation level. We conduct experiments for sentiment classification on two benchmark datasets: a tweet dataset and a review dataset.

### **Polarity Consistency Checking for Domain Independent Sentiment Dictionaries**

Polarity classification of words is important for applications such as Opinion Mining and Sentiment Analysis. A number of sentiment word/sense dictionaries have been manually or (semi)automatically constructed. We notice that these sentiment dictionaries have numerous inaccuracies. Besides obvious instances, where the same word appears with different polarities in different dictionaries, the dictionaries exhibit complex cases of polarity inconsistency, which cannot be detected by mere manual inspection. We introduce the concept of polarity consistency of words/senses in sentiment dictionaries in this paper. We show that the consistency problem is NP-complete. We reduce the polarity consistency problem to the satisfiability problem and utilize two fast SAT solvers to detect inconsistencies in a sentiment dictionary. We perform experiments on five sentiment dictionaries and WordNet to show inter and intra-dictionaries inconsistencies.

### **Estimating the Helpfulness and Economic Impact of Product Reviews: Mining Text and Reviewer Characteristics**

the rapid growth of the Internet, the ability of users to create and publish content has created active electronic communities that provide a wealth of product information. However, the high volume of reviews that are typically published for a single product makes harder for individuals as well as manufacturers to locate the best reviews and understand the true underlying quality of a product. In this paper, we reexamine the impact of reviews on economic outcomes like product sales and see how different factors affect social outcomes such as their perceived usefulness. Our approach explores multiple aspects of review text, such as subjectivity levels, various measures of readability and extent of spelling errors to identify important text-based features. In addition, we also examine multiple reviewer-level features such as average usefulness of past reviews and the self-disclosed identity measures of reviewers that are displayed next to a review. Our econometric analysis reveals that the extent of subjectivity, informativeness, readability, and linguistic correctness in reviews matters in influencing sales and perceived usefulness. Reviews that have a mixture of

objective, and highly subjective sentences are negatively associated with product sales, compared to reviews that tend to include only subjective or only objective information. However, such reviews are rated more informative (or helpful) by other users. By using Random Forest-based classifiers, we show that we can accurately predict the impact of reviews on sales and their perceived usefulness. We examine the relative importance of the three broad feature categories: “reviewer-related” features, “review subjectivity” features, and “review readability” features, and find that using any of the three feature sets results in a statistically equivalent performance as in the case of using all available features.

### **Rating Prediction Based on Social Sentiment from Textual Reviews**

In recent years, we have witnessed a flourish of review websites. It presents a great opportunity to share our viewpoints for various products we purchase. However, we face an information overloading problem. How to mine valuable information from reviews to understand a user's preferences and make an accurate recommendation is crucial. Traditional recommender systems (RS) consider some factors, such as user's purchase records, product category, and

geographic location. In this work, we propose a sentiment-based rating prediction method (RPS) to improve prediction accuracy in recommender systems. Firstly, we propose a social user sentimental measurement approach and calculate each user's sentiment on items/products. Secondly, we not only consider a user's own sentimental attributes but also take interpersonal sentimental influence into consideration. Then, we consider product reputation, which can be inferred by the sentimental distributions of a user set that reflect customers' comprehensive evaluation. At last, we fuse three factors-user sentiment similarity, interpersonal sentimental influence, and item's reputation similarity-into our recommender system to make an accurate rating prediction. We conduct a performance evaluation of the three sentimental factors on a real-world dataset collected from Yelp.

### **Weakly-Supervised Deep Embedding for Product Review Sentiment Analysis**

Product reviews are valuable for upcoming buyers in helping them make decisions. To this end, different opinion mining techniques have been proposed, where judging a review

sentence's orientation (e.g., positive or negative) is one of their key challenges. Recently, deep learning has emerged as an effective means for solving sentiment classification problems. A neural network intrinsically learns a useful representation automatically without human efforts. However, the success of deep learning highly relies on the availability of large-scale training data. We propose a novel deep learning framework for product review sentiment classification which employs prevalently available ratings as weak supervision signals. The framework consists of two steps: (1) learning a high level representation (an embedding space) which captures the general sentiment distribution of sentences through rating information; and (2) adding a classification layer on top of the embedding layer and use labeled sentences for supervised fine-tuning. We explore two kinds of low level network structure for modeling review sentences, namely, convolutional feature extractors and long short-term memory. To evaluate the proposed framework, we construct a dataset containing 1.1M weakly labeled review sentences and 11,754 labeled review sentences from Amazon

### 3. PROPOSED SYSTEM

Nowadays, there are several websites that allow customers to buy and post reviews of purchased product, which results in incremental accumulation of a lot of reviews written in natural language. Moreover, conversance with E-commerce and social media has raised the level of sophistication of online shoppers and it is common practice for them to compare competing brands of product before making a purchase. Prevailing factors such as availability of online reviews and raised end-user expectations have motivated the development of opinion mining systems that can automatically classify and summarize users' reviews. This work proposes an opinion mining system that can be used for sentiment classifications of user reviews. Feature-based sentiment classification is a multistep process that involves preprocessing to remove noise, extraction of features and corresponding descriptors, and tagging their polarity. The proposed technique extends the feature-based classification approach to incorporate the effect of various linguistic hedges by using ISODATA clustering and implement support vector classification to classify the

reviews from trained opinion words and recommend the product reviews to users.

#### Algorithm

Support Vector Machine(SVM) is a machine learning tool that is based on the idea of large margin data classification. The tool has strong theoretical foundation and the classification algorithms based on it give good generalization performance. Standard implementations, though provide good classification accuracy, are slow and do not scale well. Hence they cannot be applied to large-scale data mining applications. They typically need large number of support vectors. Hence the training as well as the classification times is high.

Input: Input data matrix, class information

Output: Set of Basis vectors

Begin

Repeat

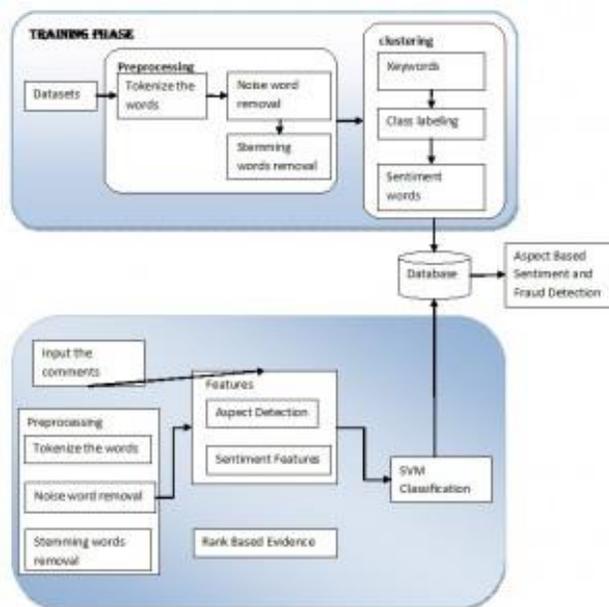
For every candidate example - examples not in current set of BVs

Include it in the model efficiently Observe the generalization performance on the remaining points

end for candidate examples

Add that point to the BVs' list that gave better test error

Till the stopping criterion  
End



**Fig. Online products recommendation system**

#### 4.1 MODULES

##### Review datasets

Opinion is person viewpoint about an object whereas mining is the extraction of knowledge from facts or raw data. Thus, in another word it is a technique which detects intelligent information from data accessible on web. The people who express their opinion on web has dramatically day by day. They can express their opinion almost based on User Generated Content eg review sites,

forums, discussions groups, blogs, products etc. Based on above web site, we can collect user reviews about products.

##### Preprocessing

In this module, we can eliminate stop words and stemming words based on POS tagger. In corpus linguistics, part-of-speech tagging (POS tagging or POST), also called grammatical tagging or word-category disambiguation, is the process of marking up a word in a text (corpus) as corresponding to a particular part of speech, based on both its definition and its context—i.e., its relationship with adjacent and related words in a phrase, sentence, or paragraph. In computing, stop words are words which are filtered out before or after processing of natural language data (text). Though stop words usually refer to the most common words in a language, there is no single universal list of stop words used by all natural language processing tools, and indeed not all tools even use such a list. In computational linguistics, a *stem* is the part of the *word* that never changes even when morphologically inflected, and a lemma is the base form of the *word*. Stemming words are also removed from user reviews. Then implement POS tagger that reads text in some language and assigns parts of speech

to each word (and other token), such as noun, verb, adjective, etc., although generally computational applications use more fine-grained POS tags like 'noun-plural'.

### **Product recommendation**

This recommend the product based on classification by using SVM approach. The input space is planned into a high dimensional feature space. Then, the hyper plane that exploits the margin of separation between classes is constructed. The points that lie closest to the decision surface are called support vectors directly involves its location. When the classes are non-separable, the optimal hyper plane is the one that minimizes the probability of classification error. Initially input image is formulated in feature vectors. Then these feature vectors mapped with the help of kernel function in the feature space. And finally division is computed in the feature space to separate out the classes for training data. A global hyper plane is required by the SVM in order to divide both the program of examples in training set and avoid over fitting. This phenomenon of SVM is higher in comparison to other machine learning techniques which are based on artificial intelligence. Here the important feature for

the classification is the width of the vessels. With the help of SVM classifier we can easily separate out the vessels into arteries and veins. The SVMs demonstrate various attractive features such as good generalization ability compared to other classifiers.

### **Sentiment labeling**

In this module, we only get the keywords and cluster the keywords as positive or negative using ISODATA clustering. ISODATA algorithm is an unsupervised classification [1, 5, 6]. It is similar in principle to the K-means algorithm. However, the ISODATA algorithm determines the number of clusters dynamically. To run the ISODATA algorithm, a lot of parameters such as initial cluster means, splitting parameters, lumping parameters, the minimum number of pixels in a cluster and the number of iterations must be specified. Once these parameters are defined, each sample of the feature space is grouped to the nearest cluster center. The total number of grouped samples in each cluster must meet the minimum required amount. The cluster is eliminated if the minimum amount cannot be reached. After that, compute the mean of the rest grouped samples to update each cluster center.

Finally cluster the positive and Negative reviews.

## **Fraud Detection**

### **1. Ranking Based Evidences**

By analyzing the Apps' historical ranking records, we observe that Apps' ranking behaviors in a leading event always satisfy a specific ranking pattern, which consists of three different ranking phases, namely, rising phase, maintaining phase and recession phase. Specifically, in each leading event, an App's ranking first increases to a peak position in the leader board (i.e., rising phase), then keeps such peak position for a period (i.e., maintaining phase), and finally decreases till the end of the event (i.e., recession phase). Fig. 3 shows an example of different ranking phases of a leading event. Indeed, such a ranking pattern shows an important understanding of leading event. In the following, we formally define the three ranking phases of a leading event.

### **2. Rating Based Evidences**

The ranking based evidences are useful for ranking fraud detection. However, sometimes, it is not sufficient to only use ranking based evidences. For example, some

Apps created by the famous developers, such as Game loft, may have some leading events with large values of  $u1$  due to the developers' credibility and the "word-of-mouth" advertising effect. Moreover, some of the legal marketing services, such as "limited-time discount", may also result in significant ranking based evidences. To solve this issue, we also study how to extract fraud evidences from Apps' historical rating records.

## **CONCLUSION AND FUTURE ENHANCEMENT**

We can conclude that determined the polarity of the customer reviews of product. System performs the product based opinion mining on the given reviews and the feature wise summarized results generated by the system will be helpful for the user in taking the decision. Experimental results indicate that the reviews based Sentiment Orientation System' perform well and has achieved the accuracy. Aspect based opinion mining is necessary because nowadays everyone is busy and they don't have a time to read all the positive or negative reviews if someone just wants to know about some feature of the product. Aspect based opinion mining has proved to be helpful in these situations as

compared to simple opinion mining. Using existing performance benchmarks, the empirical evaluation results show that even when our method does not achieve the best results for all the measures, it does obtain the best precision results, and the F1 results are close to those achieved by the best performance benchmark used in our comparison. Results obtained for the ranking of aspects are also encouraging. In future work, efforts would be done to make some enhancements in this technique in such a way that it can identify the repeated reviews and classify those reviews only once. It would deal with the sentences contain relative clauses like not only-but also and the sentences contain clauses neither-nor, either-or etc.

## REFERENCE

- [1] Duyu Tang, Bing Qin 2015, A Joint Segmentation and Classification Framework for Sentence Level Sentiment Classification.
- [2] Eduard C. Dragut., Kelsey G., (2015). Polarity Consistency Checking for Domain Independent Sentiment Dictionaries. Press Release, November.
- [3] Anindya Ghose., Leilei Z., Daniel K., Dongwon L., (2011). Estimating the Helpfulness and Economic Impact of Product Reviews: Mining Text and Reviewer Characteristics. In Proceedings of AAAI Conference on Artificial Intelligence.
- [4] Xiaojiang Lei, Xueming Qian., Bernardo A., Huberman, (2016). Rating Prediction Based on Social Sentiment From Textual Reviews. Arxiv preprint arXiv:1003.5699.
- [5] Shenghua Liu., Dipanjan D., Kevin G., Noah A. S., (2015). TASC: Topic-Adaptive Sentiment Classification on Dynamic Tweets. Proceedings of the North American Chapter of the association for computational Linguistics Human Language Technologies Conference (NAACL).
- [6] Wei Zhao., Parameswaran A., Petros V., (2018). Weakly-Supervised Deep Embedding for Product Review Sentiment Analysis. In Proceedings of the Third International Conference on Weblogs and Social Media (ICWSM).
- [7] Com S., Korap Ö., (2007). Online consumer-generated reviews have significant impact on offline purchase

behaviour: Türkiye’de Bir Anket Çall mass  
Journal Of Yasar University, 2014 9(35)  
6099-6260.

[8] Joshi M., Nizam H., (2010). Movie reviews and revenues: An experiment in text regression. International Artificial Intelligence and Data Processing Symposium (IDAP'16).

[9] Junting Ye and Leman Akoglu. Discovering opinion spammer groups by network footprints. In Machine Learning and Knowledge Discovery in Databases, pages 267–282. Springer, 2015

[10]. Iker Burguera, Urko Zurutuza, and Simin Nadjm-Tehrani. Crowdroid: Behavior-Based Malware Detection System for Android. In Proceedings of ACM SPSM, pages 15–26. ACM, 2011.

[11] Michael Grace, Yajin Zhou, Qiang Zhang, Shihong Zou, and Xuxian Jiang. Riskranker: Scalable and Accurate Zero-day Android Malware Detection. In Proceedings of ACM MobiSys, 2012.

[12] Hao Peng, Chris Gates, Bhaskar Sarma, Ninghui Li, Yuan Qi, Rahul Potharaju,

Cristina Nita-Rotaru, and Ian Molloy. Using Probabilistic Generative Models for Ranking Risks of Android Apps. In Proceedings of ACM CCS, 2012.

[13] S.Y. Yerima, S. Sezer, and I. Muttik. Android Malware Detection Using Parallel Machine Learning Classifiers. In Proceedings of NGMAST, Sept 2014.

[14] Etsuji Tomita, Akira Tanaka, and Haruhisa Takahashi. The worstcase time complexity for generating all maximal cliques and computational experiments. Theor. Comput. Sci., 363(1):28–42, October 2006

[15] Yajin Zhou and Xuxian Jiang. Dissecting Android Malware: Characterization and Evolution. In Proceedings of the IEEE S&P, pages 95–109. IEEE, 2012.