

# MULTITIER CLASSIFIERS FOR SECURITY OF BIG DATA

<sup>1</sup>G.L.Kale, <sup>2</sup>R.R. Rathod, <sup>3</sup>V.S. Konghe <sup>4</sup>Nilesh N. Shingne,

<sup>1,2,3</sup>Student, <sup>4</sup>Assistant Professor

<sup>1</sup>Department Of Computer Science and Engineering  
Sanmati Engineering College Washim, Maharashtra, India

**Abstract :** Big data achieves a lot of attention from researchers in recent years as a result of it's become in varied application domains. The SVM (Support Vector Machine) and J48 classifiers with base classifiers for rising performance of classification. SVM is higher accuracy and it will turn out powerful leads to vary from wonderful. The planed LIME classifier is giant as a result of it's tailored for handling massive knowledge. during this ensemble classifiers are combined at every tier. Next tier can collect outputs from previous tier, analyses and mix them and send their output to the following tier. Here multitier are used due to several tiers, work is split into every of those tiers in order that speed and accuracy will increase. it's simple to run. It includes totally different ensemble classifiers on many levels, combining strengths of their ways. This classifier is additionally concern for security of huge knowledge. they're generated mechanically as a results of many iterations in applying ensemble Meta classifiers. The ensemble meta classifiers into many tiers at the same time and mix them into one mechanically generated unvaried system in order that several ensemble meta classifiers operate as integral components of alternative ensemble meta classifiers at higher tiers.

**IndexTerms -** Big data , SVM and J48, LIME classifier.

## I. INTRODUCTION

This article introduces five-tier large iterative Multitier Ensemble (LIME) classifiers specifically designed for applications regarding the knowledge security of big knowledge and generate. The most aim of this paper is to develop LIME classifiers as a general technique that will be helpful for the analysis of big data in varied application domains. the technology to extract the knowledge from the pre-existing databases. It's wont to explore and analyze the identical. the information that is to be strip-mined varies from a tiny low data-set to an oversized data-set i.e. Big Data. Big Data is thus massive that it doesn't slot in the most memory of one machine, and it must method massive knowledge by economical algorithms. Trendy computing has entered the age of big data. The investigation of this new construction is very important, as a result of the role of algorithms for analysis of massive knowledge has been growing. It conjointly helps in up security of big data. The most aim of this paper is to develop the classifier as a general technique that will be helpful for the analysis of big data in varied application domains. This construction is illustrated in Figure 1.

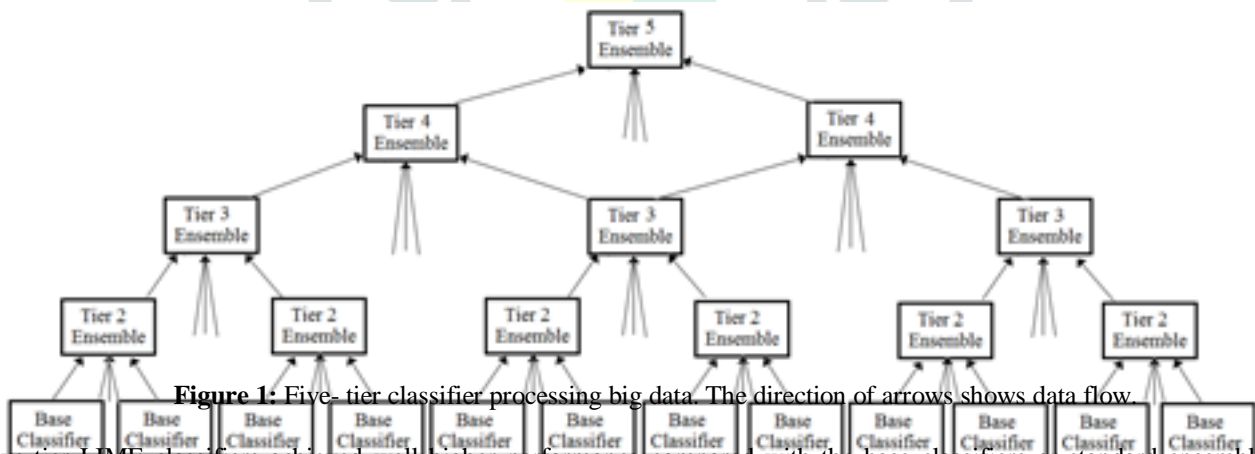


Figure 1: Five-tier classifier processing big data. The direction of arrows shows data flow.

Five-tier LIME classifiers achieved well higher performance compared with the base classifiers or standard ensemble Meta classifiers. This demonstrates that our new technique of combining diverse ensemble Meta classifiers into one unified five-tier ensemble incorporating diverse ensemble Meta classifiers as elements of different ensemble Meta classifiers can be applied to enhance classifications. This paper is organized as follows. Contains brief overview of previous related work. Describes five-tier classifier investigated in this paper. Describes the base classifier which is used in this planned classifier and deals with the ensemble Meta classifiers used in this classifier.

## II. RELATED WORK

This Major security challenges facing the analysis of Big Data and the Cloud have been considered in [11], [13].

Researchers in [1] proposed the four-tier Large Iterative Multitier Ensemble (LIME) classifier which is used for security of the big data. This classifier puts the idea of combining multiple classifiers at several levels.

The paper [3] investigates an iterative hierarchical key exchange scheme for secure scheduling of big data applications in cloud computing. The privacy preservation over big data on cloud is considered in [4].

The first classification method integrating static and dynamic features into a single test was presented in [5]. The approach proposed there improved on previous results using individual features collected separately. The time required for the test was reduced by half.

### III. PROPOSED WORK

We model online user-generated review and overall rating pairs, and aim to identify semantic aspects and aspect-level sentiments from review texts as well as to predict overall sentiments of reviews.

User-generated reviews are different from ordinary text documents. For example, when people read a product re-view, they often care about which specific aspects of the product are commented on, and what sentiment orientations (e.g., positive or negative) have been expressed on the aspects. Instead of employing bag-of-words representation, which is typically adopted for processing usual text documents, we represent each review in an intuitive form of opinion pairs, where each opinion pair consists of an aspect term and related opinion word in the review. Probabilistic topic models, notably latent Dirichlet allocation (LDA) [8], have been widely used for analyzing semantic topical structure of text data. Based on the basic LDA, we introduce an additional aspect-level sentiment identification layer, and construct a probabilistic joint aspect and sentiments framework to model the textual bag-of-opinion-pair data. Online user-generated reviews often come with overall ratings (sentiment labels), which provides us with great flexibility to develop supervised unification topic model. Then, on top of the constructed probabilistic framework, we introduce a new supervised learning layer via normal linear model to jointly model the overall rating data. Thus, we propose a novel supervised joint aspect and sentiment model (SJASM), which can cope with the overall and aspect-based sentiment analysis problems in one go under a unified framework.

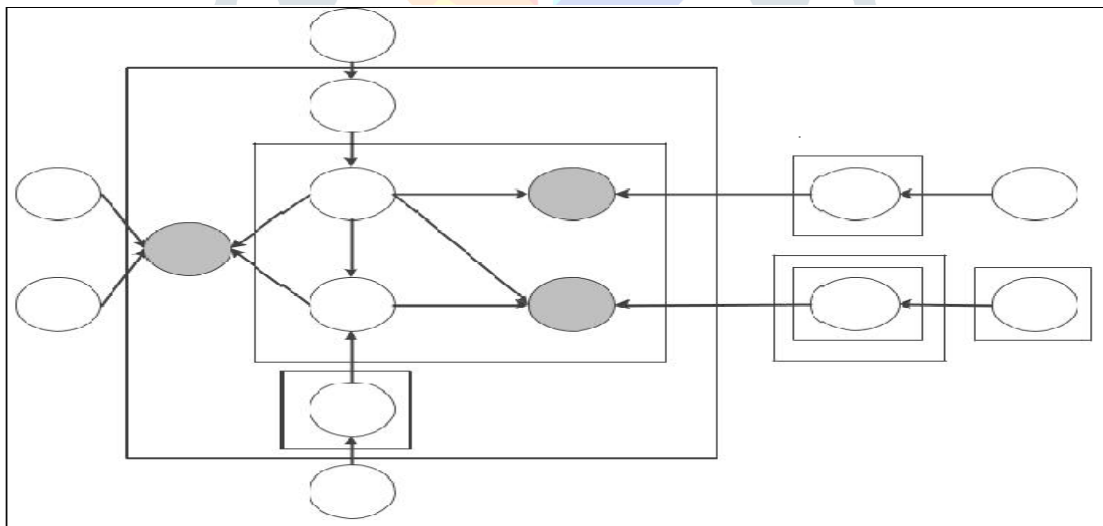
#### 3.1 Supervised Joint Aspect and Sentiment Model

We make the following assumptions about our proposed SJASM model. The generation for aspect-specific sentiments depends on the aspects. This means that we first generate latent aspects, on which we subsequently generate corresponding sentiment orientations.

The generation for aspect terms depends on the aspects, while the generation for opinion words relies on the sentiment orientations and semantic aspects. The formulation is intuitive, for example, to generate an opinion word “beautiful”, we need to know its sentiments orientation positive and related semantic aspect appearance.

The generation for overall ratings of reviews depends on the semantic aspect-level sentiments in the reviews. Based on the model assumptions, to generate a review document and its overall rating, we first draw hidden semantic aspects conditioned on document-specific aspect distribution; We then draw the sentiment orientations on the aspects conditioned on the per document aspect-specific sentiment distribution; Next, we draw each opinion pair, which contains an aspect term and corresponding opinion word, conditioned on aspect and sentiment specific word distributions; We lastly draw the overall rating response based on the generated aspect and sentiment assignments in the review document.

The graphical representation of the proposed model is shown in Figure 1. The notations used in



**Figure 1:** Graphical representation the boxes refer to plates that indicate replicates.

The outer plate refers to review documents, while the inner plate refers to the repeated selection of latent aspects and sentiment orientations as well as random forest aspect terms and user opinion words within each review document.

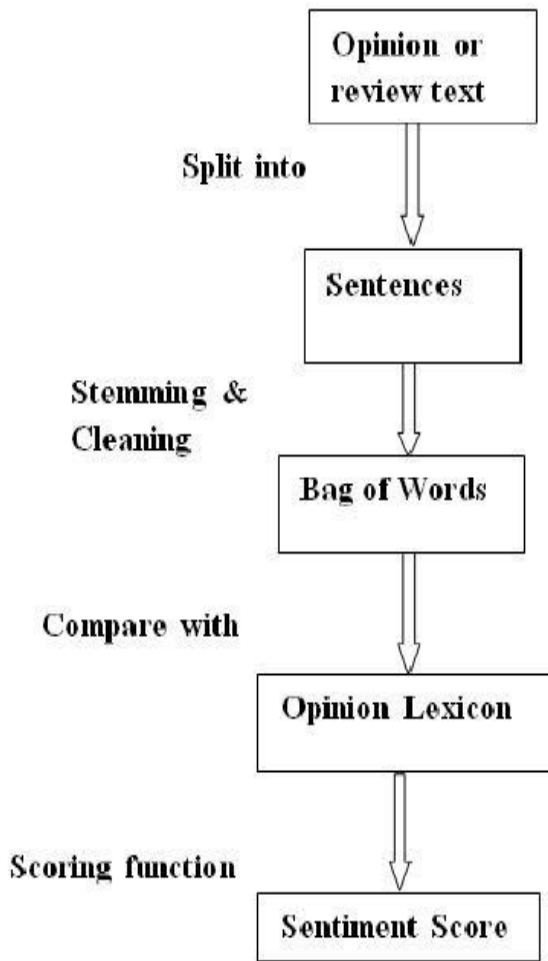
Sentiment analysis and opinion mining are sub fields of machine learning. They are very important in the current development because, lots of user opinionated texts are obtainable in the web now. This is a hard problem to be solved because natural language is highly unstructured in nature. The interpretation of the meaning of a particular sentence by a machine is exasperating. But the usefulness of the sentiment analysis is increasing day by day. Mining public sentiments and analysis of them on E-commerce website data has provided easy way to expose public opinion, which helps for decision making in various domains. E-commerce website is important and popular platforms for people’s interaction. By using E-commerce website platform number of users share their views and opinions. For making important decision it is necessary to mine public opinions and to find reasons behind variation of sentiments is valuable. For example, a company can analyze opinions of public for obtaining users feedback about its products in tweets. In general, opinion mining helps to collect information about the positive and negative aspects of a particular topic. Finally, the positive and highly scored opinions obtained about a particular product are recommended to the user. The main objectives of the work is as below

- To develop a framework that can be used to find users opinions about any product or a person by classifying the tweets into positive or negative polarity.
- To detect the sentiments for tweets obtained in the E-commerce website regarding movie or product.
- To analyze the sentiment which are in non textual formats (Emoticons)
- To enhance the nature of product in accurate way to customer.

**IV. SYSTEM DESIGN**

**4.1 Data Flow Diagram**

This is stage of the project when the theoretical design is turned out in working system. Thus it can be considered to be the most critical stage in achieving a successful new system and in giving the user, confidence that the new system will work and be effective. The framework involves careful planning, investigation of existing system and its constraints on implementation, designing of methods to achieve.



**Figure 2:** Data flow diagram



**Figure 3:** User case diagrams

**4.2 Use-case Diagram**

A use case diagram in the Unified Modeling Language (UML) is a type of behavioral diagram defined by and created from a Use-case analysis. Its purpose is to present a graphical overview of the functionality provided by a system in terms of actors, their goals (represented as use cases), and any dependencies between those use cases. The main purpose of a use case diagram is to show what system functions are performed for which actor. Roles of the actors in the system can be depicted. The purpose of use case diagram is to capture the dynamic aspect of a system. But this definition is too generic to describe the purpose. Use case diagrams are used to gather the requirements of a system including internal and external influences. These requirements are mostly design requirements. So when a system is analyzed to gather its functionalities use cases are prepared and actors are identified. Now when the initial task is complete use case diagrams are modeled to present the outside view. In this case, Use case diagrams are behavior diagrams used to describe a set of actions (use cases) that some system or systems (subject) should or can perform in collaboration with one or more external users of the system (actors). Each use case should provide some observable and valuable result to the actors or other stakeholders of the system. The Actors are E-commerce website, System and user. The actions are E-commerce website dataset (Tweet), Input File, Replacing Emoticons by Aliases, Conversion Hindi to English, POS tagging, Feature Extraction, Sentiment Analysis, Textual Results and Graphical Results

### 4.3 Sequence Diagram.

A Sequence diagram is an interaction diagram that shows how processes operate with one another and in what order. It is a construct of a Message Sequence Chart. A sequence diagram shows object interactions arranged in time sequence. It depicts the objects and classes involved in the scenario and the sequence of messages exchanged between the objects needed to carry out the functionality of the scenario. Sequence diagrams are typically associated with use case realizations in the Logical View of the system under development. Sequence diagrams are sometimes called event diagrams or event scenarios. A sequence diagram shows, as parallel vertical lines (lifelines), different processes or objects that live simultaneously, and, as horizontal arrows, the messages exchanged between them, in the order in which they occur.

A sequence diagram is a type of interaction diagram. In particular it shows the objects participating in the interaction by their lifelines and messages that they are exchanged.

Step 1: Collect reviews from E-commerce Dataset.

Step 2: Store Input in file system (Storage).

Step 3: Process reviews such as

Step 3.1. Replace Emoticons by Aliases

Step 4: harvest and match aspect word of review with dataset predefined aspect

Step 5: Apply Sentimental Analysis as assigning polarity (Very Poor, Poor, Neutral, Good and Very Good).

Step 6: Generate the textual and graphical Results..

Step 7: Exit: User will Exit from system.

## V. CONCLUSION

Current In this work, we focus on modeling online user-generated review data, and aim to identify random forest algorithm in multitier classifier aspects and sentiments on the aspects, as well as to predict over-all ratings of reviews. We have developed a novel supervised joint aspect and web site to deal with the problems in one goes under a unified framework. Multitier treats review documents in the form of user opinion pairs, and can simultaneously model aspect terms and their corresponding words of the reviews for semantic aspect and sentiment detection. Moreover, multitier also leverages overall ratings of reviews as supervision and constraint data, and can jointly infer hidden aspects and sentiments that are not only meaningful but also predictive of overall of the review documents. We conducted experiments using publicly available real-world review data, and extensively compared gender, age, products with seven well-established representative baseline methods. For se-mantic aspect detection and aspect-level sentiment identification problems, LIME classifier outperforms all the generative benchmark models.

## REFERENCES

- [1] Jemal H. Abawajy, Andrei Kelarev, Morshed Chowdhury, "Large Iterative Multitier Ensemble Classifiers for Security of Big Data", in IEEE TRANSACTIONS ON EMERGING TOPICS IN COMPUTING, 30 October 2014.
- [2]. Laura Auria, Rouslan A. Moro, "Support Vector Machines (SVM) as a Technique for Solvency Analysis", in Berlin, August 2008.
3. C. Liu et al. "An iterative hierarchical key exchange scheme for secure scheduling of big data applications in cloud computing," in Proc. 12th IEEE Int. Conf. Trust Security Privacy Comput. Commun. Melbourne, Australia, Jul. 2013, pp. 9-16.
4. R. Islam, J. Abawajy, and M. Warren, "Multi-tier phishing email classification with an impact of classifier rescheduling," in Proc. 10th ISPAN, 2009, pp. 789-793.
5. R. Islam and J. Abawajy, "A multi-tier phishing detection and filtering approach," J. Netw. Comput. Appl., vol. 36, no. 1, pp. 324-335, 2013.
6. Tan, Shulong, et al. "Interpreting the public sentiment variations on E-commerce website." IEEE transactions on knowledge and data engineering 26.5 (2014): 1158-1170.
7. Gautam, Geetika, and DivakarYadav. "Sentiment analysis of E-commerce website data using machine learning approaches and semantic analysis." Contemporary Computing (IC3), 2014 Seventh International Conference on. IEEE, 2014.
8. Jha, Vandana, et al. "HOMS: Hindi opinion mining system." Recent Trends in Information Systems (ReTIS), 2015 IEEE 2nd International Conference on. IEEE, 2015.
9. Larsen, Mark E., et al. "We Feel: mapping emotion on E-commerce website." IEEE journal of biomedical and health informatics 19.4 (2015): 1246-1252.