



e-ISSN: 2319-8753 | p-ISSN: 2347-6710

# IJRSET

International Journal of Innovative Research in  
**SCIENCE | ENGINEERING | TECHNOLOGY**

# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 13, Issue 5, May 2024

**ISSN** INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA

Impact Factor: 8.423

9940 572 462

6381 907 438

[ijirset@gmail.com](mailto:ijirset@gmail.com)

[www.ijirset.com](http://www.ijirset.com)

# Increase the Productivity of Internet Marketing by Using Unsupervised Machine Learning

**Kousar Ghalwade, Dr. Mahesh Ghaikwad, Prof. Amruta Sarudkar**

PG Student, Department of AIDS, Sanjay Ghodwat University, Kolhapur, India

Associate Professor, Department of CSE, Sanjay Ghodawat University, Kolhapur, Maharashtra, India

Assistant Professor, Department of Computer Engineering, Trinity College of Engineering & Research, Pune, India

**ABSTRACT:** Internet marketing is the use of Internet means to increase product sales. This marketing means, with a variety of characteristics such as sincerity, applicability, and facility, is not only an efficient and important marketing tool but also an important core of enterprise product sales. Marketing of e-commerce networks has low cost, good connections, and other advantages in the worldwide traffic. So, the advantage is that it attracts consumers, decreases maintenance, and increases the revenue of internet marketing. So, to increase the productivity of internet marketing we are using unsupervised machine learning by extraction of popular product attributes from E-commerce websites. By using popular product attributes from the product specification pages, we create a graphical mode from a variety of new domains and websites. Without using samples of labeling data and current techniques for extracting information by taking account of the popularity of product characteristics, it allows to mapping of the popular product properties and the related product characteristics from a number of customer reviews. The important thing is that it fills this vocabulary gap between the text on product summary pages and the text on the client's reviews. This paper going to create a graphic model based on hidden random conditions fields, which gives unfair information and apply extensive tests to shows how effective and stable our system is. Our task is to create an unattended method of learning using data mining e-commerce sites. By defining a product, it can classify the review as either positive or negative. It describes the two key concepts in this project, namely features and attribute. The attribute and function relate to the element of a specific domain product. Here the "attribute" in the web-pages defined as a product and the "popular feature" refer to the concepts of the reviews. This model automatically generates associations between the attributes and the popular features in order to identify the popular product attributes.

**KEYWORDS:** Machine Learning, Supervised, Unsupervised ML, preprocessing, Extraction Attributes, Features

## I. INTRODUCTION

The mobile device (application software), therefore, produces more and more regular online business transactions, converts a large number of users from the traditional internet to the mobile Internet. The Internet financial growth itself depends on Internet technology and the financial management. In order to achieve the spectrum of distance control, timely accounting and online payment processing, it is able to handle the electronic trading and the creation of electronic goods settlement papers. This financial model has also become an important factor limiting e-commerce innovations and has become the center of the growth of the e-commerce industry. The referred network is neither the conventional LAN and the WAN, nor the basic Internet, but the Cloud-based open-resource network. Due to constant growth and the rapid popularization of technology, the world's major countries and regions have kept up a fast pace in recent years in the internet marketing. As the e-commerce trend is growing, At the same time, e-commerce companies have expanded and become a strong force in the global e-commerce industry in developing countries. The network funding of the company is therefore the overall result of IT infrastructure and sophisticated internet management. This establishes, applies and innovates in the area of e-commerce financial activities. In this paper the following points are used as the most appropriate technique of web system review:

## INTERNET SYSTEM REVIEW AND ANALYSIS

Based on the IP addresses of communication patterns from the end to end of the world, the Web has become difficult to adapt to the Web as the center of social development. The internet at fixed terminal access has also become difficult. Dynamical protection can be implemented in series or in parallel with high flexibilities based on the "flow" monitoring process, a defense system focused on traffic analysis disturbance detection and a statistical reporting linkage alert control. Experimental results show that the system data gain rates are high and can boost the reliability and protection of the network system operations by connecting each module in an accurate manner and real-time monitoring and testing of the traffic information.

## MOTIVATION

In this project, we are designing an unsupervised learning framework to draw popular product attributes from the product description pages of various websites. This framework can not only detect popular product characteristics from a set of customer reviews, but also map these popular features to the corresponding product characteristics simultaneously as compared to existing systems.

## GOAL

The aim of the proposed system is to build a system which does not experiment the number of popular product attributes in advance. It can be extracted by automatically recognizing hidden concepts, derived from a set of customer reviews and by addressing the vocabulary from the text in the web-based product attributes and the text in the customer review process. This allows consumers to easily view similarities between the products needed and decide quickly to purchase products online.

## BACKGROUND

In today's technology world, data science, machine learning and artificial intelligence is one of the top trend subjects. There are developments in data mining and Bayesian analytics and this raises the need for computer education. Machine learning is a programming discipline to allow them to learn and increase their knowledge automatically. Training in this context includes the identification and appreciation of the data provided and informed decisions based on the data provided. Certain decisions are based on all relevant feedback are very difficult to take into account. In order to address this issue, algorithms are built, that draw on knowledge and experiences by applying statistical science concepts, probability, logic, mathematical optimization, improved learning and control theory.

## MACHINE LEARNING TASKS

Machine Learning (ML) is the scientific study of algorithms and mathematical models using patterns and deductions instead of computer systems for a certain purpose. It is viewed as part of the artificial intelligence. Machine learning algorithms create a mathematical model based on simulated data called training data to take predictions or decisions without being specifically programmed to perform the task. In a wide range of applications, such as e-mail and computer vision, machine learning algorithms are used, where it is impossible or inefficient to create a traditional algorithm to work properly.

Machine learning is closely linked to machine science, based on predicting using machines. The work on mathematical optimization provides the fields of machine learning techniques, theory and implementation. Data mining is an area of research in machine learning and concentrates on the study of exploratory data by non-controlled learning. Machine learning is also referred to as predictive analytics in its application to business problems.

## PURPOSE OF MACHINE LEARNING

The ability to turn the professional knowledge into expertise or identify patterns in complex data is a symbol of human or animal intelligence can be considered as an AI or Artificial Intelligence branch. Machine learning, as a field of science, has similar principles in other fields including statistics, information processing, game theory and optimization. As an area of IT, the aim is to program machines so that they can understand. The aim of machine learning is not however to automate intelligent conduct replication, but to supplement and complement the knowledge and understanding of human beings using the power of computers. Machine learning algorithms, for example, will search and process huge databases identifying patterns outside the scope.

We usually come across many raw data in the real world, which are not eagerly processed through the algorithm of machine learning. Before entering various machine learning algorithms, we need to pre-process the raw data.

### Applications of ML Algorithms

The machine learning algorithms are used in many applications, for example:

- Image processing
- Language processing
- Forecasting such items as stock market movements, economy
- Pattern Identification

- Gaming
- Data analysis
- Robotics systems

**Steps Involved in ML**

The following stages include:

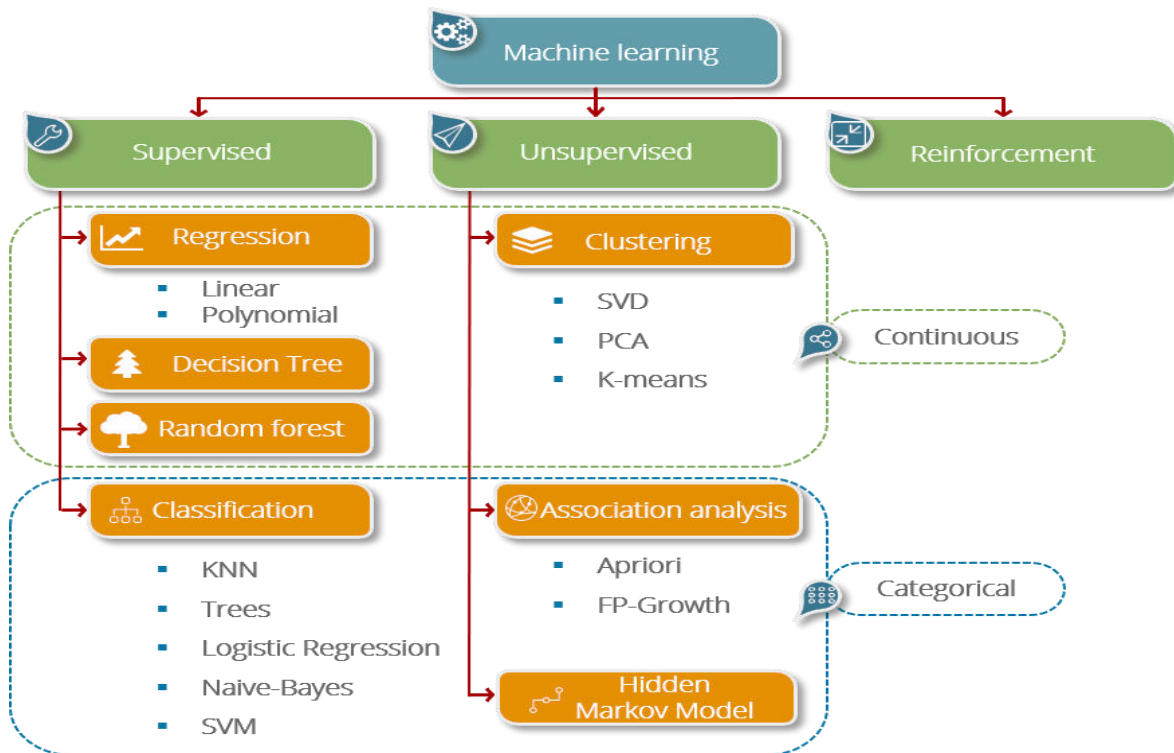
- Problem description
- Software preparation
- algorithms evaluation
- results enhancement
- results presentation

**Machine Learning (ML)** is an automatic learning that requires little or no interaction in humans. It requires programming computers to learn from the inputs that are available. Machine learning aims primarily to explore and develop algorithms that can draw lessons from past data and predict new inputs. The best method to start with Python is to work through an end-to-end project, covering key steps, such as data loading, data synchronization, algorithms evaluation and some predictions. It helps you to use a replicable process after the dataset. Further data can also be applied and performance enhanced.

**CONCEPTS OF LEARNING**

Training is the method by which skill or intelligence is transformed. As described below, ML can be loosely categorized according to the essence of the learning data and the relationship between the learner and the environment in the following categories.

- Supervised learning algorithm
- Unsupervised learning algorithm
- Semi-supervised learning algorithm
- Reinforcement learning algorithm



## Supervised Learning

It is controlled learning is often used in real-world applications like face and voice recognition, products or movies, and sales forecasts. Overseen learning can also be divided into two types: **regression** and **classification**.

- **Regression** trains, for example, the estimation of real estate prices and forecasts a continuous response.
- **Classification** efforts to find an appropriate class mark such as positive / negative feelings, men and women, benign and malignant tumors, healthy and uncertain loans, etc.

Training data is provided in supervised training with desired goals or desired outputs and the goal is to find a general rule that map outputs. This type of data is known as marked data. The learned rule is then used for the detection of new data with uncertain effects. Training under supervision requires the development of a model for machine learning based on examples. For example, when we build a system for estimating the cost of an area or a house, we must first build up a database, and label it, based on various features such as size, location, and so forth. The algorithm needs to be taught what features suit the prices.

On the basis of this knowledge, the algorithm how to measure the real estate price using the input characteristics values. Supervised learning involves learning a task from the training data available. An algorithm analyzes the data of the training and produces a derived function for mapping new instances.

**Several supervised learning algorithms: Logistic Regression, Neural Nets, SVM and Naive Bayes Classifiers.**  
**Supervised learning examples:** Tree of decision, regression, KNN, random forest, logistic regression etc.

### 1.Unsupervised Learning

Unsupervised learning is used for identifying differences, outliers such as fraud and faulty products or to group consumers who conduct a sales operation in a similar way. It's the opposite of supervised learning.

Whether data is used to understand, it is the coder or algorithm's duty to identify the structure of the underlying data and to discover secret patterns or to decide how the data are to be represented, without a description or labels. It is referred as unlabeled data. Unsupervised learning algorithms are powerful tools for data analysis and pattern and trend recognition. Most generally, they are used to organize related inputs into logical classes. **Unsupervised learning algorithms:** Kmeans, Random Forests, Hierarchical clustering.

### 2.Semi-supervised Learning

It is half-controlled learning when some learning samples are labelled but others are not identified. This uses a large number of unmarked training data and a small number of classified test data. Semi-supervised learning occurs in situations where the finding of a completely labelled and functional data set is more costly. It often requires qualified experts for the identification of certain remote sensing, pictures and many field experiments for the location of fuel, for instance, while it is relatively easy to acquire undecided data.

### Reinforcement Learning

The data provided here for learning provides feedback in order for the program to respond to changing situations to achieve a certain goal. Based on the feedback, the machine measures its output and reacts accordingly. The most common examples are self-driving cars and the algorithm for chess master.

## DIFFERENT METHODOLOGIES FOR PYTHON ML PREPROCESSING RESULTS

### 1. Naïve Bayes Classifier Technique

Naive Bayes Classifier is a simple technique for classification. Methods involve classification methods. This is not an algorithm, but a set of algorithms for the training of such classifiers. A classifier for Bayes creates models for problem classification. These classifications are based on the data available. One important aspect of the Bayes classifier is the need to only measure the required parameters for classification by a small quantity of training data. Naive Bayes graders can be very effectively trained in a managed learning environment in certain types of models. Although the naive Bayes classifiers were super-simplified, they worked effectively in many complicated real-world situations.

These have performed well in the detection of spam filters and papers.

## 2. Clustering

Clusters are referred to as classes of similar observations. A common unattended task is the finding of clusters in the training information. Clustering can be characterized also as a way of organizing objects in a particular collection into groups based on certain similar characteristics.

### Applications of Clustering

Applications can be found in many areas, such as market research, model recognition, data analysis and image processing. Assists marketers in discovering different customer groups and in characterizing their customer groups based on shopping patterns.

- In genetics, plant and animal taxonomies can be created, genes with similar characteristics classified and populations can be investigating
- Support for recognition in the earth observation database of areas of the same land use
- Helps to identify web documents for discovery of knowledge
- Software for fraud detection, such as credit card fraud detection
- Cluster analysis gives an insight into the distribution of data to observing properties of each cluster

## 3. Reinforcement Learning

The machine is trained to make certain choices using this algorithm. The algorithm is continuously trained here by means of testing and error methods and feedback methods. The computer learns from past experience and attempts to gain the best possible knowledge to make correct business decisions.

### List of Common ML Algorithms

The commonly used algorithms for machine learning that can apply to virtually any data problem:

- Linear Regression
- Logistic Regression
- Decision Tree
- SVM
- Naive Bayes
- KNN
- K-Means
- Random Forest

## II. LITERATURE SURVEY

1. Innovation of E-commerce Fresh Agricultural Products Marketing Based on Big Internet Data Platform [Lan Li1, Ying zhang Miao ] E-commerce is suggested as the best method for maximizing performance of this program. In the context of the electronic commerce setting, the research is required to carry out the study procedure for the electronic commerce patterns in a new model of electronic commerce that provides the way in which the invention of the electronic commerce pattern is helpful for us, too, to formulate the electronic trade pattern in an enterprise. Experimental results show that the rate of system data acquisition is high and can boost network traffic efficiency and protection in real time, through the connections of each module
2. Quantitative Analysis of the Internal Quality Control and Financing Constraints in Electric Power Enterprises based on SEM model [ Daxiang Suo1, Shuibin Jiang2, Ming Qi1, Li Chen3] : The quantitative research approach of internal quality control and funding limitations for electricity companies based on the SEM model proposes. Internal control and corporate governance are two important aspects of the modern enterprise system. The structure of the internal control system and mechanism of corporate governance is rational and efficient to improve the competitiveness of power companies and the necessary precondition for improved electricity industries, especially the scientific and effective internal control system, improving the company's core strength and enhancing the competitiveness of the company sector. The internal control quality is described in this survey as: integrated internal control capacity assessment and operational effectiveness.
3. Research on the Growth Mechanism of Agricultural Production Enterprises from the Perspective of Resource Dynamic Supply [ Yan Xu1, Hong Chen1, Chang Liu2 ]: Agricultural production is experiencing marked change in the United States as consumer demands, input costs, food safety issues and environmental consequences move rapidly. Multidimensional components and drivers that work in a complex way to influence production

sustainability compose agricultural production systems. In a mixed method approach, we combine quantitative and qualitative data to create and simulate a system dynamic model that examines the structural relationship between these drivers with the economy, climate, and agricultural production society

4. Product aspect ranking using sentiment analysis and TOPSIS [ Saif A. Ahmad Ali Alrababah<sup>1</sup>, Keng Hoon Gan<sup>2</sup>, Tien-Ping Tan<sup>3</sup>]: The rapid expansion of client evaluations on e-commerce sites motivated many scientists to explore the issue of recognizing the product aspects listed in online reviews. Many studies in the study are based on three main criteria: 1) the drawing up of aspects that have been reported on regularly in online reviews, 2) the determination of important aspects that many consumers in their reviews have defined positively and negatively and 3) the relation between the product aspect of the domain (e.g. the camera) and other aspects. This paper uses sentiment analysis and TOPSIS (Technique for Order Quality by Similarity to Ideal Solution) to propose a new product element ranking system in its response to these proposals. The proposed work is divided into two stages: the extraction of aspects and the classification of aspects. In the extraction aspect point, sentiment analysis is applied to understand the product aspects of customer reviews based on the 3 extraction criteria. In the second stage, TOPSIS simultaneously included the product aspects derived from the preceding parameters to generate the most representative product aspects. The analytical assessment of the research proposed using four items online tests indicates its efficacy in defining representative aspects
5. Unsupervised Extraction of Popular Product Attributes from E-Commerce Web Sites by Considering Customer Reviews [ Lidong Bing, Taklam Wong, Wailam]: An unregulated learning system for removing attributes from product description pages was developed at various Web sites for e-commerce. Unlike current extracted information approaches that take into account the popularity of product attributes, our proposed system can detect popular product features from a collection of customer reviews as well as map popular product features. One of our new features is to address a language gap between the product overview page text and the customer feedback site. We are theoretically constructing a segregated graphic model based on secret random conditional fields. As an unsupervised model, it is easy to apply our system to various new domains and websites without requiring marking samples.
6. A Sentiment Classification Model Based On Multiple Classifiers [ Cagatay Catal<sup>1</sup>]: Customer reviews have become a critical factor in consumer decisions through the extensive use of social networks, forums, and blogs. Since the early 2000s, researchers have been focusing on such comments to categorize them automatically as polarities, such as positive, negative and neutral. This problem is known as the description of sentiments. The goal of this study was to investigate and propose a new classification methodology into the potential advantage of the multi-classified system definition for the Turkish sentiment classification. For 3 classifiers, viz. Naive Bayes, the SVM and Bagging, voter algorithms were used. SVM parameters were optimized when utilized as an individual classifier. Experimental results demonstrated that multiple classifier systems increase the efficiency of individual classifiers in Turkish sentiment classifying datasets and that these multiple classifier systems contribute to their strength. The method suggested achieved better efficiency as Naive Bayes, the best individual classification for these datasets, and Support for Vector Machines. A strong sentiments classification method is offered by multiples classifier systems (MCSs), which should be considered when designing MCS-based prediction system for the parameter optimization of individual classifier modes.
7. Extracting Attribute: Value Pairs from Product Specifications on the Web [ Petar Petrovski, Christian Bizer] : The author provided the system of derived value pairs from specifications of products on the internet. The author. Supervised learning is used to separate or not the HTML table and HTML list into a webpage. We again use controlled learning to identify columns as columns of attribute or value column to extract attribute value pairs from the HTML-fragments found in the specification detector. We add several new features for specification detection compared with DEXTER, the current state-of - the-art method for extracting value paired attributes of product specification and support the extraction by pairs of specifications having more than two columns. The developers of the report suit a Bing search engine dataset. Developers use historical knowledge in their approach to build attributes and to fit schemas.
8. Aspect extraction in sentiment analysis- comparative analysis and survey [ Toqir A. Rana<sup>1</sup>, Yu-N Cheah<sup>1</sup>]: Sentiment analytical (SS), due to the expansion of the World Wide Web (WWW), has become one of the most important and increasingly common areas of information and text mining. SA deals with the measurement of the look, thoughts and emotions of users concealed in the text. Extraction of elements is SA's most critical and thoroughly studied process in order to accurately identify sentiments. Such methods were graded according to the method adopted. A thorough comparative analysis is carried out among various approaches of aspects extraction, despite being a traditional survey, which not only elaborates the performance of any technique but also leads readers to compare precision to other state-of - the-art and most recent approaches.
9. A Feature Terms Extraction Process based on Customer User Recommendation Polarity Analysis [Tomofumi Yoshida, Daisuke Kitayama]: Paper provides a method for the extraction of feature conditions that reflect the

feelings over the use of a customer reviews product on websites a recommendation based on content. In view of previous research, which indicates that negative and impressive experiences have more effect than positive ones, we describe the conditions about factors on which consumers disagree in comments about the advantages and disadvantages of sensations relating to product usage. Our methodology consists in extracting sentences from consumer review that express opinions and considers each assessed word as a candidate for product features. Using each candidate's positive opinion, we extract feature words for the chosen product by looking at a feature score based upon a positive evaluation ratio, in order to determine how divided the opinions of the reviewers are. We are presenting an experiment to determine the usefulness of function words extracted using our method.

10. Innovation of E-commerce Fresh Agricultural Products Marketing Based on Big Internet Data Platform [Lan Li<sup>\*</sup>, Ying zhang Miao<sup>2</sup>]: In practice, agricultural production in our country was a family as a small production unit, the rural household often relies not on old prices to choose the project, and decides the size of the production between agricultural products. E-commerce is expected to increase efficiency. In the sense of electronic trade, the research into an examination of electronic business trends is relevant, as the new electronic trade trend is very important to us, because it offers the way for electronic commerce innovation to also be useful in formulating a specific e-commerce pattern Beginning with this impetus and using the large Internet network, this paper proposes the innovative idea of the marketing of electronic commercial fresh agricultural products.

### III. EXISTING SYSTEM

We describe a proposed structure with two important components in this existing system. This first aspect is the common element extraction function that finds to eliminate textual segments from those in the product description web pages. Web pages are considered to be semi-structured text documents with a mixture of standardized content such as HTML tags and free phrases that are either ungrammatical or comprise of short sentences. The buying rate of e-commerce firms is lower. E-commerce company finding is performed on the network with actual benefits compared to the buying activity of other firms. The electronic commerce purchasing enterprise is founded in network contact, negotiation and order completion contracts and will not be limited to other constraints on completion, transaction costs, saving time, and energy by using online banking payment Process to complete the transaction on the network. However, the conventional supply is communicating by telephone, fax and so on, and at last both sides need to meet and discuss communicating, trading, and this is a slow mechanism rather than rapid contact between businesses and electronic commerce firms. The benefits of e-commerce are not only reflected in the remainder of the inventory, they can be monitored on a timely basis through the network itself, effectively reducing inventory costs, reducing product cost and cost. E-commerce enterprises based on network infrastructure use reflects the regional logistical strengths, create a broad and rational logistics system, make it more convenient to transport the product between company and upstream providers and downstream customers, directly on corporate performance, speed up corporate capita activities.

### IV. PROBLEM DEFINATION

Features and attribute, both feature and attribute refer to a characteristic that characterizes a specific domain object. Using "attribute" in the product description to refer to such element and use "popular feature" to refer to the secret aspects / concepts of the feedback that have been uncovered. This model generates associations between the attributes and the popular features automatically to classify popular attributes. Therefore, it is supposed that every feature discovered in the reviews is a popular feature.

Let  $A = \{A_1, A_2, \dots\}$  in a different domain.

The characteristics of the products in this field are set include, for instance, "panel," "multi-media," and the netbook domain product characteristics.

Given a product's web page  $W$  in that domain,

$W$  may be viewed as a tokens sequence (tok 1, ..., tok  $N(W)$ ), with the tokens number referred to in  $N(W)$ .

Tokl,  $k$  is also known as a text fragment consisting of successive tokens, from tokl to tokk in  $W$ , where  $1 \leq \text{tokl} \leq \text{tokk} \leq N(W)$ .

Let the layout features and contents of the text tokl,  $k$  respectively are  $L(\text{tokl}, k)$  and  $C(\text{tokl}, k)$ . We refer to APOP as we are interested in the set of popular product features. Note that APOP is associated with  $C(R)$ , discovered from a selection of customer reviews, common attributes.

Given a web page  $W$  of the company and a collection of customer reviews, we strive to automatically identify all text fragments tokl,  $K$  in  $W$  in order for  $V(\text{tokl}, k) = A_j$  and  $A_j$  APOP, taking into account  $L(\text{tokl}, k)$  and  $C(\text{tokl}, k)$  and common characteristics  $C(R)$ . Currently, we have a number of customer reviews.



Note: APOP is derived automated in advance from C(R) and is not required to be previously pre-specified. In fact, R does not necessarily show the service concerned on page W.

## V. OVERVIEW OF THE FRAMEWORK

In this proposed work, two main components are included. The first element is the common extractor attribute feature to extract text fragments from the product description Web pages that match the typical attributes. Web pages are considered to be a type of semi-structured text document that includes a mixture of structured materials such as HTML tags and free text that can be either ungrammatical or short sentences. Given the unique tokens page in  $W$  ( $tok1, \dots, tokN(W)$ ) within the domain, our objective is to classify all  $tok1, k$  text fragments such that  $V(tok1, k) = A_j$  and  $A_j = tokN(W)$  where APOP is known as the APOP. This role can be described as a problem with sequence labelling. We are marking each token with two sets of labels, specifically in  $(tok1, -tokN(W))$ . The first set includes the "B," "I" and "O" labels that signify an attribute's beginning, both inside and outside an attribute.  $A_j \in APOP$ , that is, the form of popular features, is the second set of labels. To order to address sequence marking issues CRFs have been introduced as the latest model. However, the current standard CRF models are unsatisfactory. First, each token is labelled with two types of labels at the same time, whereas standard CRFs only recognize one type of label. The second reason is because, by the secret principles developed in customer reviews, the common characteristics are related and unknown in advance.

It results in a failure to use supervised instruction in traditional CRFs. We have created a graphic model based on secret CRFs to address this problem. The graphical model proposed can use the hidden concepts derived and the clues from the features of the layout and text contents. There is also an unregulated learning algorithm to remove the common attributes. The second part is designed to extract APOP automatically from a client analysis set (R). This portion produces a set of documents derived from R. A popular feature generates the terms "panel," "resolution," "screen," and so on while a popular feature produces "camera," "speaker" etc. Our graphical model can extract the text fragments associated with common attributes by using such details on words.

## VI. SCOPE

1. One of the features of our system is the lack of knowledge about the number of popular product attributes but can be gained from a set of customer review definitions, which take into account the customers' interest.
2. The capacity to deal with vocabulary differences between Web explanations of texts and customers' feedback or comments are another feature of our system.
3. Objectives
  - To increase the productivity of internet marketing
  - To Extract the Popular Product Attributes from E-Commerce Websites
  - To provide the unsupervised framework for extracting product attributes
  - To extract various features of product from e-commerce side
  - To collect reviews of high quality products from customers
  - To apply algorithm for extracting features of product
  - To do pre-processing and then use classification of dataset

## VII. PROPOSED SOLUTION

### Algorithms

#### A. DOM (Document Objective Model)

First, analyzing the composition of the DOM to decompose a web page into a token ( $tok1, tokN(W)$ ) we perform some basic preprocessing analysis. Specifically, the text background for a page is extracted by crossing the pre-order tree. Therefore, other layout features, such as fonts on each token, are removed from the DOM tree if the current sentence is an entity in the list.

We will cover Document Object Model (DOM) in this article, along with its features and methods for document handling. Introduction: The Data Object Model is the HTML and XML data programming interface. It determines the documents' logical structure and the access and handling of a text.

Note: It is called a logical structure since no relation between objects is specified by DOM.

DOM offers organized hierarchical representation of the website in order to facilitate gliding through the documents for programmers and users. For DOM, tags, IDs, classes, attributes and elements can easily be accessed and manipulated using commands or methods supplied by the object text.

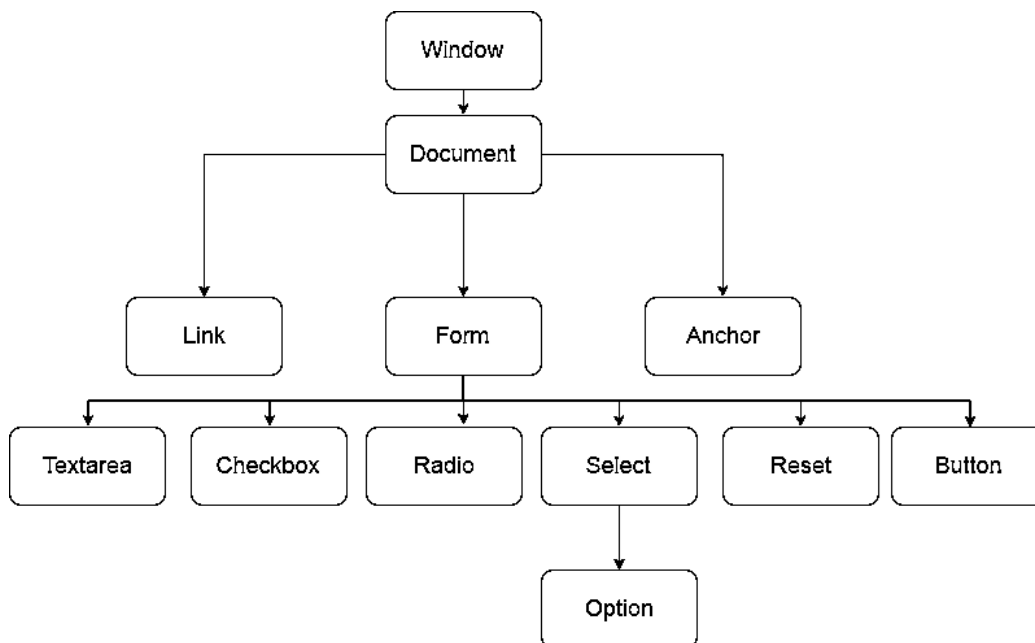
**Structure of DOM**

A Tree or Forest (more than a single tree) can be regarded as the DOM. The treelike representation of a text is commonly used to define the term structure model. Structural isomorphism is an important aspect of DOM structure models: when using two DOM frameworks to create a picture of the same document, precisely the same structure model will be formed with the same objects and relationships.

**Why called as Object Model?**

Documents are represented by artifacts, and the model does not only contain a database structure but also a behavior of a document, which is made up of tags with HTML attributes.

Properties of DOM : Properties of the document object which the document object can access and modify are shown in figure



**Properties of DOM**

- Windows Object: Window Object is always in the top of the hierarchy.
- Document object: When HTML document is loaded into a window, it is turned into an object for documents.
- Form Object: It is indicated by form tags.
- Link Objects: It is indicated by link tags.
- Anchor Objects: It is indicated by a tag.
- Form Control Elements: form can have various control elements including text fields, arrows, radio and check boxes, etc.

**B. Conditional Random Fields (CRFs)**

The product attributes can be extracted from the product overview pages, taking into account the words associated with popular features and web page layout details. In order to address the sequence marking issue, CRFs were adopted as the state-of - the-art standard. Nonetheless, for several reasons, existing standard CRF models are not adequate to accomplish this mission. Second, two kinds of labels label each token at the same time, whereas standard CRFs only

recognize one type of label. The second is that the specific attributes in hidden concepts derived from customer reviews are connected to an unknown by the second variable. It is therefore difficult to use supervised instruction in structured CRFs. We have created a graphic design using secret CRFs to address this problem.

### VIII. CONDITIONAL RANDOM FIELDS

Conditional Random Fields is a type of discriminatory model best suited to prediction tasks in which the current prediction is influenced by background or neighborhood conditions. In the called entities, parts of language marking, genetic analysis, noise reduction and object detection, CRFs consider their application for some names.

This article will first describe Markov Random Fields, the fundamental mathematical and terminology important to the interpretation of the CRF. Then I will present and explain in detail a simple model of random conditions that shows why they are ideal for sequential problems of prediction. I will then discuss the issue of optimizing the probability of occurrence in the context of the CRF model and related derivatives. Ultimately, I will demonstrate the CRF model by placing it for training and inference in the handwriting recognition method.

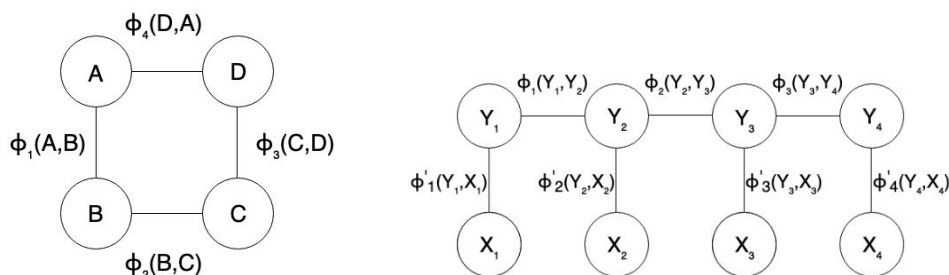
#### A. MARKOV RANDOM FIELDS

A Markov Random Field or Markov Network is a class of undirected graphic models between random versions as shown in figure 3. The graph structure determines whether the random variables are dependent or independent.

##### Markov Network

A Markov Network is shown in Graph  $G = (V, E)$ , the vertices or nodes that define vertices and the nodes that collectively indicate the correlations between these factors.

1. The map can be factorized into various cliques or factors in  $J$ , each of which is dominated by the factor  $\phi_j$ , with a random variables  $D_j$  in its scope. For all  $d_j$  potential values,  $\phi_j(d_j)$  should be strictly positive.
2. All of them should be included with the graph to indicate a subset of random variables as a variable or clique. In



addition, all nodes in the graph should be equal to all cliques ' scope union.

3. The unnormalized joint probability of the variables is the addition of all the factor functions that is for the above-mentioned MRF with  $V = (A, B, C, D)$ , the joint probability can be defined as: -

The denominator is an average of the product of the elements over all possible values that can be calculated by random variables as a standard factor number. It is possible to write the joint likelihood as a normalized factor product.

#### B. CONDITIONAL RANDOM FIELD MODEL

Let us take the Markov Random Field for a while and split it into two different sets of  $Y$  and  $X$  random variables. The Conditional Random Field is a particular Markov Random field where the graph fulfills the property: "If we set the graph for  $X$  globally, i.e. if random  $X$  variables are fixed or given, all random  $Y$  variables follow the Markov  $p(Y_u / X, Y_v, uv) = p(Y_u / X, Y_x, Y_u \sim Y_x)$ , where  $Y_u \sim Y_x$  means  $Y_u$  and  $Y_x$  are in the graph. The chain-structured graph shown below in figure 4 is one of those graphs that satisfies the above property.

##### Conditional Random Field Structure

CRF is, thus, a discriminatory model, i.e. it models  $P(Y / X)$  conditions, i.e.  $X$  is always performed or observed. Figure 5 shows it as,

##### CRF model conditioned on $X$

Because we condition  $X$  and try to find the appropriate  $Y_i$  for every  $X_i$ ,  $X$  and  $Y$  are also referred respectively to as proofs and label variables.

The "factor reduced," as shown in variable  $Y_2$  below, can be checked in accordance with the Markov Principle. As shown, the conditions for  $Y_2$  depends finally only on its neighboring nodes, given all other variables.

The conditional variable  $Y_2$  satisfying the Markov property depends only on the adjacent variables.

## CRF THEORY AND LIKELIHOOD OPTIMIZATION

Let us initially describe the parameters and then construct the joint (and conditional) probabilities equations with the Gibbs notation.

1. Label domain: Suppose that a domain is used for random variables in set  $Y: \{m \in \mathbb{N} \mid 1 \leq m \leq M\}$  i.e. the first  $M$  natural numbers.
2. Domain and structure of evidence: Makes the assumption that in set  $X$  random variables are valued real vector size  $F$  i.e.  $\forall X_i \in X, X_i \in \mathbb{R}^s$ .
3. May the CRF chain length be  $L$  i.e. Variable  $L$  labels and  $L$  proof.
4. Let  $\beta_i(Y_i, Y_j) = W_{cc'}$  if  $Y_i = c, Y_j = c'$  and  $j = i+1, 0$  otherwise.
5. Let  $\beta'_i(Y_i, X_i) = W'c \cdot X_i$ , if  $Y_i = c$  and  $0$  otherwise, where represents the dot product i.e.  $W'c \in \mathbb{R}^s$ .
6. Consider that the total number of parameters is  $M \times M + M \times S$ , i.e. one parameter ( $M \times M$  possible label transitions) and  $S$  parameters for each label ( $M$  potential labels) are to be used for each label transformation. Multiplied to the observation variable of that mark (a vector of size  $S$ ).
7. Let  $D = \{(x_n, y_n)\}$  be the  $N$ -example training data for  $n=1$  to  $N$ .
8. Taking into account the above, the energy and probability can be expressed as follows:-Definition of likelihood for CRF model,

$$\beta'_i(x_i, y_i) = \sum_{c=1}^M \sum_{s=1}^S [y_i = c] W'_{cs} x_{is}$$

$$\beta_i(x_i, y_i) = \sum_{c=1}^M \sum_{c'=1}^M [y_i = c][y_{i+1} = c']$$

$$E(x_n, y_n) = - \sum_{i=1}^L \beta'_i(x_{ni}, y_{ni}) - \sum_{i=1}^{L-1} \beta_i(x_{ni}, y_{ni})$$

$$P(y/x) = \frac{P(y, x)}{P(x)} = \frac{e^{E-(x,y)} / Z}{\sum_{y' \in Y} e^{-E(x,y')}} = \frac{e^{E-(x,y)}}{Z(x)}, \text{ where } Z(x) = \sum_{y' \in Y} e^{-E(x,y')}$$

$$LL(D/\theta) = \frac{1}{N} \sum_{n=1}^N \log(P(y_n/x_n)), \text{ where } \theta \text{ are the model params}$$

## Broyden–Fletcher–Goldfarb–Shanno (BFGS)

The optimal condition can be ideal locally. The solution obtained is influenced by the BFGS algorithm starting point. We may initialize the algorithm with carefully built starting points to improve the quality of the output. In the description of products, we observed that the most popular product value is the noun phrases of "good smell reduction" and "quiet and fast clean fan mode" on the web pages. In order to achieve better performance, the initialization of features associated with a noun phrases is useful. For example, if  $x$  refers to the observation that the underlying token's speech part is a noun, the feature weight  $\mu_k$  for  $g_k(v, y|v, x)$  will be set to a higher value.

## BFGS code

### • First code set

- BFGS: constructor, `dfpmin()`, `lnsrch()`
- `MathFunction4`: interface with `func()`, `df()` [gradient]
- `DemandModelBFGS`: Almost same as in Newton
- `DemandModelBFGSTest`: Almost same as in Newton

### • Second code set

- Same BFGS, `MathFunction4`, `DemandModelBFGSTest`
- `DemandModelBFGS` `df()` method computes gradient Numerically, without analytic expressions for derivatives
- Both BFGS implementations converge quickly on Same logit model coefficients as Newton
- Code is subtle, especially interaction between `dfpmin()` and `lnsrch()`.

## Naive Bayes Classifier Technique

The Naive Bayes Classifier, a basic technique for building classifiers, incorporates classification techniques. It is not an algorithm, but a group of algorithms to train these classifiers. A classifier in Bayes constructs models for the classification of problem cases. The details are used to render these classifications.



An important feature of the naive Bayes classification is that to test the classification parameters, it needs only a small amount of training data. Naive Bayes classifiers can be trained very effectively in some types of models in a supervised learning environment.

Given the overly simplistic assumptions, naive Bayes classifiers have been productive in many complicated situations in the real world. This worked well in the processing of spam filters and papers.

**Multinomial Naive Bayes Classifier**

The combination of probability distribution of P with a fraction of documents per class,

For class **j**, word **i** at a word frequency of **f**:

$$Pr(j) \propto \pi_j \prod_{i=1}^{|V|} Pr(i|j)^{f_i}$$

In order to avoid underflow, we will use the sum of logs:

$$Pr(j) \propto \log(\pi_j \prod_{i=1}^{|V|} Pr(i|j)^{f_i})$$

$$Pr(j) = \log \pi_j + \sum_{i=1}^{|V|} f_i \log(Pr(i|j))$$

One problem is that if a word appears again, the probability that it will appear again will increase.

To smooth this, we take the frequency log:

$$Pr(j) = \log \pi_j + \sum_{i=1}^{|V|} \log(1 + f_i) \log(Pr(i|j))$$

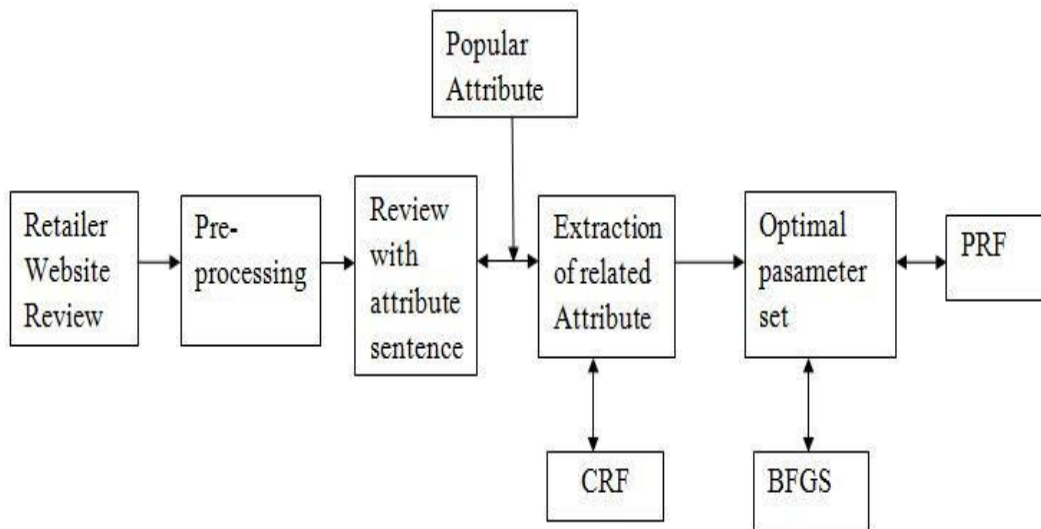
We will also apply an Inverse Document Frequency (IDF) weight to each word in order to take stop

Words into account.

$$t_i = \log\left(\frac{\sum_{n=1}^N doc_n}{doc_i}\right)$$

$$Pr(j) = \log \pi_j + \sum_{i=1}^{|V|} f_i \log(t_i Pr(i|j))$$

IX. SYSTEM ARCHITECTURE



There are several contributions to this work. First, we have developed an unexperienced approach based upon hidden conditional random fields (CRFs) to extract product attributes from the product description pages by considering terms related to popular features as a problem with unknown attributes. The first contribution to this is that we are using the common attribute extraction as an extraction issue with unknown features. The second goal is to close the language gap between the features found in feedback and the features found in the product description. In that context even common attributes can be extracted. Thirdly, on a wide range of product description pages obtained from 13 different fields we have carried out detailed studies. We also compared some current models that can fairly solve this unregulated common extraction problem. The experimental results will demonstrate our framework's efficiency and solidity.

Two major components form this conceptual structure. The first component is the popular extraction attribute, which extracts text fragments from the web pages that are popular attributes. Websites are considered to be a type of semi-structured text documents with a mix of structured contents such as HTML tags and free texts, ungrammatical or composed of only short sentences. Firstly, every token is simultaneously branded with two types of labels, whereas standard CRFs recognize only one type of label. The second reason for this is that the common attributes are related to and unknown in advance by secret principles derived from customer reviews. The second component seeks to extract APOP automatically from a customer feedback set. This portion produces a number of documents derived from R.

In addition to the hints of the interface features, the proposed graphic models will manipulate derived secret concepts. An unregulated learning algorithm for the removal of common attributes is also created. This graphical model can extract fragments of the text relating to the popular attributes by using such information on terms.



X. RESULT ANALYSIS

Product	CRF & BFGS			Multinomial Naive Bayes		
	P	R	F1	P	R	F1
P1	0.692	0.774	0.730	0.73	0.778	0.695
P2	0.570	0.737	0.661	0.695	0.78	0.659
P3	0.641	0.830	0.741	0.662	0.91	0.719
P4	0.789	0.706	0.745	0.798	0.74	0.718
P5	0.641	0.830	0.741	0.67	0.87	0.688
P6	0.570	0.737	0.661	0.65	0.79	0.636
P7						

For analysis on the Amazon E-commerce site we are looking at 7 brands. The common text fragments from explanations can be accurately extracted from our hidden CRF model by using the hidden concept information, the material information and the layout information of each token.

A further comment is that our approach is much more accurate because our hidden CRF model uses the delicate customer reviews, i.e. the derived concepts specified by the customer. While we tested the Multinomial Naive Bayes method using derived concepts to recognize the common attributes, the rich features such as layout characteristics cannot be implemented. Furthermore, our CRF-hidden model is more robust than the Multinomial Naive Bayes ad hoc method, with the sequential labelling on the token set.

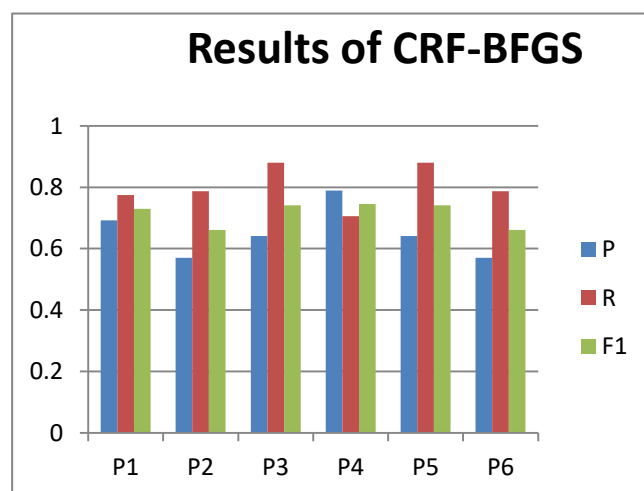
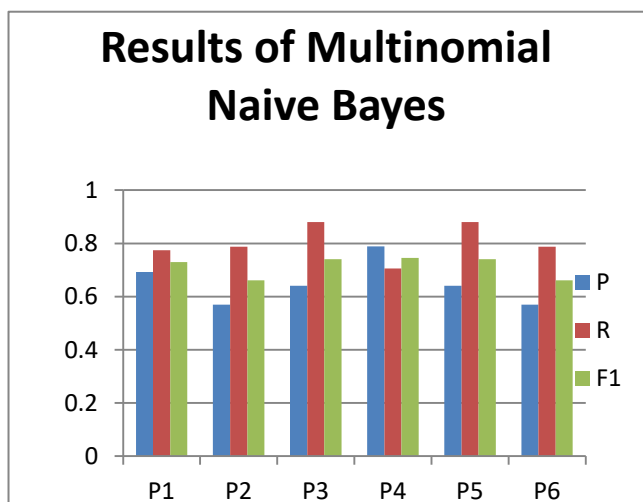


CORPUS & SENTIMENT ANALYSIS

Sr.no	Corpus	Sentiment	Score		
			Positive	Negative	Compound
1	It's friendly to skin very soft and comfortable stuff. Color not fade after wash size is as per expected.	Strongly Positive	0.537	0	0.7003
2	live it. Perfect fit and great looks. <del>Color</del> absolutely <del>Soft</del> and dignified. Good for all occasion both official and informal.	Positive	0.647	0	0.6209333333333333
3	Mostly I buy cloth offline. But this time Amazon rocks. This t-shirts fabric is so soft and fit is very good.	Neutral	0.551	0.7145	0.07692500000000001
4	Awesome product quality. material is thick. Good for both winter and summer.	Positive	0.632	0	0.53265000000000001
5	Good	Positive	1.0	0	0.4404
6	Fitting not as expected.	Neutral	0	0	0
7	<del>Quality</del> quality is <del>Good</del> and never fade after wash	Neutral	0	0	0
8	Products quality was good.	Positive	0.592	0	0.4404
9	Quality is <del>good</del> Mark.... good product.	Positive	0.42	0	0.4404
10	Cloth quality is good.	Positive	0.592	0	0.4404
11	Cloth quality is very good.	Positive	0.592	0	0.4404
12	<del>get</del> so many breaking compliments on this shirt, and i am just like ' <del>chokko</del> ' <del>go</del> . Great shirt nice fit, vibrant color deserves 5	Positive	0.576	0	0.62763

GRAPHICAL RESULT

Results of Multinomial Naive Bayes and CRF - BFGS





## XI. CONCLUSION

This proposal creates new possibilities for companies to develop e-commerce and not just for countries or regions, but also for internationalization, accuracy, computerization, network path. In addition, we are proposing that the e-commerce increase would create new opportunities. Financial management research is more attentive and developable in the new environment. There is therefore a very significant practical value to explore financial management dependent on e-commerce. Electronic commerce can bring together pre-natal agricultural production, production and various post-natal links to a solution, agricultural production and information on the markets, not as symmetrical issue which enables agricultural producers to promptly understand the information on the market and produce fair organization that avoids the price h according to market requirements. The proposed model will address current challenges and provide the world with the new marketing scenario for new agricultural products. In this project, we are able to learn about the strategic climate, strategic goals, strategic material and method of financial management for electronic commerce.

## REFERENCES

1. Akash U. Suryawanshi, P. D. N. K. (2018). Review on Methods of Privacy-Preserving auditing for storing data security in the cloud. *International Journal of Advanced Research in Computer and Communication Engineering (IJARCCE)*, ISSN, 7(4), 247–251.
2. Archana, R. C., Naveenkumar, J., & Patil, S. H. (2011). Iris Image Pre-Processing And Minutiae
3. Points Extraction. *International Journal of Computer Science and Data Security*, 9(6), 171–176.
4. Ayush Khare, D. N. J. (2017). Perspective Analysis Recommendation System in Machine Learning. *International Journal of Emerging Trends & Technology in Computer Science*, 6(2), 184–187.
5. AyushKhare Nitish Bhatt, Dr Naveen Kumar, J. G. (2017). Raspberry Pi Home Automation System Using Mobile App to Control Devices.
6. *International Journal of Innovative Research in Science, Engineering and Technology*, 6(5), 7997– 8003.
7. AyushKhare, J. G., Bhatt, N., & Kumar, N. (2017). Raspberry Pi Home Automation System Using
8. Mobile App to Control Devices. *International Journal of Innovative Research in Science, Engineering and Technology*, 6(5), 7997–8003.
9. R.Udayakumar, R.Elankavi, V.R.Vimal, R.Sugumar, Improved Particle Swarm Optimization With Deep Learning-Based Municipal Solid Waste Management In Smart Cities, *Revista de Gestao Social e Ambiental*, Volume 17, Issue 4, July 2023
10. Bhore, P. R., Joshi, S. D., & Jayakumar, N. (2016). A Survey on the Anomalies in System Design: A Novel Approach. *International Journal of Control Theory and Applications*, 9(44), 443–455.
11. Bhore, P. R., Joshi, S. D., & Jayakumar, N. (2017a). A Stochastic Software Development Process Improvement Model To Identify And Resolve The Anomalies In System Design. *Institute of Integrative Omics and Applied Biotechnology Journal*, 8(2), 154–161.
12. Bhore, P. R., Joshi, S. D., & Jayakumar, N. (2017b). Handling Anomalies in the System Design: A Unique Methodology and Solution. *International Journal of Computer Science Trends and Technology*, 5(2), 409–413.
13. Desai, P. R., & Jayakumar, N. K. (2017). A Survey on Mobile Agents. *International Journal for Research in Applied Science & Engineering Technology (IJRASET)*, 5(XI), 2915–2918.
14. Divyansh Shrivastava Amol K. Kadam, Aarushi Chhibber, Naveenkumar Jayakumar, S. K. (2017). Online Student Feedback Analysis System with Sentiment Analysis. *International Journal of Innovative Research in Science, Engineering and Technology*, 6(5), 8445–8451.
15. Gawade, M. S. S., & Kumar, N. (2016). Three Effective Frameworks for semi-supervised feature selection. *International Journal of Research in Management & Technology*, 6(2), 107–110.
16. Jaiswal, U., Pandey, R., Rana, R., Thakore, D. M., & JayaKumar, N. (2017). Direct Assessment Automator for Outcome Based System. *International Journal of Computer Science Trends and Technology (IJCS T)*, 5(2), 337–340.
17. Jayakumar, D. T., & Naveenkumar, R. (2012). SDjoshi,“. *International Journal of Advanced Research in Computer Science and Software*
18. *Engineering,” Int. J.*, 2(9), 62–70.
19. Jayakumar, N. (2014). Decreases and Discretization Concepts, tools for Predicting Student’s Performance. *Int. J. Eng. Sci. Innov. Technol*, 3(2), 7–15.
20. Jayakumar, N. (2015). Active storage framework leveraging processing capabilities of the embedded storage array.
21. Jayakumar, N., Bhardwaj, T., Pant, K., Joshi, S. D., & Patil, S. H. (2015). A Holistic Approach for Performance Analysis of Embedded Storage Array. *Int. J. Sci. Technol. Eng*, 1(12), 247–250.

22. Jayakumar, N., Iyer, M. S., Joshi, S. D., & Patil, S. H. (2016). A Mathematical Model in Support of Efficient Offloading for Active Storage Architectures. In International Conference on
23. Electronics, Electrical Engineering, Computer Science (ECS) : Innovation and Convergence (Vol. 2, p. 103).
24. Jayakumar, N., & Kulkarni, A. M. (2017). A Simple Measuring Model for Evaluating the Performance of Small Block Size Accesses in Lustre File System. *Engineering, Technology & Applied Science Research*, 7(6), 2313–2318.
25. R.Udayakumar, Suvarna Yogesh, Yogesh Manohar, V.R.Vimal, R.Sugumar, User Activity Analysis Via Network Traffic Using DNN and Optimized Federated Learning based Privacy Preserving Method in Mobile Wireless Networks, *Journal of Wireless Mobile Networks, Ubiquitous Computing, and Dependable Applications*, 14(2), pp. 66–81, June 2023.
26. Jayakumar, N., Singh, S., Patil, S. H., & Joshi, S. D. (2015). Evaluation Parameters of Infrastructure Resources Required for Integrating Parallel Computing Algorithm and Distributed File System. *IJSTE-Int. J. Sci. Technol. Eng.* 1(12), 251–254.
27. Khare, A., & Jayakumar, N. (2017). Perspective
28. Analysis Recommendation System in Machine Learning. *International Journal of Emerging Trends & Technology in Computer Science (IJETTCS)*, 6(2), 184–187.
29. Komalavalli, R., Kumari, P., Navya, S., & Naveenkumar, J. (2017). Reliability Modeling and Analysis of Service-Oriented Architectures. *International Journal of Engineering Science*, 5591.
30. Kumar, N., Angral, S., & Sharma, R. (2014). Integrating Intrusion Detection System with Network Monitoring. *International Journal of Scientific and Research Publications*, 4, 1–4.
31. Kumar, N., Kumar, J., Salunkhe, R. B., & Kadam, A. D. (2016). A Scalable Record Retrieval Methodology Using Relational Keyword Search System. In *Proceedings of the Second International Conference on Data and Communication Technology for Competitive Strategies* (p. 32). ACM.
32. Namdeo, J., & Jayakumar, N. (2014). Predicting Student's Performance Using Data Mining Technique with Rough Set Theory Concepts. *International*
33. *Journal*, 2(2).
34. Osho Tripathi Dr. Naveen Kumar Jayakumar, P. G. (2017). GARDUINO- The Garden Arduino. *International Journal of Computer Science and Technology*, 8(2), 145–147.
35. Singh, A. K., Pati, S. H., & Jayakumar, N. (2017). A Treatment for I/O Latency in I/O Stack. *International Journal of Computer Science Trends and Technology (IJCS T)*, 5(2), 424–427.
36. Yogesh Sawant, P. D. N. Kumar. (2016). Crisp Literature Review One and Scalable Framework: Active Model to Create Synthetic Electrocardiogram Signals. *International Journal of Application or Innovation in Engineering & Management*, 5(11), 73–80.
37. R. Udayakumar, Muhammad Abul Kalam, R. Sugumar, R. Elankavi, Assessing Learning Behaviors Using Gaussian Hybrid Fuzzy Clustering (GHFC) in Special Education Classrooms, *Journal of Wireless Mobile Networks, Ubiquitous Computing, and Dependable Applications (JoWUA)*, Volume 14, Number 1, March 2023, pp. 118-125.
38. Zaeimfar, S. (2014). Workload Characteristics Impacts on File System Benchmarking. *Int. J. Adv.* 39–
39. Gaikwad Mahesh Parasharam , Dr. Harsh Lohiya," Outlier Identification Based on Machine Learning for Medical. *Mathematical Statistician and Engineering Applications*. ISSN:2094-0343 2326-9865
40. Gaikwad Mahesh Parasharam , Dr. Harsh Lohiya,"An empirical Reasearch on the use of Medical Learning Algorithm for Medical Image Classification"ISSN:2394-5125 ,VOL7,ISSUE 10, 2020
41. Gaikwad Mahesh Parasharam , Dr. Harsh Lohiya,"An evaluation of Outlier Detection Using Machine Learning in Medicine" *Mathematical Statistician and Engineering Applications*. ISSN:2094-0343
42. KORE, MS AMRUTA, DR PROF DM THAKORE, and DRAK KADAM. "Innovation of E-commerce Fresh Agricultural Products Marketing Based on Big Internet Data Platform." (2019).
43. Kore, A. A., D. M. Thakore, and A. K. Kadam. "Unsupervised extraction of common product attributes from E-commerce websites by considering client suggestion." *Int. J. Innov. Technol. Explor. Eng.* 8.11 (2019): 1199-1203



INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  [ijirset@gmail.com](mailto:ijirset@gmail.com)



[www.ijirset.com](http://www.ijirset.com)

Scan to save the contact details