

---

# A SURVEY OF MACHINE LEARNING AND DATA MINING TECHNIQUES USED IN MULTIMEDIA SYSTEM

---

A PREPRINT

**Hieu Tran**

Department of Computer Science  
The University of Texas at Dallas  
Richardson, TX 75080  
trunghieu.tran@utdallas.edu

September 23, 2019

## ABSTRACT

Machine learning and data mining are research areas of computer science whose quick development is due to the advances in data analysis research, growth in the database industry and the resulting market needs for methods that are capable of extracting valuable knowledge from large data stores. A vast amount of research work has been done in the multimedia area, targeting different aspects of data analytics, such as the capture, storage, indexing, mining, and retrieval of multimedia big data. However, very few research work provides a complete survey of the whole pine-line of the methods used in machine learning and data mining in the research problems. In this survey paper, we conduct a comprehensive overview of the state-of-the-art methods, algorithms machine learning and data mining in multimedia systems.

**Keywords** Big Data Analytics, Multimedia analysis, Multimedia databases, Indexing, Retrieval, Machine Learning, Data Mining, Mobile Multimedia, Survey, Multimedia System.

## 1 INTRODUCTION

Machine learning is a paradigm that may refer to learning from past experience (which in this case is previous data) to improve future performance. The sole focus of this field is automatic learning methods. Learning refers to modification or improvement of algorithm based on past *experiences* automatically without any external assistance from human. While designing a machine (a software system), the programmer always has a specific purpose in mind. For instance, consider J. K. Rowling's Harry Potter Series and Robert Galbraith's Cormoran Strike Series. To confirm the claim that it was indeed Rowling who had written those books under the name Galbraith, two experts were engaged by The London Sunday Times and using Forensic Machine Learning they were able to prove that the claim was true. They develop a machine learning algorithm and *trained* it with Rowling's as well as other writers writing examples to seek and learn the underlying patterns and then *test* the books by Galbraith. The algorithm concluded that Rowling's and Galbraith's writing matched the most in several aspects. So instead of designing an algorithm to address the problem directly, using Machine Learning, a researcher seek an approach through which the machine, i.e., the algorithm will come up with its own solution based on the example or training data set provided to it initially [3].

The search result performance of the given query is predicted by Query difficulty estimation. It is a powerful tool which is used for multimedia retrieval and now it is becoming more popular. There are several techniques proposed to estimate the query difficulty in the textual information retrieval, but directly they cannot apply for image search, since it will result in poor performance. Existing research on query difficulty estimation focuses on the text-based queries, while the difficulty of image and video retrieval related to multimedia queries has not been yet studied so far. In current years, the prevalence of social media systems, e.g., Flickr, Facebook, and YouTube, has largely increase Internet's multimedia database. These enrich database triggers may leads to the growth of large number of multimedia research scenarios .The success of these social media system also benefits the content based image retrieval [4]. Various

content-based multimedia retrieval methods have been introduced by a large number of researchers. Beyond the methods for content-based image retrieval, audio retrieval and video retrieval, there also have been a wide variety of innovative content-based retrieval methods for new media types, such as content-based retrieval of 3D model, culture artefacts, motion data, biological data, etc. [5]. Data mining can be defined as a method of extracting or "mining" knowledge from large amount of data. It refers to the process of discovering interesting knowledge from huge amounts of data stored in databases, data warehouses, or other information repositories [6]. A multimedia data mining is the process that includes the construction of a multimedia data cube which facilitates multiple dimensional analyses of multimedia data, primarily based on visual content, and the mining of multiple kinds of knowledge, including summarization, comparison, classification, association, and clustering [5]. Thus, in multimedia system the knowledge discovery deals with non-structured information. In order to improve the results of the multimedia files, a database must be first pre-processed, followed by feature extraction. The significant patterns may be discovered with the help of generated features, using various data mining techniques [7].

Multimedia data mining is used for extracting interesting information for multimedia data sets, such as audio, video, images, graphics, speech, text and combination of several types of data set which are all converted from different formats into digital media [8]. Multimedia mining is a subfield of data mining which is used to find interesting information of implicit knowledge from multimedia databases. Multimedia data are classified into five types; they are (i) text data, (ii) Image data (iii) audio data (iv) video data and (v) electronic and digital ink [9]. Text data can be used in web browsers, messages like MMS and SMS. Image data can be used in art work and pictures with text still images taken by a digital camera. Audio data contains sound, MP3 songs, speech and music. Video data include time aligned sequence of frames, MPEG videos from desktops, cell phones, video cameras [10]. Electronic and digital ink its sequence of time aligned 2D or 3D coordinates of stylus, a light pen, data glove sensors, graphical, similar devices are stored in a multimedia database and use to develop a multimedia system.

Since 1960s the research in the field of multimedia has initiated for combining different multimedia data into one application when text and images were combined in a document. During the research and development process of video synchronization of audio and animation was completed using a timeline to specify when they should be played [9]. The difficulties of multimedia data capture, storage, transmission and presentation have been explored in the middle of 1990s where the multimedia standards MPEG-4, X3D, MPEG-7 and MX have continued to grow. These are reformed and clearly handled sound, images, videos, and 3-D (three-dimension) objects that combined by events, synchronization, scripting languages which describe the content of any multimedia object [11]. For multimedia distribution and database applications different algorithms are required. Such a database can be queried, for example, with the SQL multimedia and application packages known as SQL/MM. Multimedia database system includes a multimedia database management system (MMDBMS) which handles and provides foundation for storing, manipulating and retrieving multimedia data from multimedia database [12]. Multimedia data consists of structured data and unstructured data such as audio, video, graphs, images and text media.

## 2 DATA MINING AND MACHINE LEARNING OVERVIEW

### 2.1 Data Mining

Data mining is the extraction of present information from high volume of data sets, it is a modern technology. The main intention of the mining is to extract the information from a large no of data set and convert it into a reasonable structure for further use. The social media websites like Facebook, twitter, instagram enclosed the billions of unrefined raw data. The various techniques in data mining process after analyzing the raw data, new information can be obtained. Since this data is active and unstructured, conventional data mining techniques may not be suitable.

Roughly, Data mining is basically about interpreting any kind of data, but it lays the foundation for machine learning [3]. In practice, it not only sample information from various sources but it analyses and recognises pattern and correlations that exists in those information that would have been difficult to interpret manually. Hence, data mining is not a mere method to prove a hypothesis but method for drawing relevant hypotheses. That mined data and the corresponding patterns and hypotheses may be utilised the basis for machine learning.

The multimedia data mining is classified into two broad categories as static media and dynamic media. Static media contains text (digital library, creating SMS and MMS) and images (photos and medical images). Dynamic media contains Audio (music and MP3 sounds) and Video (movies). Multimedia mining refers to analysis of large amount of multimedia information in order to extract patterns based on their statistical relationships. Figure 1 shows the categories of multimedia data mining

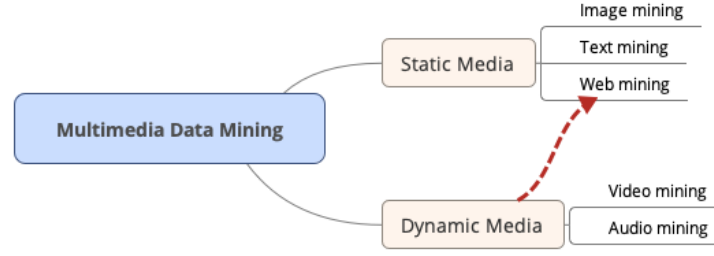


Figure 1: Multimedia Data Mining

## 2.2 Machine Learning

Machine learning takes promote the approach to an advanced level by providing the data essential for a machine to train and modify suitably when exposed to new data. This is known as *training*. It focuses on extracting information from considerably largesets of data, and then detects and identifies underlying patterns using various statistical measures to improve its ability to interpret new data and produce more effective results. Evidently, some parameters should be *tuned* at the incipient level for better productivity.

Machine learning is the foothold of artificial intelligence. It is improbable to design any machine having abilities associated with intelligence, like language or vision, to get there at once. That task would have been almost impossible to solve. Moreover, a system can not be considered completely intelligent if it lacked the ability to learn and improve from its previous exposures.

An overwhelming number of ML algorithm have been designed and introduced over past years. Not everyone of them are widely known. Some of them did not satisfy or solve the problem, so another was introduced in its place. Here the algorithms are broadly grouped into two category and those two groups are further sub-divided. This section try to name most popular ML algorithms and the next section compares three most widely used ML algorithms. Figure 2 shows an overview of machine learning algorithms which are used popularly in data mining.

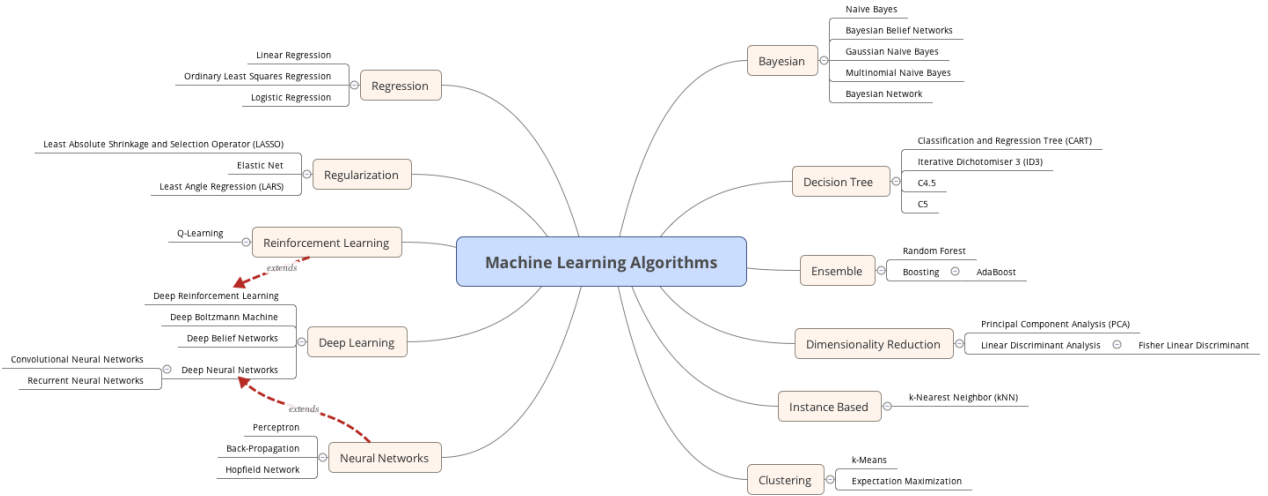


Figure 2: Machine Learning Algorithms used in Multimedia Mining

## 3 MACHINE LEARNING TECHNIQUES USED IN MULTIMEDIA SYSTEM

### 3.1 Regression Algorithms

Regression analysis is part of predictive analytics and exploits the co-relation between dependent (target) and independent variables. The notable regression models are: Linear Regression, Logistic Regression, Stepwise Regression ,

Ordinary Least Squares Regression (OLSR), Multivariate Adaptive Regression Splines (MARS) , Locally Estimated Scatterplot Smoothing (LOESS) etc.

### 3.1.1 Linear Regression model

Simple linear regression is a statistical method that enables users to summarise and study relationships between two continuous (quantitative) variables. Linear regression is a linear model wherein a model that assumes a linear relationship between the input  $variables_x$  and the single output  $variable_y$ . Here the  $y$  can be calculated from a linear combination of the input  $variables_x$ . When there is a single input  $variable_x$ , the method is called a simple linear regression. When there are multiple input variables, the procedure is referred as multiple linear regression. Linear regression is applicable in some of the most popular applications of Linear regression algorithm are in financial portfolio prediction, salary forecasting, real estate predictions and in traffic in arriving at ETAs. Figure 3 shows an example of linear regression model.

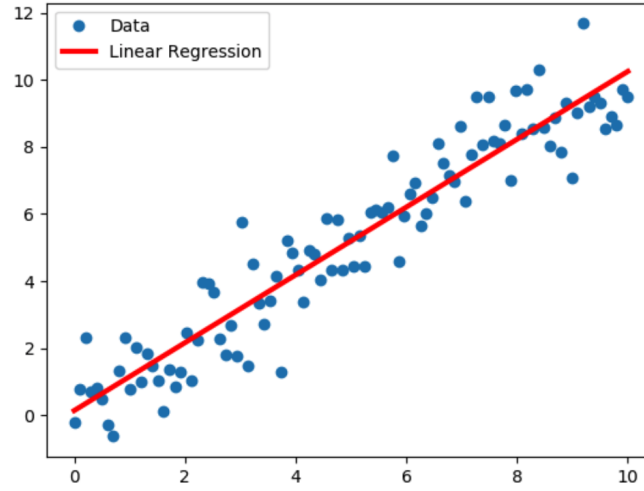


Figure 3: Linear Regression model sample illustration

### 3.1.2 Logistic regression.

One of the most commonly used regression techniques in the industry which is extensively applied across fraud detection, credit card scoring and clinical trials, wherever the response is binary has a major advantage. One of the major upsides of this popular algorithm is that one can include more than one dependent variable which can be continuous or dichotomous. The other major advantage of this supervised machine learning algorithm is that it provides a quantified value to measure the strength of association according to the rest of variables. Despite its popularity, researchers have drawn out its limitations, citing a lack of robust technique and also a great model dependency. Today enterprises deploy Logistic Regression to predict house values in real estate business, customer lifetime value in the insurance sector and are leveraged to produce a continuous outcome such as whether a customer can buy/will buy scenario.

### 3.1.3 Multivariate Regression algorithm.

This technique is used when there is more than one predictor variable in a multivariate regression model and the model is called a multivariate multiple regression. Termed as one of the simplest supervised machine learning algorithms by researchers, this regression algorithm is used to predict the response variable for a set of explanatory variables. This regression technique can be implemented efficiently with the help of matrix operations and in Python, it can be implemented via the "numpy" library which contains definitions and operations for matrix object. Industry application of Multivariate Regression algorithm is seen heavily in the retail sector where customers make a choice on a number of variables such as brand, price and product. The multivariate analysis helps decision makers to find the best combination of factors to increase footfalls in the store.

### 3.1.4 Ordinary least-squares (OLS) regression.

This is a generalized linear modelling technique that may be used to model a single response variable which has been recorded on at least an interval scale. The technique may be applied to single or multiple explanatory variables and also categorical explanatory variables that have been appropriately coded. OLS regression is one of the major techniques used to analyse data and forms the basis of many other techniques [13]. The usefulness of the technique can be greatly extended with the use of dummy variable coding to include grouped explanatory variables [14] and data transformation methods [15]. OLS regression is particularly powerful as it is relatively easy to also check the model assumption such as linearity, constant variance and the effect of outliers using simple graphical methods [16].

## 3.2 Instance-based Algorithms

Instance-based [17] or memory-based learning model stores instances of training data instead of developing an precise definition of target function. Whenever a new problem or example is encountered, it is examined in accordance with the stored instances in order to determine or predict the target function value. It can simply replace a stored instance by a new one if that is a better fit than the former. Due to this, they are also known as winner-take-all method. Examples: K-Nearest Neighbour (KNN), Learning Vector Quantisation (LVQ), Self-Organising Map (SOM), Locally Weighted Learning (LWL) etc.

Instance-based learning generates classification predictions using only specific instances. Instance-based learning algorithms do not maintain a set of abstractions derived from specific instances. This approach extends the nearest neighbor algorithm, which has large storage requirements. It describes how storage requirements can be significantly reduced with, at most, minor sacrifices in learning rate and classification accuracy. While the storage-reducing algorithm performs well on several realworld databases, its performance degrades rapidly with the level of attribute noise in training instances. Therefore, researchers extended it with a significance test to distinguish noisy instances. This extended algorithm's performance degrades gracefully with increasing noise levels and compares favorably with a noise-tolerant decision tree algorithm [17].

### 3.2.1 K-Nearest Neighbour (KNN)

[18] The K-nearest neighbor decision rule has often been used in these pattern recognition problems. One of the difficulties that arises when utilizing this technique is that each of the labeled samples is given equal importance in deciding the class memberships of the pattern to be classified, regardless of their typicalness. Three methods of assigning fuzzy memberships to the labeled samples are proposed, and experimental results and comparisons to the crisp version are presented. The k-Nearest-Neighbors method of classification is one of the simplest methods in machine learning, and is a great way to introduce yourself to machine learning and classification in general. At its most basic level, it is essentially classification by finding the most similar data points in the training data, and making an educated guess based on their classifications. Although very simple to understand and implement, this method has seen wide application in many domains, such as in recommendation systems, semantic searching, and anomaly detection.

KNN Algorithm is based on feature similarity: How closely out-of-sample features resemble our training set determines how we classify a given data point as shown in the Figure 4. KNN can be used for classification - the output is a class membership (predicts a class - a discrete value). An object is classified by a majority vote of its neighbors, with the object being assigned to the class most common among its k nearest neighbors. It can also be used for regression - output is the value for the object (predicts continuous values). This value is the average (or median) of the values of its k nearest neighbors.

## 3.3 Regularisation Algorithm

One of the major aspects of training your machine learning model is avoiding overfitting. The model will have a low accuracy if it is overfitting. This happens because your model is trying too hard to capture the noise in your training dataset. By noise we mean the data points that don't really represent the true properties of your data, but random chance. Learning such data points, makes your model more flexible, at the risk of overfitting. This is a form of regression, that constrains/regularizes or shrinks the coefficient estimates towards zero. In other words, this technique discourages learning a more complex or flexible model, so as to avoid the risk of overfitting.

Regularisation is simply the process of counteracting overfitting or abate the outliers. Regularisation is just a simple yet powerful modification that is augmented with other existing ML models typically Regressive Models. It smoothes up the regression line by castigating any bent of the curve that tries to match the outliers. Examples: Ridge Regression, Least Absolute Shrinkage and Selection Operator (LASSO), Elastic Net, Least-Angle Regression (LARS) etc.

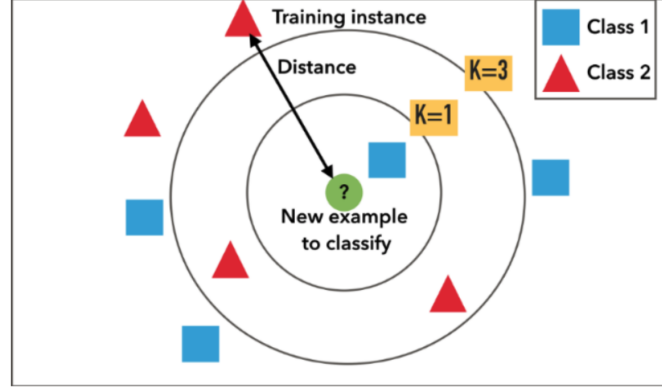


Figure 4: Example of k-NN classification. The test sample (inside circle) should be classified either to the first class of blue squares or to the second class of red triangles. If  $k = 3$  (outside circle) it is assigned to the second class because there are 2 triangles and only 1 square inside the inner circle. If, for example  $k = 5$  it is assigned to the first class (3 squares vs. 2 triangles outside the outer circle).

### 3.4 Decision Tree Algorithms

Decision tree builds classification or regression models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes. A decision node (e.g., Outlook) has two or more branches (e.g., Sunny, Overcast and Rainy). Leaf node (e.g., Play) represents a classification or decision as shown in Figure 5. The topmost decision node in a tree which corresponds to the best predictor called root node. Decision trees can handle both categorical and numerical data.

A decision tree constructs a tree like structure involving of possible solutions to a problem based on certain constraints. It is so named for it begins with a single simple decision or root, which then forks off into a number of branches until a decision or prediction is made, forming a tree. They are favoured for its ability to formalise the problem in hand process that in turn helps identifying potential solutions faster and more accurately than others. Examples: Classification and Regression Tree (CART), Iterative Dichotomiser 3 (ID3), C4.5 and C5.0, Chi-squared Automatic Interaction Detection (CHAID), Decision Stump, M5, Conditional Decision Trees etc.

Decision tree algorithm is a kind of data mining model to make induction learning algorithm based on examples [19]. It is easy to extract display rule, has smaller computation amount, and could display important decision property and own higher classification precision. For the study of data mining algorithm based on decision tree, this article put forward specific solution for the problems of property value vacancy, multiple-valued property selection, property selection criteria, propose to introduce weighted and simplified entropy into decision tree algorithm so as to achieve the improvement of ID3 algorithm. The experimental results show that the improved algorithm is better than widely used ID3 algorithm at present on overall performance.

### 3.5 Bayesian Algorithms

A group of ML algorithms employ Bayes Theorem to solve classification and regression problems. Examples: Naive Bayes, Gaussian Naive Bayes, Multinomial Naive Bayes, Averaged One-Dependence Estimators (AODE), Bayesian Belief Network (BBN), Bayesian Network (BN) etc.

#### 3.5.1 Naive Bayes.

The complexity of the above Bayesian classifier needs to be reduced, for it to be practical. The naive Bayes algorithm does that by making an assumption of conditional independence over the training dataset. This drastically reduces the complexity of above mentioned problem to just  $2n$ . The assumption of conditional independence states that, given random variables  $X$ ,  $Y$  and  $Z$ , we say  $X$  is conditionally independent of  $Y$  given  $Z$ , if and only if the probability distribution governing  $X$  is independent of the value of  $Y$  given  $Z$ . In other words,  $X$  and  $Y$  are conditionally independent given  $Z$  if and only if, given knowledge that  $Z$  occurs, knowledge of whether  $X$  occurs provides no information on the likelihood of  $Y$  occurring, and knowledge of whether  $Y$  occurs provides no information on the likelihood of  $X$  occurring. Assume that  $A$  is the response variable and  $B$  is the input attribute. So according to the equation, we have:



Figure 5: Decision Tree sample

$P(A|B)$  : conditional probability of response variable belonging to a particular value, given the input attributes. This is also known as the posterior probability.  $P(A)$  : The prior probability of the response variable.  $P(B)$  : The probability of training data or the evidence.  $P(B|A)$  : This is known as the likelihood of the training data. Therefore, the above equation can be rewritten as  $Posterior = \frac{Prior \times Likelihood}{Evidence}$

Naive Bayes classifiers work really well in complex situations, despite the simplified assumptions and naivety. The advantage of these classifiers is that they require small number of training data for estimating the parameters necessary for classification. This is the algorithm of choice for text categorization. This is the basic idea behind naive Bayes classifiers, that you need to start experimenting with the algorithm.

### 3.5.2 Bayesian Networks.

Bayesian networks (BNs), also known as belief networks (or Bayes nets for short), belong to the family of probabilistic graphical models (GMs). These graphical structures are used to represent knowledge about an uncertain domain. In particular, each node in the graph represents a random variable, while the edges between the nodes represent probabilistic dependencies among the corresponding random variables. These conditional dependencies in the graph are often estimated by using known statistical and computational methods. Hence, BNs combine principles from graph theory, probability theory, computer science, and statistics. GMs with undirected edges are generally called Markov random fields or Markov networks. These networks provide a simple definition of independence between any two distinct nodes based on the concept of a Markov blanket. Markov networks are popular in fields such as statistical physics and computer vision [20] [21].

BNs became extremely popular models in the last decade. They have been used for applications in various areas, such as machine learning, text mining, natural language processing, speech recognition, signal processing, bioinformatics, error-control codes, medical diagnosis, weather forecasting, and cellular networks. The name BNs might be misleading. Although the use of Bayesian statistics in conjunction with BN provides an efficient approach for avoiding data overfitting, the use of BN models does not necessarily imply a commitment to Bayesian statistics. In fact, practitioners often follow frequentists' methods to estimate the parameters of the BN. On the other hand, in a general form of the graph, the nodes can represent not only random variables but also hypotheses, beliefs, and latent variables [22]. Such a structure is intuitively appealing and convenient for the representation of both causal and probabilistic semantics. As indicated by David [23], this structure is ideal for combining prior knowledge, which often comes in causal form, and observed data. BN can be used, even in the case of missing data, to learn the causal relationships and gain an understanding of the various problem domains and to predict future events.

### 3.6 Support Vector Machine

Support Vector Machine (SVM) is another most powerful algorithm with strong theoretical foundations based on the Vapnik-Chervonenkis theory, as defined by Oracle docs. This supervised machine learning algorithm has strong regularization and can be leveraged both for classification or regression challenges. They are characterized by usage of kernels, the sparseness of the solution and the capacity control gained by acting on the margin, or on number of support vectors, etc. The capacity of the system is controlled by parameters that do not depend on the dimensionality of feature space. Since the SVM algorithm operates natively on numeric attributes, it uses a z-score normalization on numeric

attributes. In regression, Support Vector Machines algorithms use epsilon-insensitivity (margin of tolerance) loss function to solve regression problems. Support vector machines regression algorithms has found several applications in the oil and gas industry, classification of images and text and hypertext categorization. In the oilfields, it is specifically leveraged for exploration to understand the position of layers of rocks and create 2D and 3D models as a representation of the subsoil.

SVM is so popular a ML technique that it can be a group of its own. It uses a separating hyperplane or a decision plane to demarcate decision boundaries among a set of data points classified with different labels. It is a strictly supervised classification algorithm. In other words, the algorithm develops an optimal hyperplane utilising input data or training data and this decision plane in turns categories new examples. Based on the kernel in use, SVM can perform both linear and nonlinear classification.

### 3.7 Clustering Algorithms

Clustering is a Machine Learning technique that involves the grouping of data points. Given a set of data points, we can use a clustering algorithm to classify each data point into a specific group. In theory, data points that are in the same group should have similar properties and/or features, while data points in different groups should have highly dissimilar properties and/or features. Clustering is a method of unsupervised learning and is a common technique for statistical data analysis used in many fields.

Clustering is concerned with using ingrained pattern in datasets to classify and label the data accordingly. Examples: K-Means, K-Medians, Affinity Propagation, Spectral Clustering, Ward hierarchical clustering, Agglomerative clustering, DBSCAN, Gaussian Mixtures, Birch, Mean Shift, Expectation Maximisation (EM) etc.

#### 3.7.1 K-Means Clustering.

K-means clustering is an extensively used technique for data cluster analysis. To process the learning data, the K-means algorithm in data mining starts with a first group of randomly selected centroids, which are used as the beginning points for every cluster, and then performs iterative (repetitive) calculations to optimize the positions of the centroids. It halts creating and optimizing clusters when either: 1) The centroids have stabilized - there is no change in their values because the clustering has been successful. 2) The defined number of iterations has been achieved. The demonstration of the algorithm is described as in Figure 6.

K-Means has the advantage that it's pretty fast, as all we are really doing is computing the distances between points and group centers; very few computations. It thus has a linear complexity  $O(n)$ . On the other hand, K-Means has a couple of disadvantages. Firstly, you have to select how many groups/classes there are. This is not always trivial and ideally with a clustering algorithm we would want it to figure those out for us because the point of it is to gain some insight from the data. K-means also starts with a random choice of cluster centers and therefore it may yield different clustering results on different runs of the algorithm. Thus, the results may not be repeatable and lack consistency. Other cluster methods are more consistent.

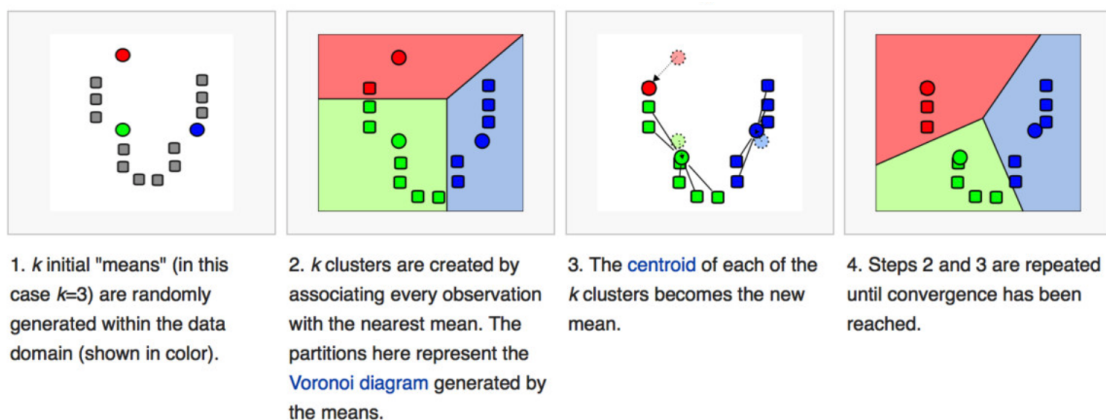


Figure 6: Demonstration of K-means algorithms



### 3.7.2 K-Medians.

This is another clustering algorithm related to K-Means, except instead of recomputing the group center points using the mean we use the median vector of the group. This method is less sensitive to outliers (because of using the Median) but is much slower for larger datasets as sorting is required on each iteration when computing the Median vector.

## 3.8 Association Rule Learning Algorithms

Most Machine learning algorithms in data science work with numeric data and tend to be quite mathematical. However, association rule is perfect for categorical data and involves little more than simple counting. It is a rule-based ML method for discovering interesting relations between variables in large databases. Identifies strong rules using some measures of interestingness.

Association rules help discover correlation between apparently unassociated data. They are widely used by ecommerce websites to predict customer behaviours and future needs to promote certain appealing products to him. Examples: Apriori algorithm, Eclat algorithm etc.

### 3.8.1 Apriori algorithm.

Apriori is an algorithm for frequent item set mining and association rule learning over transactional databases. It proceeds by identifying the frequent individual items in the database and extending them to larger and larger item sets as long as those item sets appear sufficiently often in the database. The frequent item sets determined by Apriori can be used to determine association rules which highlight general trends in the database: this has applications in domains such as market basket analysis. It uses a breadth-first search strategy to count the support of itemsets and uses a candidate generation function which exploits the downward closure property of support.

There are several advantages of association rule in comparison with others machine learning algorithms such as least memory consumption, easy implementation. Using Apriori property for pruning, therefore, itemsets left for further support checking remain less. However, there exist some drawbacks because the algorithm requires many scans of database, it allows only single minimum support threshold and it is favourable for small databases.

## 3.9 Artificial Neural Network (ANN) Algorithms

### 3.9.1 Neural Network.

An Artificial Neural Network (ANN) is an information processing paradigm that is inspired by the way biological nervous systems, such as the brain, process information. The key element of this paradigm is the novel structure of the information processing system. It is composed of a large number of highly interconnected processing elements (neurons) working in unison to solve specific problems. ANNs, like people, learn by example. An ANN is configured for a specific application, such as pattern recognition or data classification, through a learning process. Learning in biological systems involves adjustments to the synaptic connections that exist between the neurons. Neural networks, with their remarkable ability to derive meaning from complicated or imprecise data, can be used to extract patterns and detect trends that are too complex to be noticed by either humans or other computer techniques. A trained neural network can be thought of as an "expert" in the category of information it has been given to analyse. This expert can then be used to provide projections given new situations of interest and answer "what if" questions. Other advantages include:

*Adaptive learning:* An ability to learn how to do tasks based on the data given for training or initial experience.

*Self-Organisation:* An ANN can create its own organisation or representation of the information it receives during learning time.

*Real Time Operation:* ANN computations may be carried out in parallel, and special hardware devices are being designed and manufactured which take advantage of this capability.

*Fault Tolerance via Redundant Information Coding:* Partial destruction of a network leads to the corresponding degradation of performance. However, some network capabilities may be retained even with major network damage.

A model based on the built and operations of actual neural networks of humans or animals. ANNs are regarded as non-linear models as it tries to discover complex associations between input and output data. But it draws sample from data rather than considering the entire set and thereby reducing cost and time. Examples: Perceptron, Back-Propagation, Hop-field Network, Radial Basis Function Network (RBFN) etc. The Figure 7 shows the Artificial Neural Network and Deep Learning Neural Network.

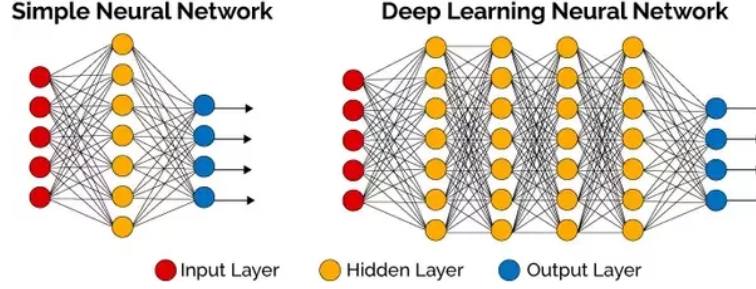


Figure 7: Artificial Neural Network and Deep Learning Neural Network

### 3.9.2 Deep learning algorithms.

Deep learning algorithms are more modernised versions of ANNs that capitalise on the profuse supply of data today. They utilise larger neural networks to solve semi-supervised problems where major portion of an abundant data is unlabelled or not classified. Examples: Deep Boltzmann Machine (DBM), Deep Belief Networks (DBN), Convolutional Neural Network (CNN), Stacked Auto-Encoders etc.

### 3.10 Dimensionality Reduction Algorithms

Dimensionality reduction is typically employed to reduce a larger data set to its most discriminative components to contain relevant information and describe it with fewer features. This gives a proper visualisation for data with numerous features or of high dimensionality and helps in implementing supervised classification more efficiently. Examples: Principal Component Analysis (PCA), Principal Component Regression (PCR), Partial Least Squares Regression (PLSR), Sammon Mapping, Multidimensional Scaling (MDS), Projection Pursuit, Linear Discriminant Analysis (LDA), Mixture Discriminant Analysis (MDA), Quadratic Discriminant Analysis (QDA), Flexible Discriminant Analysis (FDA) etc.

Dimensionality reduction techniques are being applied in many scientific areas ranging from biomedical research to text mining and computer science. In this review we have covered different families of methodologies; each of them based on different criteria but all chasing the same goal: reduce the complexity of the data structure while at the same time delivering a more understandable representation of the same information. The field is still very active and ever more powerful methods are continuously appearing providing an excellent application test bed for applied mathematicians. The Figure 8 below shows an illustration of dimension reduction techniques.

### 3.11 Ensemble Algorithms

Ensemble methods are meta-algorithms that combine several machine learning techniques into one predictive model in order to decrease variance (bagging), bias (boosting), or improve predictions (stacking) [26]. Ensemble methods can be divided into two groups:

1. *sequential* ensemble methods where the base learners are generated sequentially (e.g. AdaBoost). The basic motivation of sequential methods is to exploit the dependence between the base learners. The overall performance can be boosted by weighing previously mislabeled examples with higher weight.
2. *parallel* ensemble methods where the base learners are generated in parallel (e.g. Random Forest). The basic motivation of parallel methods is to exploit independence between the base learners since the error can be reduced dramatically by averaging.

Most ensemble methods use a single base learning algorithm to produce homogeneous base learners, i.e. learners of the same type, leading to homogeneous ensembles. There are also some methods that use heterogeneous learners, i.e. learners of different types, leading to heterogeneous ensembles. In order for ensemble methods to be more accurate than any of its individual members, the base learners have to be as accurate as possible and as diverse as possible. The main purpose of an ensemble method is to integrate the projections of several weaker estimators that are singly trained in order to boost up or enhance generalisability or robustness over a single estimator. The types of learners and the means to incorporate them is carefully chosen as to maximise the accuracy [24] [25]. Examples: Boosting, Bootstrapped Aggregation (Bagging), AdaBoost, Stacked Generalisation (blending), Gradient Boosting Machines (GBM), Gradient Boosted Regression Trees (GBRT), Random Forest, Extremely Randomised Trees etc. Ensembling contains 2 main techniques as shown in the Figure 9.

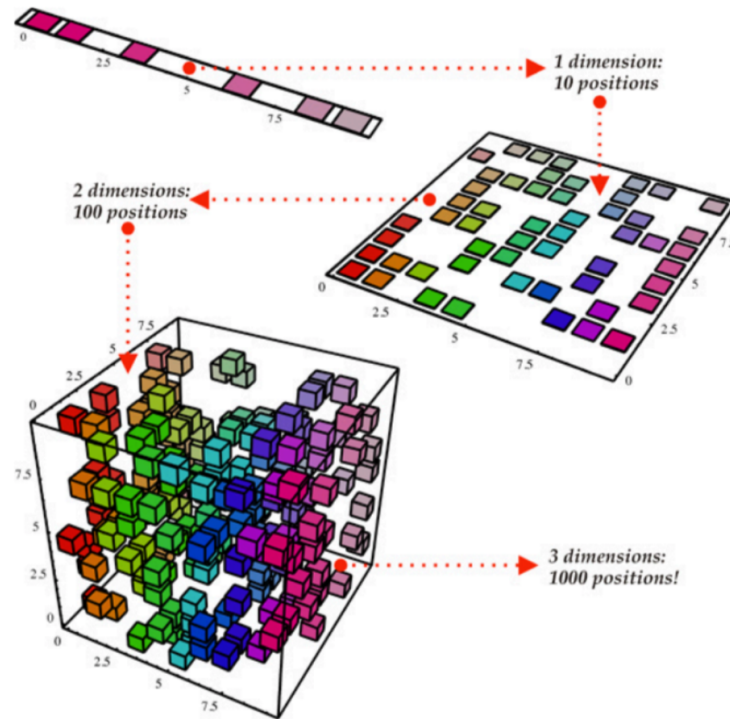


Figure 8: Dimension Reduction illustration

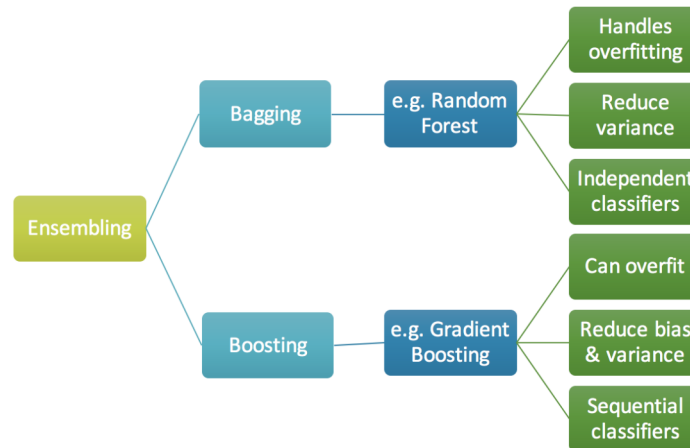


Figure 9: Ensembling

### 3.11.1 Random Forest.

Random Forest is a flexible, easy to use machine learning algorithm that produces, even without hyper-parameter tuning, a great result most of the time. It is also one of the most used algorithms, because of its simplicity and the fact that it can be used for both classification and regression tasks. In this post, you are going to learn, how the random forest algorithm works and several other important things about it [27] [28]. Random Forest is a supervised learning algorithm. It creates a forest and makes it somehow random. The forest it builds, is an ensemble of Decision Trees, most of the time trained with the bagging method. The general idea of the bagging method is that a combination of learning models increases the overall result [29]. One big advantage of random forest is, that it can be used for both classification and regression problems, which form the majority of current machine learning systems. I will talk about

random forest in classification, since classification is sometimes considered the building block of machine learning. Below you can see how a random forest would look like as Figure 10.

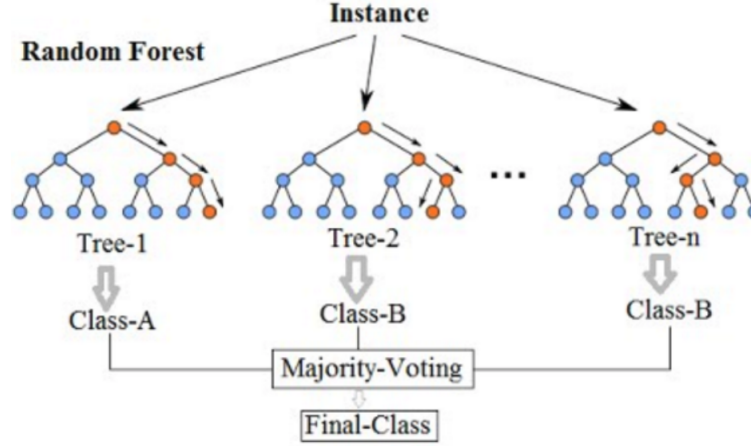


Figure 10: Random Forest Illustration

Random Forest is also considered as a very handy and easy to use algorithm, because its default hyperparameters often produce a good prediction result. The number of hyperparameters is also not that high and they are straightforward to understand. One of the big problems in machine learning is overfitting, but most of the time this will not happen that easy to a random forest classifier. That is because if there are enough trees in the forest, the classifier will not overfit the model. The main limitation of Random Forest is that a large number of trees can make the algorithm to slow and ineffective for real-time predictions. In general, these algorithms are fast to train, but quite slow to create predictions once they are trained. A more accurate prediction requires more trees, which results in a slower model. In most real-world applications the random forest algorithm is fast enough, but there can certainly be situations where run-time performance is important and other approaches would be preferred. And of course Random Forest is a predictive modeling tool and not a descriptive tool. That means, if you are looking for a description of the relationships in your data, other approaches would be preferred.

### 3.11.2 Gradient Boosting Machines.

Gradient Boosting Machines are a family of powerful machine-learning techniques that have shown considerable success in a wide range of practical applications. They are highly customizable to the particular needs of the application, like being learned with respect to different loss functions. The principle idea behind this algorithm is to construct the new base-learners to be maximally correlated with the negative gradient of the loss function, associated with the whole ensemble. The loss functions applied can be arbitrary, but to give a better intuition, if the error function is the classic squared-error loss, the learning procedure would result in consecutive error-fitting. In general, the choice of the loss function is up to the researcher, with both a rich variety of loss functions derived so far and with the possibility of implementing one's own task-specific loss.

This high flexibility makes the GBMs highly customizable to any particular data-driven task [32]. It introduces a lot of freedom into the model design thus making the choice of the most appropriate loss function a matter of trial and error [30] [31]. However, boosting algorithms are relatively simple to implement, which allows one to experiment with different model designs. Moreover the GBMs have shown considerable success in not only practical applications, but also in various machine-learning and data-mining challenges [33]. The capabilities of the GBMs were investigated on a set of real-world practical applications. In every case, GBMs provided excellent results in terms of accuracy and generalization [36] [37]. In addition, the GBMs offered additional insights into the resulting model design, allowing for deeper investigation and analysis of the modeled effects [34] [35].

## 4 APPLICATIONS OF TECHNIQUES IN REAL DATA MINING

### 4.1 Mining of Multimedia database

Multimedia data mining is the method of discovering, interesting patterns from multimedia databases that store and manage large collections of multimedia objects, such as video data, audio data, image data, sequence data, and hypertext data which includes text, text markups, and linkages [6]. Some of the important issues in multimedia data mining include similarity search, content based retrieval, and multidimensional analysis. In mining of multimedia process, data collection is the first step of a learning system, as the overall achievable performance depends upon the quality of raw data. The next steps are the data preprocessing which is responsible to discover important features from raw data. It includes data cleaning, normalization, transformation, feature selection, etc. Learning largely depends on the informative features identified at pre-processing stage. The output of data pre-processing is the training set. Thus a training set has to select a learning model to learn from it and make multimedia mining model more iterative [38]. The Figure 11 below shows the multimedia mining processing:

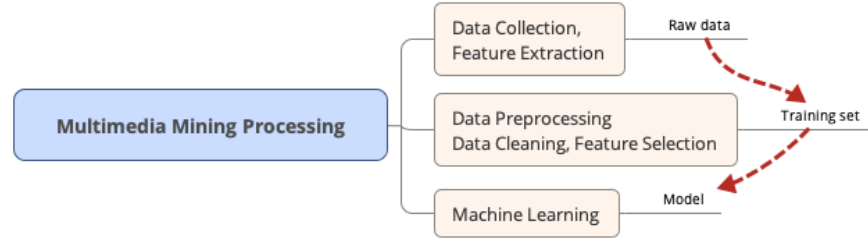


Figure 11: Multimedia Mining Processing

Data collection is the very first step in multimedia mining process. It acts as a raw data which are further input to the data preprocessing stage, which includes several task such as data cleaning and feature selection. After data preprocessing, a training set is obtained, which are further refine by applying machine learning to obtain a realistic model [3].

### 4.2 Web Mining

The Internet has now become a tool to support the planning of activities that require information gathering and reasoning. A major activity on the Internet is the retrieval and browsing of multimedia information [39]. It consist of a huge, widely, distributed, advertisements, consumer information, education, government, e-commerce, and many other services. The main task of web mining includes mining of web contents, web access patterns and web linkage structures as shown in Figure 12. This involves mining the web page layout structure, mining the web's link structures to identify authorize web pages, mining multimedia data on the web, automatic classification of web documents, and web usage mining [38].

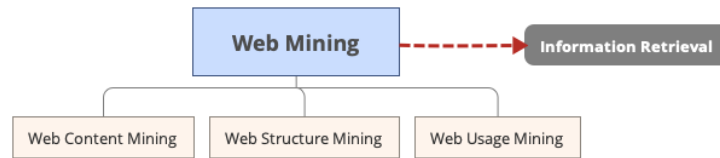


Figure 12: Web Mining Process

*A. Web Content Mining:* The process of extracting useful information from the contents of web documents is called web content mining. It may consist of text, images, audio, video information which is used to convey to the users about that documents [38].

*B. Web Structure Mining:* The process of discovering structure information from the Web is known as web structure mining. It consists of Web pages as nodes, and hyperlinks as edges connecting between two related pages [38].

*C. Web Usage Mining:* It is the method of discovering interesting usage patterns from web data, in order to understand and better serve the needs of web-based applications [38].

However the web also poses great challenges for effective resource and knowledge discovery, some of them are as follows:

1. The web seems to be too huge for effective data warehousing and data mining.
2. The complexity of web pages is far greater than that of any traditional text document collection.
3. The web is a highly dynamic information source.
4. The web serves a broad diversity of user communities.
5. Only a small portion of the information on the web is truly relevant or useful.

These challenges have promoted research into efficient and effective discovery and uses of resources on the internet. There are many search index -based web search engines. These search the web, index web pages, and build and store huge keyword-based indices that help locate sets of web pages containing certain keywords [6]

**Link Analysis** Web mining technique provides the additional information through hyperlinks where different documents are connected [40]. We can view the web as a directed labelled graph whose nodes are the documents or pages and edges are the hyperlinks between them. This directed graph structure is known as web graph. There are number of algorithms proposed based on link analysis. Three important algorithms are:

1. Page Rank [41],
2. Weighted Page Rank [42]
3. HITS (Hyper-link Induced Topic Search) [43].

Web mining is the powerful technique used to extract the information from past behaviour of users. There are various algorithms that are used for the web mining. Web mining is classified as three categories. In web structure mining we are using PageRank, Weighted Page Rank, and HITS. In web usage mining we used Clustering, K-Means Algorithm. In web content mining we are used correlation algorithm for ranking [44].

### 4.3 Text mining

Text mining is the process of finding useful or interesting patterns, models, directions, trends, or rule from unstructured text [38]. Latent Semantic Indexing is the latest and efficient approach to access databases that contain text. LSI associates with each document a limited-size vector  $vec(d)$  that contains frequency terms. Thus, the storage of documents in database becomes equivalent with storing the associated vectors. The basic idea of the technique is that similar documents have similar frequencies of words. The technique allows both the elimination of words and phrases that do not allow distinguishing between various documents and the identification of the ones that do so. It can identify also the similar words [45]. In text mining there are two open problems: Polysemy, synonymy. Polysemy means that a single word can have multiple meanings where as Synonymy means that multiple words can have the same/similar meaning [38].

Text Mining or knowledge discovery from text (KDT) - first introduced by Fledman et al. [46] - refers to the process of extracting high quality of information from text (i.e. structured such as RDBMS data [47] [48], semi-structured such as XML and JSON [49] [50] [51], and unstructured text resources such as word documents, videos, and images). It widely covers a large set of related topics and algorithms for analyzing text, spanning various communities, including information retrieval, natural language processing, data mining, machine learning many application domains web and biomedical sciences [1] [2].

#### 4.3.1 Information Retrieval (IR).

Information Retrieval is the activity of finding information resources (usually documents) from a collection of unstructured data sets that satisfies the information need [52] [53]. Therefore information retrieval mostly focused on facilitating information access rather than analyzing information and finding hidden patterns, which is the main purpose of text mining. Information retrieval has less priority on processing or transformation of text whereas text mining can be considered as going beyond information access to further aid users to analyze and understand information and ease the decision making. Natural Language Processing (NLP): Natural Language Processing is sub-field of computer science, artificial intelligence and linguistics which aims at understanding of natural language using computers [54] [55]. Many of the text mining algorithms extensively make use of NLP techniques, such as part of speech tagging (POG), syntactic parsing and other types of linguistic analysis [56], [57] [123] [124].

#### 4.3.2 Information Extraction from text (IE).

Information Extraction is the task of automatically extracting information or facts from unstructured or semi-structured documents [58] [59]. It usually serves as a starting point for other text mining algorithms. For example extraction entities, Name Entity Recognition (NER), and their relations from text can give us useful semantic information. Text Summarization: Many text mining applications need to summarize the text documents in order to get a concise overview of a large document or a collection of documents on a topic [60] [61]. There are two categories of summarization techniques in general: extractive summarization where a summary comprises information units extracted from the original text, and in contrary abstractive summarization where a summary may contain synthesized information that may not occur in the original document [62] [63].

#### 4.3.3 Unsupervised Learning Methods.

Unsupervised learning methods are techniques trying to find hidden structure out of unlabeled data. They do not need any training phase, therefore can be applied to any text data without manual effort. Clustering and topic modeling are the two commonly used unsupervised learning algorithms used in the context of text data. Clustering is the task of segmenting a collection of documents into partitions where documents in the same group (cluster) are more similar to each other than those in other clusters. In topic modeling a probabilistic model is used to determine a soft clustering, in which every document has a probability distribution over all the clusters as opposed to hard clustering of documents. In topic models each topic can be represented as a probability distributions over words and each documents is expressed as probability distribution over topics. Thus, a topic is akin to a cluster and the membership of a document to a topic is probabilistic [64] [65].

#### 4.3.4 Supervised Learning Methods.

Supervised learning methods are machine learning techniques pertaining to infer a function or learn a classifier from the training data in order to perform predictions on unseen data. There is a broad range of supervised methods such as nearest neighbor classifiers, decision trees, rule-based classifiers and probabilistic classifiers [66] [67].

#### 4.3.5 Probabilistic Methods for Text Mining.

There are various probabilistic techniques including unsupervised topic models such as probabilistic Latent semantic analysis (pLSA) [68] and Latent Dirichlet Allocation (LDA) [69], and supervised learning methods such as conditional random fields [70] that can be used regularly in the context of text mining.

#### 4.3.6 Text Streams and Social Media Mining.

There are many different applications on the web which generate tremendous amount of streams of text data. news stream applications and aggregators such as Reuters and Google news generate huge amount of text streams which provides an invaluable source of information to mine. Social networks, particularly Facebook and Twitter create large volumes of text data continuously. They provide a platform that allows users to freely express themselves in a wide range of topics. The dynamic nature of social networks makes the process of text mining difficult which needs special ability to handle poor and non-standard language [71] [72].

#### 4.3.7 Opinion Mining and Sentiment Analysis.

With the advent of e-commerce and online shopping, a huge amount of text is created and continues to grow about different product reviews or users opinions. By mining such data we find important information and opinion about a topic which is significantly fundamental in advertising and online marketing (see [73] for an overview). Biomedical Text Mining: Biomedical text mining refers to the task of text mining on text of biomedical sciences domains. The role of text mining in biomedical domain is two fold, it enables the biomedical researchers to efficiently and effectively access and extract the knowledge out of the massive volumes of data and also facilitates and boosts up biomedical discovery by augmenting the mining of other biomedical data such as genome sequences and protein structures [74] [1].

### 4.4 Image mining

Image mining is the technique to detect unusual patterns and extract implicit and useful data from images which are stored in the large databases. It deals with making associations between different images from large image databases. The applications areas of Image mining are medical diagnosis, remote sensing, agriculture, industries and space research



and also handling hyper spectral images. Images database include maps, geological structures, and biological structures. The primary challenge in image mining is to reveal out how low-level pixel representation enclosed in a raw image or image sequence can be processed to recognize high-level image objects and relationships [38]. In today's databases, images are represented as relationships, or as spatial data structures or with help of image transformations [45]. Existing query difficulty estimation methods play an important role in the text based information retrieval. However, it cannot be directly applied to content-based image retrieval, due to the complex structure of image queries and the well known semantic gap [4]. Image mining performs several operations before generating Knowledge. Firstly the images stored in the databases are pre-process and then they undergo transformation and feature extraction. After this the mining is performed on the collected data and then after Interpretation and evaluation the Knowledge is generated. The entire process is described as the Figure 13 below.

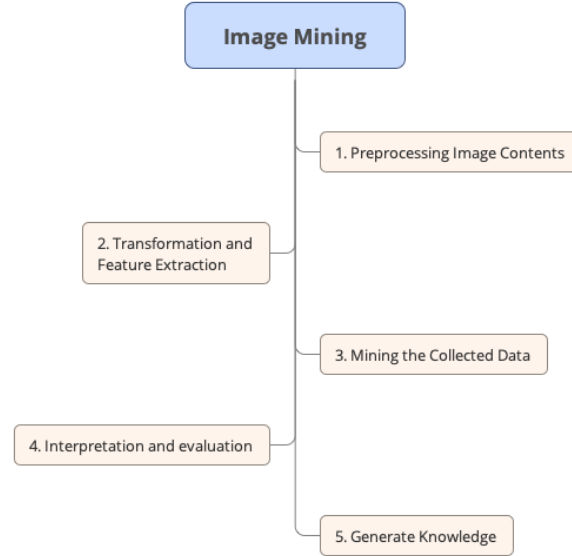


Figure 13: Image Mining Processing

Image mining includes object recognition, image indexing and retrieval, image classification and clustering, association rules mining, and neural network [75] [76].

#### 4.4.1 Object Recognition.

Using object models which might be known a priori, an object recognition technique finds objects in actuality from an image. Machine learning and purposeful information extraction can simply be realized when some objects have been identified and recognized through machine. The object recognition problem might be refer to as any supervised labelling problem according to models of known items i. e. given a target image containing a number interesting objects and a collection of labels corresponding to a collection of models known to technique, what is object recognition to assign correct product labels to regions, or a collection of regions, in the image.

#### 4.4.2 Image Retrieval.

Image mining requires that images be retrieved according to some requirement specifications. The requirement specifications can be classified into three levels of increasing complexity: (a) Level 1 comprises low level features of such as color, texture, shape or the spatial location of image elements. (b) Level 2 comprises image retrieval by derived or logical features like objects of a given type or individual objects or persons. (c) Level 3 comprises high level features of image.

#### 4.4.3 Image Indexing.

To further improve image retrieval rate, there is require of image data base using a fast and useful indexing scheme. A couple of main approaches are usually: reducing dimensionality or indexing high dimensional info. Other proposed indexing schemes concentrate on specific image features including color, shape and texture features [83] [84].



#### 4.4.4 Image Classification.

In supervised classification technique, as input a collection of labelled (Pre-classified) images are given, and here the problem is to label a newly Encountered, yet unlabeled images. Typically, the given Labelled (training) images are used to do the machine learning of the class description which in turn is use to label a new Image [81] [82].

#### 4.4.5 Image Clustering.

In unsupervised classification (or image clustering), the problem is always to group a given assortment of unlabeled images straight into Meaningful clusters based on the image content with not a priori knowledge. Clustering is often more advantage for minimizing the searching time period of images inside database. There are a variety of clustering methods: hierarchal, partitioning, density-based, grid based and fuzzy clustering methods [89] [90] [91].

#### 4.4.6 Association rules mining.

Association rule mining generates rules who have support and confidence greater than some user specific minimum support in addition to minimum confidence thresholds. A normal association rule mining algorithm works within two steps. The 1st step finds all substantial item sets that match the minimum support constraint. The second move generates rules from each of the large item sets that match the minimum confidence constraint.

#### 4.4.7 Neural network.

Neural Networks are computational systems made up of simple processing units called neurons which are usually organized into layers with fully or partially connections. The main task associated with a neuron is to receive the activation values from its neighbours (the output of other neurons), compute an output based on its weighted input parameters and send that output to its neighbours [77] [78].

Image mining is use in various fields [79] [80]. Different applications of image are;

1. In medical for diagnose diseases.(e.g. Brain tumour)
2. Satellite Cloud Imagery (e.g. Detecting copying unauthorized image on internet)
3. In Natural scene recognition
4. In Space research
5. In Remote sensing
6. In Detection of wild plant(e.g. egeria detection)
7. In Agriculture field
8. In industrial work
9. In educational field

### 4.5 Video mining

A video can be defined as scene sequences, which are composed by shots that are frame sequences. A frame is a static image, where as a shot is a clip that presents a continuous action in time and space [95] [96]. Video data mining is more typical than mining of image data. It is a collection of moving images, much like animation. Developing query and retrieval techniques for video databases are an important area which includes optimization strategies, video indexing, and query languages. The produced video (movies), the raw video (traffic video) and the medical video (ultra sound videos) are the three types of video [45]. Conventional techniques for video indexing are based on characterizing the video content using a set of computational features. These techniques often do not exploit the commonly available high-level information conveyed by textual captions. Relevant videos retrieval from large video databases has vast applications in various domains. Retrieval from broadcast news video database is vital for effective information access [97] [98]. Video mining techniques can also be used for mining digital libraries of video clips of lectures for distance learning applications. The two important issues in Video mining are to develop a representational scheme for the content and a Human friendly query/interface [92]. Figure 14 shows general framework for video data mining. There are many video mining approaches and they are roughly classified into five categories. They are: Video pattern mining, Video clustering and classification, Video association mining, Video content structure mining and Video motion mining [93] [94].

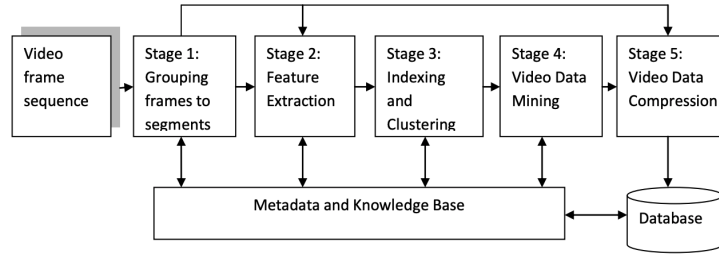


Figure 14: General framework for video data mining

#### 4.5.1 Video structure mining.

The main objective of the video structure mining is the identification of the content structure and patterns to carry out the fast random access of the video database [99] [100] [101]. Video structure mining is defined as the process of discovering the fundamental logic structure from the preprocessed video program adopting data-mining method such as classification, clustering and association rule. It is essential to analyze video content semantically and fuse multimodality information to bridge the gap between human semantic concepts and computer low-level features from both the video sequences and audio streams. Video structure mining is executed in the following steps: (1) video shot detection, (2) scene detection, (3) scene clustering, and (4) event mining. The current researches on it focus on mining object semantic information and event detection.

#### 4.5.2 Video clustering and classification.

Video clustering and classification are used to cluster and classify video units into different categories. Therefore clustering is a significant unsupervised learning technique for the discovery of certain knowledge from a dataset. Clustering video sequences in order to infer and extract activities from a single video stream is an extremely important problem and so it has a significant potential in video indexing, surveillance, activity discovery and event recognition. In the video surveillance systems, it is to find the patterns and groups of moving objects that the clustering analysis is used. Clustering algorithms are categorized into partitioning methods, hierarchical methods, density-based methods, grid based methods and model-based methods. Video classification aims at grouping videos together with similar contents and to disjoin videos with nonsimilar contents and thus categorizing or assigning class labels to a pattern set under the supervision.

#### 4.5.3 Video association mining.

Video association mining is the process of discovering associations in a given video. The video knowledge is explored in a two stages, the first being the video content processing in which the video clip is segmented into certain analysis units extracting their representative features and the second being the video association mining that extracts the knowledge from the feature descriptors. In video association mining, the video processing and the existing data-mining algorithms are seamlessly integrated into mine video knowledge.

#### 4.5.4 Video motion mining.

Motion is a key feature that essentially characterizes the contents of the video. There have been some approaches to extract camera motion and motion activity in video sequences. While dealing with the problem of object tracking, algorithms are always proposed on the basis of known object region in the frames and so the most challenging problem in the visual information retrieval is the recognition and detection of the objects in the moving videos. The camera motion having a vital role to play some of the key issues in video motion detections are, the camera placed in static location while the objects are moving (surveillance video, sports video); the camera is moving with moving objects (movie); multiple cameras are recording the same objects. The camera motion itself contains a copious knowledge related to the action of the whole match.

#### 4.5.5 Video pattern mining.

Video pattern mining detects the special patterns modeled in advance and usually characterized as video events such as dialogue, or presentation events in medical video. The existing work can be divided into two categories such as mining similar motion patterns and mining similar objects.

### 4.6 Audio mining

Since audio data is similar to video, the techniques for information processing and mining of audio data are also similar to video information retrieval and mining. Mining of audio data requires conversion of audio data into text by using speech transcription techniques. By using audio information processing techniques audio data can also be mined directly [38]. As for images or videos, audio data can be characterized in two ways: by using metadata to explain the content (objects or activities) or by extracting specific features with signal processing techniques (for instance, frequency, amplitude, vibratory period etc.). For audio data mining the most common content-based indexing technique is to segment the signal in time, to get small windows in which it can be considered homogeneous (amplitude, speed and wave length are constant). The segmentation is made with one constant step (window length) or by using a homogeneity predicate [45].

#### 4.6.1 Audio mining with feature extraction.

In multimedia application, Audio data plays vital role. Cory McKay and David Bainbridge [102] describes that Music information basically have two categories. a)Symbolic and b)Audio information. Audio is now became the continuous media type like videos. The techniques used in audio mining is similar to techniques used in video mining. audio data can be available in any form such as speech, music, radio, spoken language etc. The primary need for mining the audio data is the conversion of audio into text, using speech transcription technique this process can be done. other techniques are also available for this such as keyword extraction and then mining the text. Audio mining is that type of technique which is used to search audio files. K.A.Senthildevi, Dr.E.Chandra [103] explains that there are two main approaches of audio mining. 1) Text based indexing and 2) Phoneme based indexing. Text based indexing deals with the conversion process of speech to text. And Phoneme based indexing does not deals with conversion from speech to text, but instead works only with sound. Figure 15 shows the process of Audio mining.

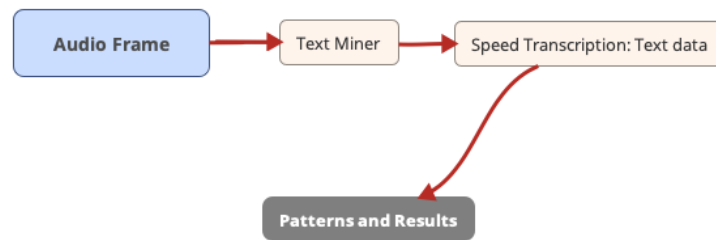


Figure 15: General framework for Audio Mining

The main objective of audio mining technology is to search through speech for identifying specific characteristics. Audio Mining can analyze and search the contents of an audio signal for identifying patterns and associations, retrieving keywords and information. Audio can be in the form of radio, speech, music, etc. Due to the continuous, dynamic and non-structured nature of audio, audio files must be represented with spectral coefficients for further processing with data mining techniques. There are variant feature types to represent speech characteristics in numeric forms. The audio mining system can run at high speed that is several times faster than that of traditional systems. Hence large quantities of audio or speech can be searched in a short time.

#### 4.6.2 Classification Of Audio Mining System.

Audio mining techniques can be roughly classified into three techniques, namely Keyword Spotting System (KWS), Wake-up-word (WUW) detection, and Spoken term Detection (STD). KWS approach aims to detect the occurrences of keywords within the test spoken utterance. WUW is related to KWS, but it uses speech commands to activate or wake up other systems by an alerting signal. Spoken Term Detection (STD) which is defined by NIST as an audio mining is employed for content based indexing. STD is aimed at open-vocabulary search over large collections of spoken documents. Similar to keyword spotting which involves finding occurrences of specific keywords in a speech utterance,

STD extends the same by finding a sequence of multiple words in the speech utterance. However, keyword spotting is considered as a part of STD and both of them have been addressed in this survey.

#### 4.6.3 Design Of Audio Mining System.

Audio mining system consists of two main phases: training and template matching [104]. During the training phase, the input speech is preprocessed and training vectors are generated from the speech signal. The training vectors extract the spectral features for distinguishing different classes of words. The feature vectors are quantized and codebook is generated for further processing. During the template matching phase, the given test keyword is matched with the generated codebook using template matching algorithm. Similarity distance is calculated and the occurrences of keyword is found with the given threshold value.

#### 4.6.4 Audio Mining Principles.

To mine audio data, first it has to be converted into text or it has to be divided as phonemes or used directly. There are three basic principles that have been developed over the years for audio mining [105] [106].

1. *VCSR Audio Mining.* In LVCSR audio mining, the entire test audio data is first transcribed into text and then the text transcription is used for searching the keywords. Hence LVCSR audio mining is a two-step process. In first phase, this method converts speech to text and in second phase, it identifies keywords in the generated dictionary that can contain several hundred thousand entries. If the keyword is not in the dictionary, the system will choose the most similar word it can find. LVCSR systems are more complex and expensive to implement.

2. *Acoustic Audio Mining.* In acoustic audio mining, the keyword search can be performed directly from untranscribed audio stream. Speech signals are represented as bark based energy or mel scale coefficients or they are represented as phoneme posteriorgrams generated by an acoustic model. This approach is easy to implement and provides some pronunciation tolerance.

3. *Phonetic Audio Mining.* Phonetic audio mining is phoneme based indexing method. This method does not convert speech to text but divides into phonemes. Phonetic audio mining is also a two-step process. In the first step, audio is processed (indexed) with a phonetic recognizer to generate a phonetic index file. In second step, system uses the generated dictionary of phonemes to compare user's search term to the correct phonetic string. This approach combines the advantages of LVCSR based keyword spotting and acoustic keyword spotting.

### 4.7 Particular social problems

There are different kinds of applications [91] [107] [108] of multimedia data mining, some of which are as follows:

#### 4.7.1 Digital Library.

The collection of digital data are stored and maintained in digital library, which is essential to convert different formats of digital data into text, images, video, audio, etc.

#### 4.7.2 Traffic Video Sequences.

In order to determine important but previously unidentified knowledge from the traffic video sequences, the detailed analysis and mining to be performed based on vehicle identification, traffic flow, and queue temporal relations of the vehicle at intersection. This provides an economic approach for regular traffic monitoring processes.

#### 4.7.3 Medical Analysis.

Multimedia mining is primarily used in the medical field and particularly for analyzing medical images. Various data mining techniques are used for image classification. For example, Automatic 3D delineation of highly aggressive brain tumors, Automatic localization and identification of vertebrae in 3D CT scans, MRI Scans, ECG and X-Ray.

#### 4.7.4 Customer Perception.

It contains details about customers opinions, products or services, customers complaints, customers preferences, and the level of customer's satisfaction of products or services which are collected together. Many companies have call centers that receives telephone calls from the customers. The audio data serves as topic detection, resource assignment and evaluation of quality of services.

#### 4.7.5 Media Making and Broadcasting.

Radio stations and TV channels creates broadcasting companies and multimedia mining can be applied to monitor their content to search for more efficient approaches and improve their quality.

#### 4.7.6 Surveillance system.

It consists of collecting, analyzing, summarizing audio, video or audio visual information about specific areas like government organizations, multi-national companies, shopping malls, banks, forest, agricultural areas and highways etc. The main use of this technology in the field of security hence it can be utilized by military, police and private companies since they provide security services.

## 5 MULTIMEDIA DATA MINING PROCESS

Figure 16 shows present architecture which includes the types of multimedia mining process [109]. Data Collection is the initial stage of the learning system; Pre-processing is to extract significant features from raw data, it includes data cleaning, transformation, normalization, feature extraction, etc. Learning can be direct, if informative types can be recognized at pre-processing stage. Complete process depends extremely on the nature of raw data and difficulty's field. The product of pre-processing is the training set. Specified training set, a learning model has to be selected to learn from it and make multimedia model is more constant.

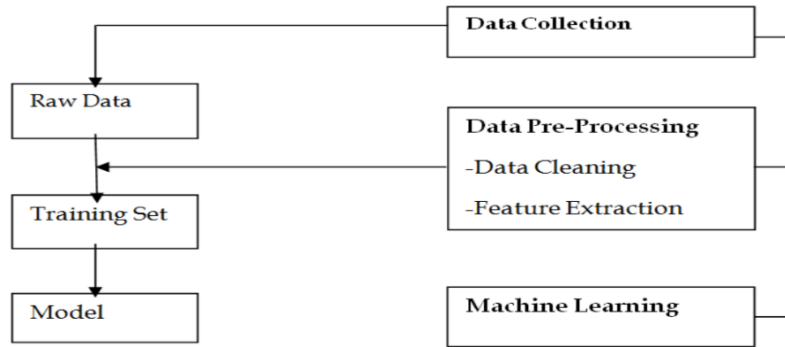


Figure 16: Multimedia data mining process

*Converting Un-structured data to structured data.* Data resides in fixed field within a record or file is called structured data and these data are stored in sequential form. Structured data has been easily entered, stored, queried and analyzed. Unstructured data is bit stream, for example pixel representation for an image, audio, video and character representation for text [1]. These sorts of files may have an internal structure, they are still considered unstructured because the data they contain does not fit neatly in a database. For example, image and video of different objects has some similarity - each represents an interpretation of a building - but then without clear structure.

Current data mining tool operate on structured data, which resides in huge volume of relational database while data in multimedia databases are semi-structured or un-structured. Hence, the semi-structured or unstructured multimedia data is converted into structured one, and then the current data mining tools are used to extract the knowledge. A difference between unstructured data and structured data mining is the sequence or time element. The architecture of converting unstructured data to structured data and it is used for extracting information from unstructured database is shown in Figure 17. Then data mining tools are applied to the stored structured databases.

## 6 ARCHITECTURES FOR MULTIMEDIA DATA MINING

Multimedia mining architecture is given in Figure 18. The architecture has several components. Important components are (1) Input (2) Multimedia Content (3) Spatiotemporal Segmentation (4) Feature Extraction (5) Finding the similar Patterns and (6) Evaluation of Results.

1. *Input* stage comprises which multimedia database is used for finding the patterns and to perform data mining process.

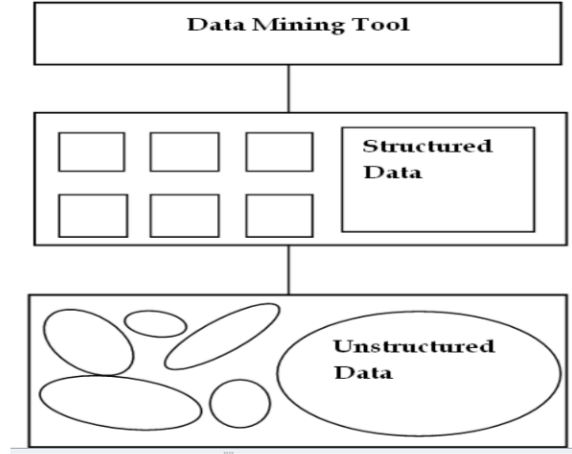


Figure 17: Unstructured data to structured data conversion

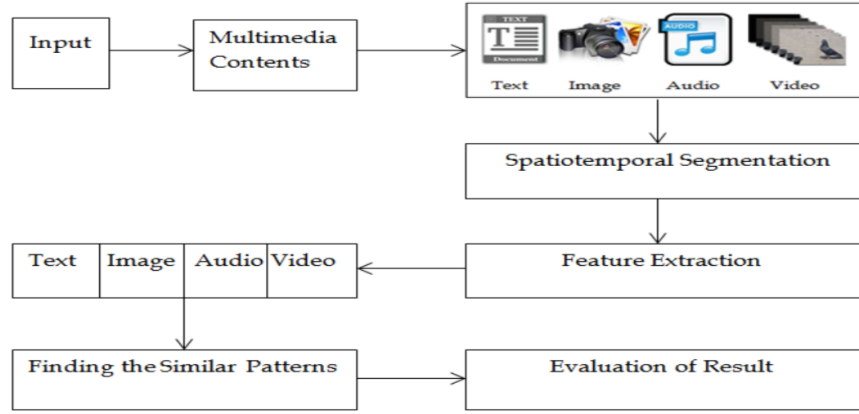


Figure 18: Multimedia data mining architecture

2. *Multimedia Content* is the data selection stage which requires the user to select the databases, subset of fields or data to be used for data mining.

3. *Spatio-temporal segmentation* is nothing but moving objects in image sequences in the videos and it is useful for object segmentation.

4. *Feature extraction* is the pre-processing step that involves integrating data from various sources and making choices regarding characterizing or coding certain data fields to serve when inputs to the pattern finding stage. Such representation of choices is required because certain fields could include data at various levels and not considered for finding the similar pattern stage. In MDM the preprocessing stage is significant since the unstructured nature of multimedia records.

5. *Finding the similar pattern stage* is the heart of the whole data mining process. The hidden patterns and trends in the data are basically uncovered in this stage. Some approaches of finding similar pattern stage contain association, classification, clustering, regression, time-series analysis and visualization.

6. *Evaluation of Results* is a data mining process used to evaluate the results and this is important to determine whether prior stage must be revisited or not. This stage consists of reporting and makes use of the extracted knowledge to produce new actions or products and services or marketing strategies [108].

## 7 RESEARCH ISSUES IN MULTIMEDIA SYSTEM DATA MINING

Major Issues in multimedia data mining contains content based retrieval, similarity search, dimensional analysis, classification, prediction analysis and mining associations in multimedia data [110].

### 7.1 Content based retrieval and Similarity search

Content based retrieval in multimedia is a stimulating problem since multimedia data is required for detailed analysis from pixel values [111]. We considered two main families of multimedia retrieval systems i.e. similarity search in multi-media data. (1) Description-based retrieval system created indices and make object retrieval, based on image descriptions, for example keywords, captions, size, and time of creation. (2) Content-based retrieval system supports retrieval on the image content, for example color histogram, texture, shape, objects and wavelet transforms. Use of content-based retrieval system: Visual features to index images and promotes object retrieval based on feature similarity; it is very desirable in various applications [112]. These applications which include diagnosis, weather prediction, TV production and internet search engines for pictures and e-commerce.

### 7.2 Multidimensional Analysis

In order to perform multidimensional analysis of large multimedia databases, multimedia data cubes may be designed and constructed in a method similar to that for traditional data cubes from relational data. A multimedia data cube can have additional-dimensions and measures for multimedia data, such as color, texture, and shape. A multimedia data cube has several dimensions. Examples are: size of the image or video in bytes; width and height of the frames, creating two dimensions, date on which image or video was created or last modified, format type of the image or video, frame sequence duration in seconds, Internet domain of pages referencing the image or video, the keywords like a color dimension and edge orientation dimension. Multimedia data mining system prototype is called MultiMediaMiner which is the extension of DBMiner system handles multimedia data. The Image Excavator component of MultiMediaMiner uses image contextual information, like HTML tags in Web pages, to derive keywords [113]. By navigating on-line directory structures, like Yahoo! directory, it is possible to build hierarchies of keywords mapped on the directories in which the image was found.

### 7.3 Classification and Prediction Analysis

Classification and predictive analysis has been used for mining multimedia data particularly in scientific analysis like astronomy, seismology, and geo-scientific analysis. Decision tree classification is an important data mining method in reported image data mining applications. For example, consider the sky images which has been carefully classified by astronomers as the training set, it can create models for the recognition of galaxies, stars and further stellar objects, based on properties like magnitudes, areas, intensity, image moments and orientation. The image data are frequently in large volumes and needs substantial processing power, for example, parallel and distributed processing. Image data mining classification and clustering are carefully connected to image analysis and scientific data mining and hence many image analysis techniques and scientific data analysis methods could be applied to image data mining [113].

### 7.4 Mining Associations in Multimedia Data

Association rules involving multimedia objects have been mined in image and video databases. Three categories can be observed: 1. Associations between image content and non-image content features 2. Associations among image contents that are not related to spatial relationships 3. Associations among image contents related to spatial relationships. The associations between multimedia objects, we can treat every image as a transaction and find commonly occurring patterns among different images. First, an image contains multiple objects, each with various features such as color, shape, texture, keyword and spatial locations, so that there can be a huge number of possible associations. Second, a picture containing multiple repeated objects is an essential feature in image analysis, recurrence of the similar objects should not be ignored in association analysis. Third, to find the associations between the spatial relationships and multimedia images and this can be used for discovering object associations and correlations [112].

## 8 FUTURE OF MACHINE LEARNING AND DATA MINING IN MULTIMEDIA SYSTEM

Multimedia mining is one of the important and challenging research domains in the field of computer science. Most of the researchers are interested to do their research work in the field of multimedia mining. Many challenging research

problems are available in multimedia mining. These problems can be solved by developing new algorithms, concepts and techniques for extracting hidden knowledge from the multimedia data bases [115] [116].

Multimedia Data Mining is the mining and analysis of various types of data, including images, video, audio, and animation. The idea of mining data which contains different kinds of information is the main objective of multimedia data mining. As multimedia data mining incorporates the areas of text mining, as well as hypertext/hypermedia mining, these fields are closely related. Much of the information describing these other areas also applies to multimedia data mining. This field is also rather new, but holds much promise for the future. Multimedia information, because its nature as a large collection of multimedia objects, must be represented differently from conventional forms of data. One approach is to create a multimedia data cube which can be used to convert multimedia-type data into a form which is suited to analysis using one of the main data mining techniques, but taking into account the unique characteristics of the data. This may include the use of measures and dimensions for texture, shape, colour, and related attributes. In essence, it is possible to create a multidimensional spatial database. Among the types of analyses which can be conducted on multimedia databases include associations, clustering, classification, and similarity search. Another developing area in multimedia data mining is that of audio data mining (mining music). The idea is basically to use audio signals to indicate the patterns of data or to represent the features of data mining results. The basic advantage of audio data mining is that while using a technique such as visual data mining may disclose interesting patterns from observing graphical displays, it does require users to concentrate on watching patterns, which can become monotonous. But when representing it as a stream of audio, it is possible to transform patterns into sound and music and listen to pitches, rhythms, tune, and melody in order to identify anything interesting or unusual.

Due to the enormous success of various application areas of data mining, the field of data mining has been establishing itself as the major discipline of computer science and has shown interest potential for the future developments. Ever increasing technology and future application areas are always posing new challenges and opportunities for data mining, the typical future trends of data mining in multimedia system includes:

1. Standardization of data mining languages
2. Data preprocessing
3. Complex objects of data
4. Computing resources
5. Web mining
6. Scientific Computing
7. Business data

### **8.1 Standardization of data mining languages**

There are various data mining tools with different syntaxes, hence it is to be standardized for making convenient of the users. Data mining applications have to concentrate more in standardization of interaction languages and flexible user interactions.

### **8.2 Data Preprocessing**

To identify useful novel patterns in distributed, large, complex and temporal data, data mining techniques have to evolve in various stages. The present techniques and algorithms of data preprocessing stage are not up to the mark compared with its significance in finding out the novel patterns of data. In future there is a great need of data mining applications with efficient data preprocessing techniques.

### **8.3 Complex object of data**

Data mining is going to penetrate in all fields of human life; the presently available data mining techniques are restricted to mine the traditional forms of data only, and in future there is a potentiality for data mining techniques for complex data objects like high dimensional, high speed data streams, sequence, noise in the time series, graph, Multi-instance objects, Multi-represented objects and temporal data.

### **8.4 Computing Resources**

The contemporary developments in high speed connectivity, parallel, distributed, grid and cloud computing have posed new challenges for data mining. The high speed internet connectivity has posed a great demand for novel and efficient



data mining techniques to analyze the massive data which is captured of IP packets at high link speeds in order to detect the Denial of Service (DoS) and other types of attacks. Distributed data mining applications demand new alternatives in different fields, such as discovery of universal strategy to configure a distributed data mining, data placement at different locations, scheduling, resource management, and transactional systems etc [119] [91]. New data mining techniques and tools are needed to facilitate seamless integration of various resources in grid based environment. Moreover, grid based data mining has to focus seriously to address the data privacy, security and governance. Cloud computing is a great area to be focused by data mining, as the Cloud computing is penetrating more and more in all ranges of business and scientific computing. Data mining techniques and applications are very much needed in cloud computing paradigm.

### 8.5 Web mining

The development of World Wide Web and its usage grows, it will continue to generate ever more content, structure, and usage data and the value of Web mining will keep increasing. Research needs to be done in developing the right set of Webmetrics, and their measurement procedures, extracting process models from usage data, understanding how different parts of the process model impact various Web metrics of interest, how the process models change in response to various changes that are made-changing stimuli to the user, developing Web mining techniques to improve various other aspects of Web services, techniques to recognize known frauds and intrusion detection.

### 8.6 Scientific Computing

In recent years data mining has attracted the research in various scientific computing applications, due to its efficient analysis of data, discovering meaningful new correlations, patterns and trends with the help of various tools and techniques. More research has to be done in mining of scientific data in particular approaches for mining astronomical, biological, chemical, and fluid dynamical data analysis. The ubiquitous use of embedded systems in sensing and actuation environments plays major impending developments in scientific computing will require a new class of techniques capable of dynamic data analysis in faulty, distributed framework. The research in data mining requires more attention in ecological and environmental information analysis to utilize our natural environment and resources. Significant data mining research has to be done in molecular biology problems [121] [122].

### 8.7 Business Trends

Business data mining needs more enhancements in the design of data mining techniques to gain significant advantages in today's competitive global market place (E-Business). The Data mining techniques hold great promises for developing new sets of tools that can be used to provide more privacy for a common man, increasing customer satisfaction, providing best, safe and useful products at reasonable and economical prices, in today's E-Business environment.

## 9 CONCLUSION

The multimedia mining, knowledge extraction plays a vital role in multimedia knowledge discovery [114] [117]. This research work provides a lot of present appraisal and updates of multimedia system analysis. This work helps to studies the relevancy of the techniques and idea of data mining in multimedia system, an comprehensive overview of the state-of-the-art methods, algorithms machine learning and data mining in multimedia systems. Multimedia mining is one of the important and challenging research domains in the field of computer science [118]. Most of the researchers are interested to do their research work in the field of multimedia mining. Many challenging research problems are available in multimedia mining. These problems can be solved by developing new algorithms, concepts and techniques for extracting hidden knowledge from the multimedia data bases. This paper discussed the multimedia mining basic concepts, essential characteristics, architectures, models and applications. Emerging and open research issues in multimedia mining are also described in this paper [120].

## References

- [1] Hieu Tran, Maxim Shcherbakov. Detection and prediction of users attitude based on real-time and batch sentiment analysis of facebook comments. 2016.
- [2] Hieu Tran, Ngoc Tran, Son Nguyen, Hoan Nguyen, Tien N Nguyen. Recovering variable names for minified code with usage contexts. ICSE 2019.
- [3] Rabi Narayan Behera Kajaree Das. A survey on machine learning: Concept, algorithms and applications. 2017.
- [4] Linjun Yang Chao Xu Wei Bian Yangxi Lian, Bo Geng. Query difficulty estimation for image retrieval. 2012.

- [5] Dong Xu Yueting Zhuang Liang-Tien Chia Yi Yang, Fei Wu. Cross-media retrieval using query dependent search methods. 2010.
- [6] Han and Kamber. Data mining: Concepts and techniques. 2006.
- [7] Pabblo Marquez Neila Carlos Pacho Manuel Barrena, Elena Jurado. A flexible framework to ease nearest neighbor search in multidimensional data spaces. 2010.
- [8] Madiha Waris Farooque Azam Abdul Wahab Muzaffar. A survey of issues in multimedia databases. 2012.
- [9] Durgesh Kumar Mishra Sarla More. Multimedia data mining: A survey. 2012.
- [10] Janusz Swierzowicz. Multimedia data mining concept.
- [11] Janusz Swierzowicz. Multimedia data mining trends and challenges.
- [12] Bhavani Thuraisingham. Managing and mining multimedia databases. 2004.
- [13] A. Rutherford. Introducing anova and ancova: a glm approach. 2001.
- [14] G. D. Hutcheson and L. Moutinho. Statistical modeling for management. 2008.
- [15] J. Fox. An r and s-plus companion to applied regression. 2002.
- [16] G. D. Hutcheson and N. Sofroniou. The multivariate social scientist. 1999.
- [17] David W. Aha, Dennis Kibler, and Marc K. Albert. Instance-based learning algorithms. *Machine Learning*, 6(1):37–66, Jan 1991.
- [18] Michael R. Gray James M. Keller and JR. James A. Givens. A fuzzy K-nearest neighbor algorithm. 1985.
- [19] Xuemin Zhang Linna Li. Study of data mining algorithm based on decision tree. 2010.
- [20] T. Stich. Bayesian networks and structure learning. 2004.
- [21] M.I. Jordan. Learning in graphical models. 1999.
- [22] K. Murphy. A brief introduction to graphical models and bayesian networks. 1998.
- [23] H. David. A tutorial on learning with bayesian networks, in learning in graphical. 1999.
- [24] Salamon P. Hansen L. Neural network ensembles. 1990.
- [25] Burn D. H. Shu C. Earthquake prediction by rbf neural network ensemble. 2004.
- [26] Burn D. H. Shu C. Artificial neural network ensembles and their application in pooled flood frequency analysis. 2004.
- [27] Breiman L. Random forests. 2001.
- [28] Gall J. Fossati A. Gool L. Fanelli G., Dantone M. Random forests for real time 3d face analysis. 2012.
- [29] Qi Y. Random forest for bioinformatics, in ensemble machine learning. 2012.
- [30] Soatto S. Bissacco A., Yang M.-H. Fast human pose estimation using appearance and motion via multi-dimensional boosting regression. 2007.
- [31] Dietterich T. G. Hutchinson R. A., Liu L.P. Incorporating boosted regression trees into ecological latent variable models. 2011.
- [32] Brown KA Pittman SJ. Multi-scale approach for predicting fish species distributions across coral reef seascapes. 2011.
- [33] Tong Zhang Johnson R. Learning nonlinear functions using regularized greedy forest. 2014.
- [34] Srebro N. Sridharan K. Cotter A., Shamir O. Better mini-batch algorithms via accelerated gradient methods. 2011.
- [35] Bulatov Y. Dietterich T. G., Ashenfelter T. D. A. Training conditional random fields via gradient tree boosting. 2004.
- [36] Friedman J. Greedy boosting approximation: a gradient boosting machine. 2001.
- [37] Zhao Y. Hu T., Li X. Gradient boosting learning of hidden markov models. 2007.
- [38] Dr. Siddu. P. Algur Pravin M. Kamde. A survey on web multimedia mining. 2011.
- [39] Harutaka Yoneyama Subhash Bhalla Tomoko Izumita Nadia Bianchi-Berthouze, Naoto Katsumi. Supporting the interaction between user and web-based multimedia information. 2003.
- [40] P. Husbands H. Zha C. Ding, X. He and H. Simon. Link analysis: Hubs and authorities on the world. 2001.

- [41] R. Kosala and H. Blockeel. Web mining research: A survey. 2000.
- [42] S. Brin and L. Page. The anatomy of a large scale hypertextual web search engine. 1998.
- [43] Wenpu Xing and Ali Ghorbani. Weighted pagerank algorithm. 2004.
- [44] T. Rajamdhangi T. Suresh Kumar. A survey of web mining algorithms. 2016.
- [45] Monica Vladoiu Catalina Negotia. Querying and information retrieval in multimedia databases. 2006.
- [46] Ronen Feldman and Ido Dagan. Knowledge discovery in textual databases. 1995.
- [47] Jiawei Han Ming-Syan Chen and Philip S. Yu. Data mining: an overview from a database perspective. 1996.
- [48] Saso Dzeroski. Relational data mining. 2009.
- [49] Azadeh Iranmehr Mahmood Doroodchi and Seyed Amin Pouriyeh. An investigation on integrating xml-based security into web services. 2009.
- [50] Seyed Amin Pouriyeh and Mahmood Doroodchi. Secure sms banking based on web services. 2009.
- [51] Mahmood Doroodchi Seyed Amin Pouriyeh and MR Rezaeinejad. Secure mobile approaches using web services. 2010.
- [52] Prabhakar Raghavan Christopher D Manning and Hinrich Schutze. Introduction to information retrieval. 2008.
- [53] Christos Faloutsos and Douglas W Oard. A survey of information retrieval and filtering methods. 1998.
- [54] Elizabeth D Liddy. Natural language processing. 2001.
- [55] et al. Christopher D Manning. Foundations of statistical natural language processing. 1999.
- [56] Martin Rajman and Romaric Besancon. Text mining: natural language techniques and text mining applications. 1998.
- [57] Anne Kao and Stephen R Poteet. Natural language processing and text mining. 2007.
- [58] Jim Cowie and Wendy Lehnert. Information extraction. 1996.
- [59] Sunita Sarawagi et al. Information extraction. 2008.
- [60] Eduard Hovy Dragomir R Radev and Kathleen McKeown. Introduction to the special issue on summarization. 2002.
- [61] Andreas Hotho. A brief survey of text mining. 2005.
- [62] Dipanjan Das. A survey on automatic text summarization. 2007.
- [63] M. Assefi S. Safaei E. D. Trippe J. B. Gutierrez M. Allahyari, S. Pouriyeh and K. Kochut. Text summarization techniques: A brief survey. 2017.
- [64] Mark Steyvers and Tom Griffiths. Probabilistic topic models. 2007.
- [65] Charu C Aggarwal and ChengXiang Zhai. Mining text data. 2012.
- [66] Tom M Mitchell. Machine learning. 1997.
- [67] Fabrizio Sebastiani. Machine learning in automated text categorization. 2002.
- [68] Thomas Hofmann. Probabilistic latent semantic indexing. 1999.
- [69] Andrew Y Ng David M Blei and Michael I Jordan. Latent dirichlet allocation. 2003.
- [70] Andrew McCallum John Lafferty and Fernando CN Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. 2001.
- [71] Ling Jiang Christopher C Yang, Haodong Yang and Mi Zhang. Social media mining for drug safety signal detection. 2012.
- [72] Pritam Gundecha and Huan Liu. Mining social media: a brief introduction. 2012.
- [73] Bo Pang and Lillian Lee. Opinion mining and sentiment analysis. 2008.
- [74] Stephen Cantrell Juan B Gutierrez, Mary R Galinski and Eberhard O Voit. From within host dynamics to the epidemiology of infectious disease: scientific overview and challenges. 2015.
- [75] Kamaljit Kaur Prabhjeet Kaur. Review of different existing image mining techniques. 2014.
- [76] K.S. Yuvaraj K. R. Yasodha. A study on image mining techniques. 2013.
- [77] Wietske Bijker Rajasekar Umamaheshwaran and Alfred Stein. Image mining for modeling of forest fires from meteosat images. 2007.

- [78] Yosio Edemir Shimabukuro. Estimating burned area in mato grosso, brazil, using an object-based classification method on a systematic sample of medium resolution satellite images. 2015.
- [79] Dr.Vijayalakshmi MN Divya TL. Envisagation and analysis of air pollution caused by forest fire using machine learning algorithm. 2015.
- [80] Rhett D. Harrison Mohan Kumar Sammathuria Yong Poh Yu, Rosli Omar and Abdul Rahim Nik. Pattern clustering of forest fires based on meteorological variables and its classification using hybrid data mining methods. 2011.
- [81] Kiyun Yu Yong Il Kim Young Gi Byun, Yong Huh. Evaluation of graph-based analysis for forest fire detections. 2005.
- [82] Ning Han Fei Van, Xing Xu. Identification method of forest fire based on color space. 2010.
- [83] Vijayalakshmi M.N Divya TL. Implementation of data mining techniques for temperature extraction for forest fire occurrence with image analysis. 2013.
- [84] Vijayalakshmi M.N Divya TL. Development of frame work for prediction of forest fire and fire spread direction using image mining. 2013.
- [85] Vijayalakshmi M.N Divya T.L. Analysis of wild fire behaviour in wild conservation area using image data mining. 2015.
- [86] Mong Li Lee Ji Zhang, Wynne Hsu. Image mining: Trends and developments. 2001.
- [87] Deepika Kishor Nagthane. Image mining techniques and applications. 2013.
- [88] P.Manikandaprabhu T.Karthikeyan. "function and information driven frameworks for image mining - a review. 2013.
- [89] S .A. Ladhake A. A. Khodaskar. Image mining: An overview of current research,. 2014.
- [90] E. Annasaro A.Hema. A survey in need of image mining techniques. 2013.
- [91] A. Sakila S. Vijayarani. Multimedia mining research - an overview. 2015.
- [92] Balamanohar Paluri Nataraj Jammalamadaka C. V. Jawahar, Balakrishna Chennupati. Video retrieval based on textual queries. 2006.
- [93] Xingquan Z. Xindong W. Sequential association mining for video summarization. 2003.
- [94] Chi-Yao Tseng Yung-Yu Chuang Ken-Hao Liu, Ming-Fang Weng. ssociation and temporal rule mining for post-filtering of semantic concept detection in video. 2008.
- [95] JeongKyu Lee JungHwan Oh and Sanjaykumar Kote. Real time video data mining for surveillance video streams. 2015.
- [96] David Moore. A real-world system for human motion detection and tracking. 2003.
- [97] Rikard Laxhammar Christoffer Brax, Lars Niklasson. An ensemble approach for increased anomaly detection performance in video surveillance data. 2009.
- [98] Pradipta Kumar Nanda Badri Narayan Subudhi. A change information based fast algorithm for video object detection and tracking. 2011.
- [99] I Anwar, F. Petrounias. Efficient periodicity mining of sequential patterns in a post-mining environment. 2008.
- [100] Morris T. Kodogiannis. V. Anwar. F., Petrounias. I. Discovery of events with negative behavior against given sequential patterns. 2010.
- [101] M. E. Khalifa Ahmed Taha, Hala H. Zayed and El-Sayed M. El-Horbaty. On behavior analysis in video surveillance. 2013.
- [102] David Bainbridge Cory McKay. A musical web mining and audio feature extraction extension to the greenstone digital library software. 2011.
- [103] Dr.E.Chandra K.A.Senthildevi. Data mining techniques and applications in speech processing - a review. 2012.
- [104] Phonexia. Keyword spotting.
- [105] Gunjan Manpreet Kaur Mand, Diana Nagpal. An analytical approach for mining audio signals. 2013.
- [106] Anna M. Kruspe. Keyword spotting in a capella singing. 2014.
- [107] Ravikumar GK Manjunath T.N, Ravindra S Hegadi. A survey on multimedia data mining and its relevance today. 2010.

- [108] Latifur Khan Valery A. Petrushin. Multimedia data mining and knowledge discovery. 2007.
- [109] P. Pintelas S. Kotsiantis, D. Kanellopoulos. Multimedia mining. 2004.
- [110] Ankush Mittal. An overview of multimedia content-based retrieval strategies. 2006.
- [111] Wolf Yu H. Scenic classification methods for image and video databases. 1995.
- [112] MichelineKamber Jiawei Han. Data mining: Concepts and techniques. 2001.
- [113] Oka R Mori Y, Takahashi H. Image-to-word transformation based on dividing and vector quantizing images with word. 1999.
- [114] Balint Zoltan Daroczy. Machine learning methods for multimedia information retrieval. 2017.
- [115] Wei Xu Yihong Gong. *Machine Learning Techniques for Multimedia*. Springer, Berlin, Heidelberg, 2008.
- [116] Oded Maimon and Lior Rokach. *Introduction to Knowledge Discovery and Data Mining*, pages 1–15. Springer US, Boston, MA, 2010.
- [117] Tay E H Francis Wang Gouhua. Data mining: Concepts, applications and techniques. 2000.
- [118] Sathishkumar R Ambeth kumar V D Elangovan D, Dr.Subedha V. A survey: Data mining techniques for social media analysis. In *Advances in Engineering Research Proceedings of the International Conference for Phoenixes on Emerging Current Trends in Engineering and Management (PECTEAM 2018)*, 2018.
- [119] SAM Rizvi Pramod Kumar Yadav. An exhaustive study on data mining techniques in mining of multimedia database. In *2014 International Conference on Issues and Challenges in Intelligent Computing Techniques (ICICT)*, 2014.
- [120] Johannes Fürnkranz, Dragan Gamberger, and Nada Lavrač. *Foundations of Rule Learning*. Springer Publishing Company, Incorporated, 2014.
- [121] Wei Xu Yihong Gong. *Machine Learning for Multimedia Content Analysis*. Springer Publishing Company, Incorporated, 2007.
- [122] Osmar R. Zaiane, Jiawei Han, Ze-Nian Li, and Jean Hou. Mining multimedia data. In *Proceedings of the 1998 Conference of the Centre for Advanced Studies on Collaborative Research, CASCON '98*, pages 24–. IBM Press, 1998.
- [123] Hoan Anh Nguyen, Tien N Nguyen, Danny Dig, Son Nguyen, Hieu Tran, Michael Hilton. *Graph-based mining of in-the-wild, fine-grained, semantic code change patterns*. ICSE 2019.
- [124] Ngoc Tran, Hieu Tran, Son Nguyen, Hoan Nguyen, Tien N Nguyen. *Does BLEU score work for code migration?*. ICPC 2019.