# WEB PAGES CATEGORIZATION BASED ON CLASSIFICATION & OUTLIER ANALYSIS THROUGH FSVM

**Geeta R.B.*[1] --- Shobha R.B.[2] --- Shashikumar G.Totad[3] --- Prasad Reddy PVGD[4]**

*[1]Department of Information Technology, GMR Institute of Technology, RAJAM, AP, India*

*[2]Department of Electronics and Communications, Basaweshwar Engineering College, Bagalkot, India*

*[3]Department of Computer Science, GMR Institute of Technology, RAJAM, AP, India*

*[4]Department of CS & SE, Andhra University, Vizag, AP, India*

## ABSTRACT

*The performance of Support Vector Machine is higher than traditional algorithms. The training process of SVM is sensitive to the outliers in the training set. Here in this Paper, a new approach called, Web Pages Categorization based on Classification and Outlier Analysis (WPC-COA), is proposed that uses a polynomial Kernel function to map web page tuples to high dimensional feature space.*

**Keywords:** Support vector machine, Outliers, Categorization, Log file, Kernel parameters, Web page.

### Contribution/ Originality

This study uses a new methodology which helps in mapping web page tuples with various attributes such as frequency, time spent on each page, in-degree, out-degree and level of a web page to high dimensional feature space. The paper's primary contribution is to categorize web pages based on classification and outlier analysis using Polynomial Kernel function.

## 1. INTRODUCTION

Intelligence has been well-defined in several ways including logical reasoning, abstract thinking power, better understanding, self-awareness, interacting with others, learning by experiences, being emotional, retaining, scheduling and better problem resolving capabilities. Intelligence is extensively observed in humans, but also perceived in animals and plants. Artificial Intelligence is the modeling of intelligence in machines. Machine learning is a branch of Artificial Intelligence which used in building systems that can learn from data. This learning is further used in building systems that can learn from test data sets.

In machine learning, Support Vector Machine (SVM) is a supervised learning simulation integrated with learning algorithms that help in analysis, and distinguish patterns used for

classification and regression analysis. The basic SVM model accepts a set of input data and predicts which of the two possible classes the output belongs, making it a non-probabilistic binary linear classifier. For a given training data set, the set is marked as belonging to one of the two categories; SVM constructs a model that earmarks a given data set into one category or the other. SVM is a prominent approach for classification problems [1] and nonlinear function approximation problems in the fields of machine learning and pattern recognition. It provides accuracy measure of fuzzy classification for solving real-world problems [2]. Section 2 explains Support Vector Machine (SVM) and discusses about Fuzzy Support Vector Machine (FSVM) and Fuzzy C-Means Clustering (FCM). Section 3 presents FSVM's Web Pages Categorizer based on Classification and Outlier Analysis (WPC-COA). Section 4 shows experimental results.

## 2. SUPPORT VECTOR MACHINE (SVM)

Support Vector Machine is a technique for classification of both linear and nonlinear data. Web page tuple is a record of one or more attributes such as in-degree, out-degree, level, frequency and utility etc..of a page. Let 'W' be the web page (data) set, collection of web pages, denoted as{(P1, B1), (P2, B2)….. (Pn, Bn)], where Pi is a tuple of page 'i' with associated class label Bi. Each Bi can take one of two values either positive class (+1/Yes) or negative class (-1/No) as shown in Fig. 1. For a web page set with 2 attributes/dimensions, infinite number of separators can be drawn. This can be generalized to n dimensions/attributes. This optimal separator is known as hyperplane [3]. Hyperplane with a larger margin is more accurate than with smaller margin.
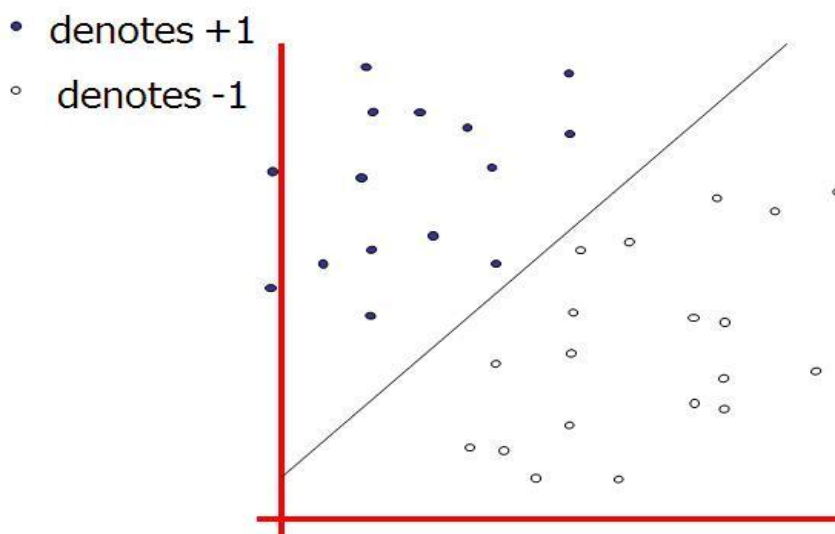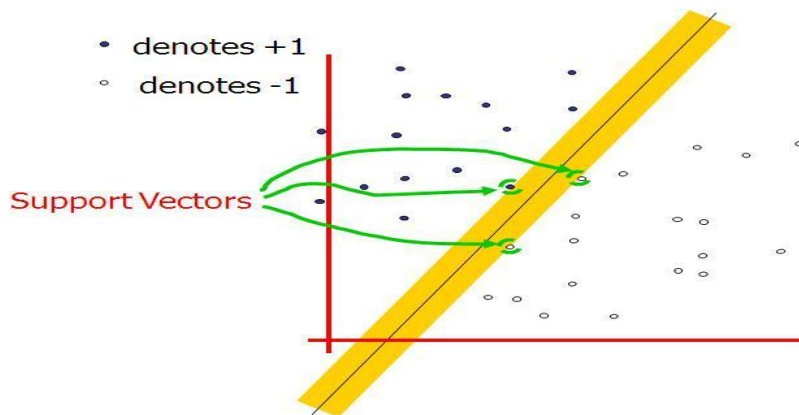
**Fig-1.**Web pages separator

**Fig-2.** Web Page as Support Vectors



A combination of few input data points, known as support vectors [3], can be an optimal solution as shown in Fig. 2.This SVM method of classification works fine when web page tuples are linear. Such a SVM is called linear SVM. For cases, where it is not possible to have a hyperplane as straight line or when web page tuples are nonlinear, extended SVM can be used. In such a case, input web page points are mapped into high dimensional feature space using nonlinear mapping, resulting in quadratic optimization problem. This problem can be solved using linear SVM. Optimal hyperplane in high dimensional feature space corresponds to nonlinear separating hyper surface in the original space. Linear and nonlinear SVMs are sensitive to noise or outliers. This is one of the main limitations of the standard SVM. In most of the real world classification problems, some training points do not belong to any of the predefined classes exactly, whenever noise or outlier exists. For example, one training point may belong 90% to one class and be 10% to none of the classes, and it may also belong 30% to one class and 70% none of the classes. SVM is a powerful tool for solving classification problems. Comparison [4] between traditional learning and SVM has shown that SVM provides better results than traditional learning in many applications.

### 2.1. Fuzzy Support Vector Machine (FSVM)

SVM draws optimal hard boundary hyperplane between two classes. It finds optimal hyperpalne using dot product functions, called Kernels. Each input web page may not completely belong to one of the two classes (+ve, -ve). This limitation can be overcome by Fuzzy Support Vector Machine (FSVM). Each input web page is assigned a fuzzy membership degree, so that input web page can make different contributions to the decision feature surface learning.

## 3. FUZZY C-MEANS CLUSTERING

Fuzzy C-means (FCM) is a clustering method. It allows one web page to belong to two or more clusters. This method proposed by Dunn in 1973, was modified in 1981 by Bezdek. FCM's aim is to minimize the following objective function

$$Fp = \sum_{i=1}^{N} \sum_{i=1}^{c} m_{ij}{}^{p} \, \|\|x_i - c_j\|\|^2$$

$$1 \leq p < \infty$$

where 'p' is a number greater than 1, $m_{ij}$ is the degree of membership of web page $x_i$ in cluster $j$, $x_i$ is the $i^{th}$ web page with 'd' attributes, $c_j$ is d–dimension center of cluster j, the center $c_j$ and any measured data's similarity is expressed using any norm $||*||$ is any norm expressing the similarity between any measured data and the center $c_j$.

Fuzzy partitioning is carried out through iterative optimization of the objective function shown above, by updating membership $m_{ij}$ and the cluster centers $c_j$ as

$$m_{ij} = \frac{1}{\sum_{k=1}^{c} \left( \|x_i - c_j\| / \|x_i - c_k\| \right)^{2 \div (m-1)}}$$

$$c_j = \frac{\sum_{i=1}^{N} m_{ij}{}^{p} \cdot x_i}{\sum_{i=1}^{N} m_{ij}{}^{p}}$$

This process is iterated till

$$max_{ij}\{|m_{ij}{}^{(k+1)} - m_{ij}{}^{(k)}| < \varphi \quad ,$$

where $\varphi$ value lies between 0 and 1 and is called as termination criterion, whereas $k$ is the iteration number. This procedure ultimately converges to a local minimum.

## 4. FSVM'S WEB PAGES CATEGORIZER BASED ON CLASSIFICATION AND OUTLIER ANALYSIS ALGORITHM

In FCM, the sum of memberships of web pages across classes is one. FCM uses probabilistic constraint. Krishnapuram and Keller relaxedthis constraint and proposed a possibilistic approach to clustering (PCM). Kernel Possibilistic C-means (KPCM) substitutes kernel-induced distance metric for the Euclidean distance [5]. KPCM algorithms avoid clusters coincidence and are less sensitive to outliers.

Consider website with n pages, n>0. Each page is associated with attributes such as in–degree (the number of web pages pointing to this page), out-degree (the number of pages to which this page is pointing), level, frequency and time spent on each page. Assuming there are 'm' web page

tuples, 80% of 'm' can be used for training and remaining 20% of 'm' can be used for testing. Polynomial Kernel is used to get a hyperplane.

$K(X_i, X_j) = (X_i \cdot X_j + b)^d$

Where 'd' is a degree of the polynomial function, and 'b' is a constant.

Given number of clusters 'C', degree of the polynomial function 'd' and constant 'b', dot product of $X_i$ and $X_j$ can be obtained. If $K(X_i, X_j) < 0$ then the tuple belongs to negative class otherwise it belongs to positive class. Outliers are not erroneous. Outliers are the tuples that do not comply with general behavior of the data. These are tuples which are far away from the other tuples in the given data space. Certain web page details are inconsistent with the remaining set of web page tuples, for example, a wrong entry of web page in-degree as -1. Noise can be of two types; attribute noise and class noise. Missing value of particular attribute (level of web page is missing) is an example of attribute noise. Class noise is of two types, namely contradictory examples and misclassifications. Example of contradictory noise, include the web page sets appearing more than once and being labeled with different classes. Examples of misclassification include two different classes with same symptoms or features. Here in this chapter, classification problems with outliers or misclassification noise are dealt with. Outliers and noise are assigned to lower membership degrees. Two farthest pair of clusters are classified to form a binary classification problem. This process degrades the effects of noise or outlier on the decision function. When the farthest pair of clusters is adopted, and outliers are degraded, FSVM is applied to obtain final classification results.

Consider a web site of 'n' web pages. The structure of website provides information about in-degree, out-degree, level, frequency of a web page and server log file helps in getting information about time spent on each web page in a specific period of time. These can be attributes of a web page. Get n pairs of web page tuples {$X_i$, $Y_i$, $S_i$} for a binary classification problem, where $X_i$ $\in R^n$ are the input web pages, $Y_i \in \{-1, +1\}$ are the corresponding binary class labels, and $S_i \in (0, 1]$ is the fuzzy membership degree of $X_i$ belonging to Y. Let $X[i]$ denote a web page with five attributes such as frequency, time spent, in-degree, out-degree and level of a web page. Sixth attribute is used as class label, which labeled as positive class (+1) or negative class (-1). 80% of log file data is used as training web pages. 20% log file data is used as testing data.

## 4.1. Algorithm WPC-COA (Server Log File, Website, Kernel Parameters, No. of Clusters)

1. Get frequency and time spent on each page from server log file.
2. Get in-degree, out-degree and level of each web page from website structure.
3. For each training web page tuple, sixth attribute of each web page is used as class label which is assigned a label either positive class (1) or negative class (-1).
4. Get n pairs of web page tuples {$X_i$, $Y_i$, $S_i$} for a binary classification problem, where $X_i$ $\in R^n$ are the input web pages, $Y_i \in \{-1, +1\}$ are the corresponding binary class labels, and $S_i \in (0, 1]$ is the fuzzy membership degree of $X_i$ belonging to Y.

5.  Read number of clusters, 'C' and termination parameter ' $\varphi$ '.

6.   Read 'b' and degree of polynomial function 'd'.

7.  Initialize cluster centroids  Xj, j = 1, 2. . . C.

8.   Map training web page tuples into high dimensional feature space using FCM and cluster them into positive class (P) or negative class (N).

9.  Decide whether pair, (Xi, Xj) of web page belongs to P or negative class N and to which cluster, using Kernel polynomial function.

    $K(Xi , Xj) = (Xi \cdot Xj + b)^d$. d is a degree of the polynomial function, and b is an constant.

10.   Let number of clusters in class P be Z+ and number of clusters in class N be Z-.

11.   Search for farthest pair of clusters, in which one cluster belongs to Z+ and the other belongs to Z-.

12.   Form a new training set with membership degrees, combining farthest pair of clusters and obtain a nonlinear classifier using FSVM with the parameters αi and b.

13.  Stop

## 5.  EXPERIMENTAL RESULTS

WPC-COA algorithm is implemented using java. Statistics collected from log file is stored in Excel with fields/attributes, frequency (f1)-number of times a web page is accessed, utility (f2)-time spent on a web page, number of backward links (f3), level (f4) of a web page. Field (f5) is used to randomly assign class label (positive +1, negative -1). Among 'n' number of log file records, 80% of web page tuples form training web page data set and remaining 20% of web page tuples are used as testing web page data set. Fig 3 to Fig. 9 shows output snapshots. Fig. 3 provides the user interface. When Browse button is clicked, it prompts for the location of file which contains web pages information. In Fig. 4 user selects the required file containing information about details of log file, with attributes such as frequency (f1), utility (f2), number of backlinks (f3), level (f4) and f5 with class label(+1 or -1). Fig. 5 shows testing data set. Fig. 6 and Fig. 7   prompt user for constant b and number of clusters respectively. Fig. 8 displays different clusters with positive and negative classes. Fig. 9 outputs distance between various clusters such as cluster distance between positive cluster 1 and negative cluster 1, positive cluster 1 and negative cluster 2, positive cluster 2 and negative cluster 1, positive cluster 2 and negative cluster 2.

**Fig-3.** GUI for Dataset selection



**Fig-4.** select the dataset



**Fig-5.** Testing Dataset

**Fig-6.** Browse the dataset Enter b value.



**Fig-7.** Enter C (number of Clusters)



**Fig-8.** Different clusters with positive and negative classes.



26

Fig-9. Clusters distance



The file statlog contains 37 (f1, f2, f3....f7) attributes and 2218 records. Fig. 10 displays different clusters with positive classes. Fig. 11 displays different clusters with negative classes. Fig. 12 outputs distance between various clusters such as cluster distance between positive cluster 1 and negative cluster 1, positive cluster 1 and negative cluster 2, positive cluster 2 and negative cluster 1, positive cluster 2 and negative cluster 2, positive cluster 1 and negative cluster 3, positive cluster 2 and negative cluster 3, positive cluster 3 and negative cluster 3.

Fig-10. Different clusters with positive class

**Fig-11.** Different clusters with negative class

**Fig-12.** Clusters distance



## 6. CONCLUSION

Here in this paper, machine intelligence is introduced. FSVM and Fuzzy C-Means algorithm along with Possibilistic K-means algorithm is discussed. WPC-COA, proposed algorithm, which helps in categorization of web pages is presented and experimented. KPCM algorithms avoid clusters coincidence and are less sensitive to outliers

## REFERENCES

[1]     J. C. Platt, *Sequential minimal optimization-A fast algorithm for training support vector machines, in advances in kernel methods-support vector learning.* Cambridge, MA: MIT Press, 1998.

[2]     D. Gomez, J. Montero, and G. Biging, "Improvements to remote sensing using fuzzy classification, graphs and accuracy statistics," *Pure Appl. Geophys.*, vol. 165, pp. 1555-1575, 2008.

[3]     E. O. Edgar, F. Robert, and G. Federico, *Support vector machines: Training and applications*: A.I. Memo No. 1602, C.B.C.L Paper No.144, March 1997, 2004.

[4]     P. Scott and H. Lutz, "Comparing the results of support vector machines with traditional data mining algorithms, supported by Amica Life Insurance Company."

[5]     W. Xiao-Hong, "Coll. of electr and inf eng, Jiangsu Univ, Zhenjiang. A possibilistic C- means clustering algorithm based on kernel methods," presented at the Computational Intelligence for Modelling, Control and Automation, 2005 and International Conference on Intelligent Agents, Web Technologies and Internet Commerce, Nov 2005.