# International Journal of Advances in Computer Science and Technology

# WAR MINING SVM BASED VISUALIZED PATTERN DISCOVERY IN MIXED DATA

**T.TamilSelvi [1], Mrs.P.S.AnnaKKodi [2]**

[1]M.Phil Scholar, Department of Computer Science, Sri Ramalinga Sowdambigai College of Science and Commerce, Vadavalli, Coimbatore-641109, Tamilnadu, India, tamil_selvi21@rediff.com.

[2]Head of the Department , Department of Information Technology , Sri Ramalinga Sowdambigai College of Science and Commerce , Vadavalli, Coimbatore-641109, Tamilnadu, India, annakkodi.ps @gmail.com

## ABSTRACT

This paper provides geometric information regarding how significant to human beings are cultured in the every part of the country and their qualifications; it affords geometric information on the security. This description is on paper in reply to frequent requests for war object statistics and listing of war dead. This description in addition cites basis of available listing of armed personnel killed in major wars and conflict proceedings. Numerous conventional classification algorithms in inductive learning knowledge use suspiciously considered categorical data spaces. They essentially cannot successfully handle continuous attributes straightforwardly. To conquer these problems proposed an attribute SVM based learning result. In this war we examination of mixed-mode data in Attribute Clustering Algorithm system. This algorithm has been provided that the result designed for this examination. The make use of numerous modes is of the majority concentration for the reason that mixed-mode scheme through unimode data gathering. SVM algorithm that trains a grouping algorithm by adjusts the item-pair correspondence quantify. The algorithm might optimize a diversity of dissimilar clustering purpose to an assortment of clustering performance procedures. Databases, particularly data warehouses and sequential databases can be converted into moderately huge. The effectiveness and usability of these databases extremely depend on how speedily data can be regained. The mixing of modes and move toward appear to be inadequate at period only by the inspiration of investigation managers. On the further hand, scheme through other than single manner of data gathering might raise the likelihood of dimension inaccuracy since the examination difficulty might show to some extent another way under diverse modes.

**Keywords:** War Mining, Department of Defense (DOD), Clustering, SVM algorithm.

## 1. INTRODUCTION

The war provides geometric information on the security information of a countryside i.e., numeral of wars take position in a specific country meant for the decade. Within the applicant datasets themselves, every individual war applicant is describing in conditions of the day with the aim of it go through and left the war. The key purpose in the categorization of wars was found on the decision of who was aggressive whom. This ranges beginning the amount of times war takes places, challenger country, year that occurred, place and as well death occurred in a particular war. In adding together it moreover present the number of conquest information regarding the country and their opponent's success.

War is a well thought-out, equipped, and frequently an extended disagreement that is approved on among sates, nation or further parties. War is supposed to be understood as a concrete, intended and extensive equipped inconsistency among supporting communities, and consequently is distinct as a type of political violence. The set of techniques second-hand by a collection to take out war is well-known as warfare. A deficiency of war is regularly called peace. War, to happen to known as one, necessity necessitate various amount of disagreement by means of weapons and further military technology and equipment through armed forces utilize military tactics and operational art inside the extensive operation topic to military logistics. War studies by martial theorists all the way through military history include required to recognize the philosophy of war, and to decrease it to a military science.

Conventional warfare is a challenge to decrease an adversary military aptitude all the way through release battle. It is a confirmed war among existing states only sees imperfect exploitation in sustain of conservative military goals and military exercises. Nuclear warfare is type of the war, where nuclear weapons are the major technique of forcing the submission of the additional side, as contrasting to a sustaining premeditated or calculated role in a conservative difference. Civil war is the type of war wherever the military in disagreement be in the right place to the similar country and is vie for manage of or self-government beginning that country. Asymmetric warfare is a disagreement among two populations of considerably dissimilar levels of armed capability. Asymmetric disagreement frequently results in guerrilla tactics person worn to conquer the occasionally huge gaps in knowledge and force size.

Support Vector Machines (SVM) method is individual of the majority authoritative classification method that was effectively practical to numerous real world problems [1]. Numerous applications require the dispensation of huge data

sets. The training time of SVM classification is a severe obstruction on behalf of this type of data sets. According to [2], it would capture years to train SVM on a data set. Numerous suggestions have been current to improve SVM to amplify its training performance [3], through estimate of the secondary classifier. Conversely, they are still not reasonable with large data sets where even numerous scans of complete information set are in addition exclusive to achieve, or they end up behind the remuneration by means of an SVM through generalization [4]. In this paper primary illustrate a new approach to preparation SVM among very big datasets which is dissimilar from the three major approaches discussed in [5].

The approach is base on the characteristic of SVM so as to only preparation data near the separating boundary is important. Therefore current a cluster method that yields just a few clusters away since the separating boundary as well as many clusters close to the boundary. This way the main information beginning the training data - specifically so as to of the preparation information close to the separating boundary is conserved whereas at the same time the range of the preparation set is efficiently decreased.

Just the once cluster have been originate, they are represent have a set of biased example which as well as the SVM preserve be present basically comprehensive to contract through the learning and review of the American Civil War have usually focused on top of the leading military officers in addition to most important battle and not an logical look on the War from the perspective of the person soldier. Use information since every soldier's forces as well as civilian experience to assemble a database from "ground up" rather than "top down", Historical information system have produced the only database of its category to facilitate can be used for algebraic and systematic examination of the War. It is at the present possible to check and measure the impact this individual soldier experience had upon regimental usefulness.

## 2. RELATED WORK

A record contains information by together continuous and definite value is called as "great mixed-mode database". In a broader exposure, the information things in the database possibly will be of an planned nature, such as a actual or an numeral rate, or ranking which might be represent as integers, or of an unordered separate nature, such as definite items prepared up of symbols, conditions, and/or interval. While it is not possible to change unordered separate data into continuous data, in mainly convenient applications, continuous records unusually will be changed into gap data, so as to every the data items in a mixed-mode information could be process as separate events, to provide a uniform support for experience pattern, organization, and rule detection tasks.

For past reason, the majority of the categorization algorithms in the mechanism learning area can just be used for categorical or nominal databases. The majority of these traditional classification algorithms are straightforwardly to handle neither databases together with continuous values [6] nor mixed-mode database straight and successfully. Though, in the actual world, a great part of data really does include either continuous or categorical ideals referred to as "mixed-mode" values.

Having the present existing inductive learn system been straightforwardly apply to these kind of mixed-mode database starting the actual world as well as all of the continuous values contained by the mixed-mode databases. newly, several researchers as well have establish to even if some learning systems are clearly considered and build for continuous attribute or databases, these systems still might manage a advanced precision than unrefined databases if continuous data are correctly discretized. As a reasonable end result, the restrictions of mainly inductive learning procedure algorithms will be beat by discretizing all permanent attributes accurately, before feed those datasets into the current knowledge system [7].

Essentially, any discretization can be thinking of as a pre-process with which partition the rate space of a permanent quality into a limited number of intervals, as well as at the same time, allocate a nominal rate to all of them. In this proposal, the investigator has describe a completely new technique for discretizing the ideals of permanent attributes inside a mixed-mode database, which is fully based on an information quantity that accurately reflects the relationship between the permanent attributes and the category attributes [8]. Conventionally, two main factors must be in use into thoughtfulness in the pre-processing phase of partitioning a permanent data space the amount of intervals, and the thickness of every interval.

As a good quality algorithm, it must usually require as a small number of inputs from the user as possible. In an exact categorization assignment, several presented class information, since the real world or domain expert, can be of essential importance in the discretization procedure [9]. The literature on discretization methods is profuse; however the majority of them are observed as univariate technique. In [10], developed a multivariate discretization schema a particular inconsistent determination be considered at the similar instance. Though, as a multivariate discretization, it is obtainable to develop into a most important way in data discretization for pattern examination.

## 3. SVM BASED VISUALIZED PATTERN DISCOVERY

Several presently-existing learn mechanism systems have been suspiciously built for processing definite attributes values. In the area of mechanism learning or data mining these inductive learning troubles are typically planned to determine classificatory pattern, or just regulations, based on a set of information samples. Categorization ruler and/or patterns are produced for individual's pre-labeled data

sample with certain preceding information set up by area experts in those areas. Thus, several traditional categorization algorithms in inductive learning make use of carefully planned categorical data spaces. They really cannot successfully handle permanent attributes straightforwardly. To apply inductive learning systems with these kinds of mixed-mode data space, the permanent variables are required to first be discretized.

However, for a mixed-mode database present is still no excellent result to the unsubstantiated knowledge task, or to clustering extremely great mixed-mode databases have even larger challenges.  It provide tables, compiled by sources at the Department of Defense (DOD), representative the number of casualties surrounded by American military personnel serving in principle wars and combat events.

The common purpose of this dissertation is to assemble this challenge. It attempt to answer the majority basic problems, primary of partitioning extremely great mixed-mode databases into less significant coherent ones, as well as then discretizing any permanent information without relying on explicit class labels. These dissertations propose a technique to answer this testing problem. In solving discretization problems, two issues have been raised. The first one is to facilitate if a leading  element actually exists, possibly will use it to make the discretization of all permanent attributes. The next one is concerning the state of the interdependence relationship among the attributes in the subgroup.

 For a extremely great mixed-mode database, unless a class label is specified, present is no motive to consider that the complete database is complete happy of a particular interrelated collection, or that it is governed by a particular attribute.  In fact, present could be some interrelated attribute groups co-existing essentially in the data set, each possibly will share more associated information among themselves than with others; thus it is not meaningful to use the method of the complete data set to take the discretization.

In outlook of this, an additional reasonable approach is to first analyze whether the database may possibly be optimally partition into a number of reasoned attribute groups or not, before discretization is functional to the complete collection or to every of the clustered groups. This is a significant notion to be exploring by this theory.  After the mixed-mode database partitioning and discretization troubles are solve as discuss above, a large mixed-mode database can be changed into a number of smaller databases, all of which possibly will contain discrete valued data, or rear into a large mixed-mode database by combine different sub-databases with discrete-valued databases. The purpose of this thesis is motivated by such realistic requirements beginning the real world. Since the majority of the information is from a diverse basis, several of them possibly will not have explicit class information. Then need pattern finding methods which may not automatically rely on class information

The learning would take in all potential to recognize the newer attributes related to war/ production enable to strategy executive to enclose the war and to encourage the cultivation. The combination of mode and approach seems to be limited at times only by the creativeness of survey manager. For example, definite survey behavior, such as prenotification and reminder messages, can be sent using one mode, but data collection might rely on a special, but particular mode or the examination data can be collect using additional than one mode

Supervise clustering is the trouble of preparation a clustering algorithm to generate desirable clusterings: specified set of items and whole clustering's over these sets; find out how to cluster future sets of items. Example applications incorporate noun-phrase co reference clustering, and clustering news article by whether they pass on to the similar topic. Support Vector Machines are based on the suggestion of map data point to a high dimensional attribute space where a separating hyper-plane can be originated. This mapping can be accepted on by apply the kernel tricks which implicitly transform the input space into an additional high dimensional attribute space. The hyper-plane is computed by maximize the expanse of the contiguous patterns. The separating hyper-plane is establishing use a quadric training routine which is computationally extremely exclusive. Furthermore, this routine based on the data set range, taking not practical time when discuss with small data sets. This paper proposes a new advance for attractive the preparation process of SVM when dealing with large data sets. It is based on the mixture of SVM and clustering investigation.

The idea is as follows: SVM compute the maximal margin separating information points; therefore, only those patterns neighboring to the border can have an effect on the computation of that border, while extra point can be discarded without affect the end answer. Those points lying close to the border are called maintain vectors. Partitioning, also called flat clustering, straight seeks a separation of the data which optimizes a predefined arithmetical measure. In partitioning clustering, the many clusters are predefined, and determining the most favorable number of clusters may involve extra computational cost than clustering itself.

Furthermore, a priori knowledge may be essential for initialization and the avoidance of local minima. Hierarchical clustering, on the additional hand, does not need a predefined number of clusters or a priori knowledge. Hence, support hierarchical clustering over flat clustering. In the course of our approach, do not use the creative data set to train SVM; instead, use the tree node references generate by the clustering algorithm. Since the size of the hierarchical tree is considerably less than the size of the compare to the original data set, the preparation process will be extremely speeding. Then produce a hierarchical clustering tree for every class up to a certain level. Then use the nodes' references in the better levels of both trees to train SVM. After preparation, compute the maintain vector references and the accurateness of the classifier.

If the accurateness is not fulfilled, de-cluster those maintain vector reference by addition their children to the preparation set. Second, include single additional search step to establish the maintain vectors by measuring the distance between nodes from both trees. This step will exclude distant nodes from equally trees, base on a certain threshold. Third, make only a single hierarchical clustering tree and decide the mainly experienced nodes to train and de-cluster nodes base on the heterogeneity of nodes. Heterogeneous nodes are individual's nodes so as to have data points assign to them since special classes.

The main contributions of this work are as follows: First, to reduce the training time of SVM, we suggest a new maintain vector selection method using clustering analysis. Second, we plan different support vector collection strategies base on combine the de-clustering (expansion of cluster) and SVM preparation phases. The algorithm computes the cluster of the set of points in every one class. It illustrates it in the case of class 1.

**Algorithm war mining (SVM for visual pattern war discovery)**

**Input:** Number of the training samples with war data
$x = \{x_1, x_2, \ldots x_n\}$ as input file for SVM classification

**Output:** Classification result along with cluster mode

Procedure ( SVM X) // input training data results from the SVM methods

Begin

Begin

Initialize the set $A_1$ of clusters of class 1

$A_1 = \{(x_i, 1) | (x_i, y_i) \in S, y_i = 1\}$

C=0

Get input file X for training

Read the number of war data $x = \{x_1, x_2, \ldots x_n\}$ as input file

For each point $(x_k, n_k) \in A$

Compute the nearest point in $A_1$ $(x_1, n_1)$ to $(x_k, n_k)$.

Let $(x_j, n_j)$ be this point and d their distance, $d = |x_k - j_k|$

Compute the center of mass of the two previous points,

$v = \frac{n_k x_k + n_j x_j}{n_k + n_j}$

Compute the distance $D$ between $v$ and the nearest training point in class $-1$.

If $\frac{d}{D} < \gamma$ delete the two previous points from $A_1$, add

$(v, n_k + n_j)$ to $A_1$, and set $C = C + 1$.

$x_i. w + b = 0$ // war data is represented as matrix and denoted by $x_i$ and w is the weight value matrix whose product is summed with bias value to give the class value.

$x_i. w + b = 1$

Decision function $f(W) = x_i. w - b$ //decision function f(w) decides the class labels for the SVM classification training examples ,

If $f(W) \geq 1$ for $x_i$ is the first class // if the F(w) is greater than or equal to the 1 is labeled as first class

Else

$f(W) \leq -1$ for $x_i$ is the first class // if the f(w) is less than or equal to the value of -1 is labeled as second class

If $C > 0$ goto step 1, otherwise stop.

The prediction result for $(i = 1, \ldots n)$ number of document

Check the result using the following function

$y_i(x_i. w - b) \geq 1$ //if the function is greater than one the results clustered as formed and patterns are predicted

Finally display values of the patterns

After the **SVM for visual pattern war discovery** algorithm has stop, $A_1$ is the group of clusters of class 1 altogether with mixed modes of attributes. The similar process is after that frequent to calculate the group of data clusters of class 2, $A_2$. Become aware of that the algorithm tends to generate large clusters of position which are far away beginning the state line among the two classes and small clusters of points subsequently to the boundary. Thus, suppose with the purpose of the points contestant to be support vectors are not robustly affected by the clustering process, whereas the others are seriously concentrated. The constraint $\gamma$ controls the significance of "near". If the dimensionality of the contribution space is not big ($say\ n \approx 10$), be expecting with the intention of the algorithm can significantly recover the working out time of the SVM.

## 4. EXPERIMENTAL RESULTS

They have predictable the accurateness results of our SVM visual pattern discovery on classification. In this work consider the American Civil War investigation Database is a relational database. This way that there are various files which are "associated" to every other. Through HDS Database the association among these numerous files is the soldier's name. These files include information get together from the dissimilar sources second-hand as converse below. Not each soldier is in each file, nor is the category of information regarding a soldier in exacting files the related for each further soldier in that comparable file. Moreover, the pace at which at HDS can enter the information regarding each soldier determination differ considerably based on the category of source documents obtainable to us. Consequently, the Database is not an absolute biographical story on each soldier, but somewhat a source of imperative war and position war events.

Since this exposition suggest a new approach to deal with the detection of patterns for mixed-mode data, should design suitable experimentation to substantiate the building and make known how sensible the proposed approach when useful to a variety of types of mixed-mode data. In this section effort primary to intend a set of experimentation through preferred data of a variety of types to test our property. Subsequently determination is appropriate to two large sets of valid world databases which are difficult; do not include class labels however are backed by sufficient field information for confirmation of the investigative outcome. Additional numeral of characteristic is obtainable in obtainable algorithms, it's boring to compute the accurateness assessment for every and each attributes, implementation time is elevated and it precedes the low altitude of accurateness. The assessment for war cannot be designed precisely. In SVM algorithm, the attributes accurateness value can be determine using testing & training

methods, it's simple to compute accurateness value for every and each attributes, carrying out time also small. Somewhat, with obtainable algorithms can discover the precise value in elevated level. The subject statistics for wars extended ended is updated occasionally, sometimes annually. This approximately forever replicate the classification of leftovers persons earlier planned as misplaced in action and persons reclassification as dead. Additional reasons, greatly rarer, consist of the finding of errors in subject records for persons of people. The result of the death percentage of the war was shown in Figure 1, Figure 2 shows the War about death and normal details and Figure 3 shows the No of death and wounded in war results. The following Table 1 shows the details of war and Table 2 What purpose the war was executed and mined details.

**Table 1:** War details

| Year | Total Deaths | Accidents | Hostile Action | Homicide | Illness | Self Inflicted | Terrorist Attack | Undetermined |
|------|------|------|------|------|------|------|------|------|
| 1980 | 2392 | 1556 |    | 174 | 419 | 231 | 1   | 11 |
| 1981 | 2380 | 1524 |    | 145 | 457 | 241 |     | 13 |
| 1982 | 2319 | 1493 |    | 108 | 446 | 245 | 2   | 16 |
| 1983 | 2465 | 1413 | 18 | 115 | 419 | 218 | 263 | 19 |
| 1984 | 1999 | 1293 | 1  | 84  | 374 | 225 | 6   | 16 |
| 1985 | 2252 | 1476 |    | 111 | 363 | 275 | 5   | 22 |
| 1986 | 1984 | 1199 | 2  | 103 | 384 | 269 |     | 27 |
| 1987 | 1983 | 1172 | 37 | 104 | 383 | 260 | 2   | 25 |
| 1988 | 1819 | 1080 |    | 90  | 321 | 285 | 17  | 26 |
| 1989 | 1636 | 1000 | 23 | 58  | 294 | 224 |     | 37 |
| 1990 | 1507 | 880  |    | 74  | 277 | 232 | 1   | 43 |
| 1991 | 1787 | 931  | 147| 112 | 308 | 256 |     | 33 |
| 1992 | 1293 | 676  |    | 109 | 252 | 238 | 1   | 17 |

**Table 2:** What purpose the war was executed

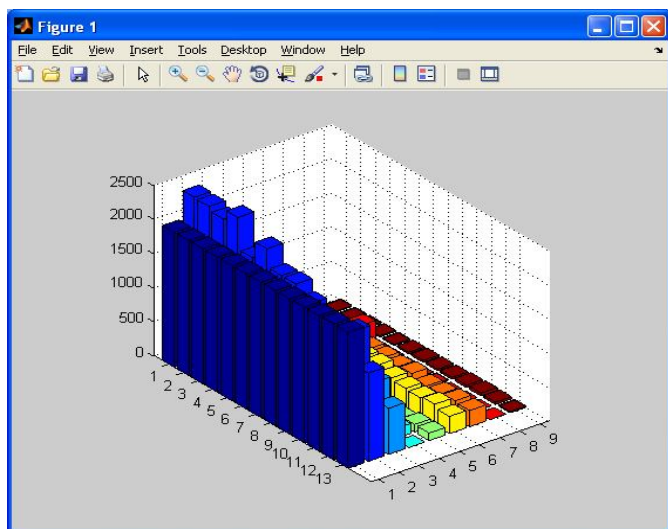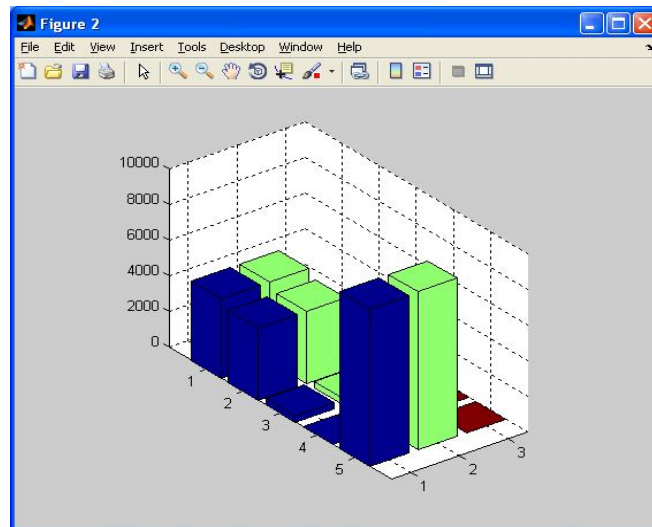| Operation | Deaths(A) | Wounded(B) | Ratio(A/B) |
|------|------|------|------|
| Iraqi Freedom | 4301 | 31430 | 1:7/3 |
| Enduring Freedom | 714 | 3162 | 1:4/4 |
| Persian Gulf | 383 | 467 | 1:1/2 |
| Vietnam | 58220 | 153303 | 1:2/6 |
| Korea | 36574 | 103284 | 1:2/8 |
| W.W-II | 405399 | 670846 | 1:1/65 |

Figure 1: Death percentage in war



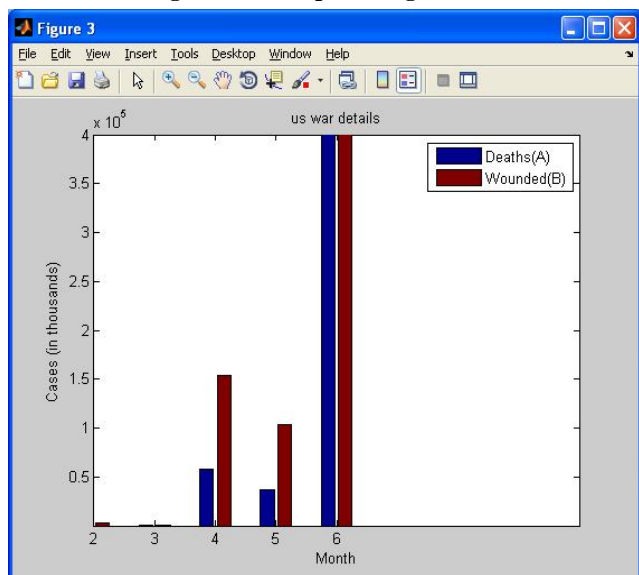Figure 2: War about death and normal details



Figure 3: No of death and wounded in war

## 5.  CONCLUSION AND FUTURE WORK

In this paper proposed an mixed mode based data analysis for effective pattern discovery methods in war mining data for improved appreciative; The reliability and the representative individuality of every of the meteorological modes revealed recommend with the intention of definite modes might provide as suggestion parameters as they renders a great deal extra accurate evaluation of the war mining . The discovery of patterns and grouping of consideration patterns in interruption process revealing scheme task and prepared distinctiveness. Its method is essentially controls the heat circumstance and serves as a trigging aspect to make active the disaster liberate reaction.

Such conclusion demonstrates the helpfulness and usefulness of the anticipated method in illuminating subtle function patterns for structure examination, manage and optimization.  In this paper, useful reduction procedure by means of clustering investigation to estimated support vectors in regulate to rapidity the training procedure of SVM. The SVM approach demonstrates to better every one other technique in expressions of accurateness, but it takes additional instance. As the training sets get better, throughout de-clustering, the inaccuracy rates go down. Though, the error rates resolve on precise stage even if the training set sizes enlarge.

There are numerous issues with the intention of be able to be added investigated in take apart study directions. Primary, an incremental version of the SVM can be anticipated to handle incremental sources of information. Subsequent, the selection of kernel function in mapping data points to characteristic freedom can be added investigated to be second-hand in both clustering examination and SVM. Lastly, other methods, such as Randomized algorithms, can be used in similar to support vectors.

### REFERENCES

1. LeCun, Y., Jackel, L. D., Bottou, L., Brunot, A., Cortes, C., Deker, J.S., Drucker, H., Guyon, I., Muller, U.A.,  Sackinger, E., Siard, P. and Vapnik, V., "**Comparison of Learning algorithms for handwritten recognition**", *International Conference on Artificial Neural Networks*, Fogelman, F. and Gallinari, P. (Ed.), pp. 53-60,1995.
2. Hwanjo Yu, Jiong Yang, and Jiawei Han, **"Classifying Large Data sets Using SVM with Hierarchical Clusters"**, *SIGKDD '03*, August 24-27, 2003, Washington, DC, USA.

3. A. Ben-Hur, D. Horn, H.T. Siegelmann and V. Vapnik**," Support Vector Clustering",** *Journal of Machine Learning Research* 2:125-137, 2001.

4. Ying Zhao and George Karypis., **"Prediction of Contact Maps Using Support Vector Machines",** *Third IEEE Symposium on Bioinformatics and Bio-Engineering (BIBE'03)* March 10 - 12, 2003 Bethesda, Maryland

5. G.Cauwenberghs and T.Poggio. **"Incremental and decremental support vector machine learning"**. *In Proc. Advances in Neural Information Processing Systems.* Vancouver, Canada, 2000.

6. Michael P. S. Brownz, William Noble Grundyz, David Lin, Nello Cristianini, Charles Sugnet, Manuel Ares Jr., David Hausslerz, **"Support Vector Machine Classification of Microarray Gene Expression Data"** , *technical report UCSC-CRL-99-09,* University of California, Santa Cruz.

7. G Li and AKC Wong, **"Pattern Distance Measures in Categorical Data for Pattern Pruning and Clustering",** *submitted to IEEE Trans. on Knowledge and Data Engineering.*

8. WH Au, KCC Chan, AKC Wong and Y Wang, **"Attribute Clustering for Grouping, Selection, and Classification of Gene Expression Data"** *IEEE/ACM Trans on Computational Biology and Bioinformatics,* Vol 2, No2, pp 83-101, 2005.

9. AKC Wong, WH Au and KCC Chan, **"Discovering High-Order Patterns of Gene Expression Levels",** *Journal of Computational Biology,* Vol. 15, No.6, 2008. revision, 2008

10. Gauthier, P. and Possamai, D.**,"** **Efficient simulation of the Wishart model",** *SSRN eLibrary,*2009.