

A Survey On Feature Selection Algorithm For High Dimensional Data Using Fuzzy Logic

T.Jaga Priya Vathana,¹ C.Saravanabhavan², Dr.J.Vellingiri³

¹ M.E-Student, Department of CSE, Kongunadu College of Engineering and Technology, Tamil Nadu, India.

² Research Scholar & Asst. Professor, Department of CSE, Kongunadu College of Engineering and Technology, Tamil Nadu, India.

³ Associate Professor, Kongunadu College of Engineering and Technology, Tamil Nadu, India.

-----ABSTRACT-----

Feature subset selection is an effective way for reducing dimensionality, removing irrelevant data, increasing learning accuracy and improving results comprehensibility. This process improved by cluster based FAST Algorithm and Fuzzy Logic. FAST Algorithm can be used to Identify and removing the irrelevant data set. This algorithm process implements using two different steps that is graph theoretic clustering methods and representative feature cluster is selected. Feature subset selection research has focused on searching for relevant features. The proposed fuzzy logic has focused on minimized redundant data set and improves the feature subset accuracy.

Date of Submission: 13, September, 2013



Date of Acceptance: 10, October 2013

I. INTRODUCTION

The performance, robustness, and usefulness of classification algorithms are improved when relatively few features are involved in the classification. Thus, selecting relevant features for the construction of classifiers has received a great deal of attention. With the aim of choosing a subset of good features with respect to the target concepts, feature subset selection is an effective way for reducing dimensionality, removing irrelevant data, increasing learning accuracy, and improving result comprehensibility. Many feature subset selection methods have been proposed and studied for machine learning applications. They can be divided into four broad categories: the Embedded, Wrapper, Filter, and Hybrid approaches. The embedded methods incorporate feature selection as a part of the training process and are usually specific to given learning algorithms, and therefore may be more efficient than the other three categories. Traditional machine learning algorithms like decision trees or artificial neural networks are examples of embedded approaches. The wrapper methods use the predictive accuracy of a predetermined learning algorithm to determine the goodness of the selected sub-sets, the accuracy of the learning algorithms is usually high. However, the generality of the selected features is limited and the computational complexity is large. The filter methods are independent of learning algorithms, with good generality.

With respect to the filter feature selection methods, the application of cluster analysis has been demonstrated to be more effective than traditional feature selection algorithms. Pereira et al. , Baker et al. , and Dillon et al. employed the distributional clustering of words to reduce the dimensionality of text data. In cluster analysis, graph-theoretic methods have been well studied and used in many applications. Their results have, sometimes, the best agreement with human performance. The general graph-theoretic clustering is simple: Compute a neighborhood graph of in-stances, then delete any edge in the graph that is much longer/shorter (according to some criterion) than its neighbors. The result is a forest and each tree in the forest represents a cluster. In our study, we apply graph-theoretic clustering methods to features. In particular, we adopt the minimum spanning tree (MST) based clustering algorithms, because they do not assume that data points are grouped around centers or separated by a regular geometric curve and have been widely used in practice. Based on the MST method, we propose a FAST clustering-Based feature Selection algorithm (FAST). The FAST algorithm works in two steps. In the first step, features are divided into clusters by using graph-theoretic clustering methods. In the second step, the most representative feature that is strongly related to target classes is selected from each cluster to form the final subset of features. Features in different clusters are relatively independent; the clustering-based strategy of FAST has a high probability of producing a subset of useful and independent features. The proposed feature subset se-lection algorithm FAST was tested upon 35 publicly

available image, microarray, and text data sets. The Experimental results show that, compared with other five different types of feature subset selection algorithms, the proposed algorithm not only reduces the number of features, but also improves the performances of the four well-known different types of classifiers.

II. LITERATURE REVIEW

2.1 Statistical Comparisons of Classifiers over Multiple Data Sets

In this method introduce some new pre- or post processing step has been proposed, and the implicit hypothesis is made that such an enhancement yields an improved performance over the existing classification algorithm. Alternatively, various solutions to a problem are proposed and the goal is to tell the successful from the failed. A number of test data sets is selected for testing, the algorithms are run and the quality of the resulting models is evaluated using an appropriate measure, most commonly classification accuracy. The remaining step, and the topic of this paper, is to statistically verify the hypothesis of improved performance. Various re-researchers have addressed the problem of comparing two classifiers on a single data set and proposed several solutions. The core of the paper is the study of the statistical tests that could be (or already are) used for comparing two or more classifiers on multiple data sets. Learning algorithms is used for the Classification purpose. The main disadvantage of this process is the problems with the multiple data set tests are quite different, even in a sense complementary.

2.2 A Features Set Measure Based On Relief

It used six real world dataset from the UCI repository have been used. Three of them have classification Problem with discrete features, the next two classifications with discrete and continuous features, and the last one is approximation problem. The learning algorithm is used to check the quality of feature selected are a classification and regression tree layer with pruning. This process and algorithms is implemented by the orange data mining System. Overall, the non-parametric tests, namely the Wilcoxon and Friedman test are suitable for our problems. They are appropriate since they assume some, but limited commensurability. They are safer than parametric tests since they do not assume normal distributions or homogeneity of variance. There is an alternative opinion among statisticians that significance tests should not be performed at all since they are often misused, either due to misinterpretation or by putting too much stress on their results. The main disadvantage of the system is its measure to low accuracy of the search process.

2.3 Feature Clustering and Mutual Information for the Selection of Variables In Spectral Data

It face many problems in spectrometry require predicting a quantitative value from measured spectra. The major issue with spectrometric data is their functional nature; they are functions discretized with a high resolution. This leads to a large number of highly-correlated features; many of which are irrelevant for the prediction. The approach for the features is to describe the spectra in a functional basis whose basis functions are local in the sense that they correspond to well-defined portions of the spectra. This process has clustering algorithm that algorithm recursively merges at each step the two most similar consecutive clusters. This algorithm return the output value associated with each cluster, its representative, is chosen to be the mean of the spectra over the range of features defined by the cluster. The main disadvantage of the problem is low number of clusters identified by the method allows the interpretation of the selected variables: several of the selected clusters include the spectral variables identified on these benchmarks as meaningful in the literature.

2.4 On Feature Selection through Clustering

This paper introduce an algorithm for feature selection that clusters attributes using a special metric and, then uses a hierarchical clustering for feature selection. Hierarchical algorithms generate clusters that are placed in a cluster tree, which is commonly known as a dendrogram. Clustering's are obtained by extracting those clusters that are situated at a given height in this tree. It use several data sets from the UCI dataset repository and, due to space limitations we discuss only the results obtained with the votes and zoo datasets, Bayes algorithms of the WEKA package were used for constructing classifiers on data sets obtained by projecting the initial data sets on the sets of representative attributes. Approach to attribute selection is the possibility of the supervision of the process allowing the user to opt between quasi-equivalent attributes. It face classification problems that involve thousands of features and relatively few examples came to the fore. We intend to apply our techniques to this type of data.

III. FUZZY BASED FEATURE SUBSET SELECTION ALGORITHMS

Irrelevant features, along with redundant features, severely affect the accuracy of the learning machines. Thus, feature subset selection should be able to identify and remove as much of the irrelevant and redundant information as possible. The cluster indexing and document assignments are repeated periodically to compensate churn and to maintain an up-to-date clustering solution. The k-means clustering technique and SPSS

Tool to develop a real time and online system for a particular supermarket to predict sales in various annual seasonal cycles. The classification was based on nearest mean.

In order to more precisely introduce the algorithm, and because our proposed feature subset selection framework involves irrelevant feature removal and redundant feature elimination.

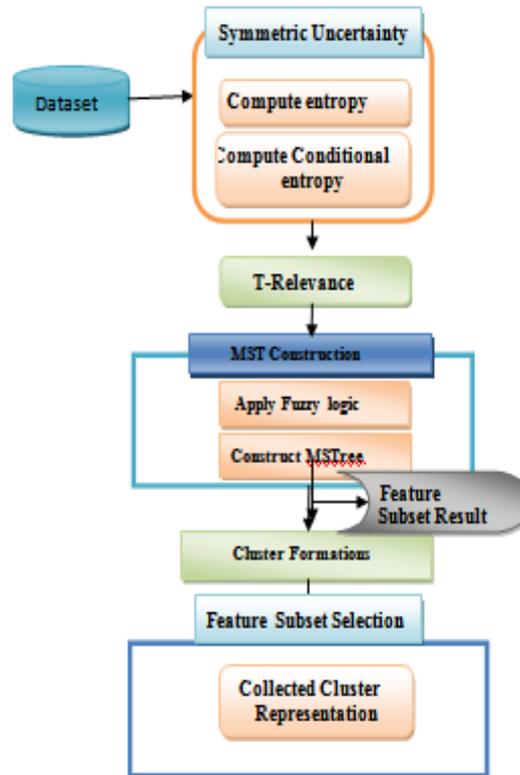


Fig. 1: Framework of the Fuzzy Based

3.1. Feature subset selection algorithm

Irrelevant features, along with redundant features, severely affect the accuracy of the learning machines. Thus, feature subset selection should be able to identify and Remove as much of the irrelevant and redundant information as possible. Moreover, “good feature subsets contain features highly correlated with (predictive of) the class, yet uncorrelated with (not predictive of) each other. Keeping these in mind, we develop a novel algorithm which can efficiently and effectively deal with both irrelevant and redundant features, and obtain a good feature subset. We achieve this through a new feature selection framework which composed of the two connected components of irrelevant feature removal and redundant feature elimination. The former obtains features relevant to the target concept by eliminating irrelevant ones, and the latter removes redundant features from relevant ones via choosing representatives from different feature clusters, and thus produces the final subset. The irrelevant feature removal is straightforward once the right relevance measure is defined or selected, while the redundant feature elimination is a bit of sophisticated. In our proposed FAST algorithm, it involves (a) the construction of the minimum spanning tree (MST) from a weighted complete graph; (b) the partitioning of the MST into a forest with each tree representing a cluster; and (c) the selection of representative features from the clusters. In order to more precisely introduce the algorithm, and because our proposed feature subset selection framework involves irrelevant feature removal and redundant feature elimination, we firstly present the traditional definitions of relevant and redundant features, then provide our definitions based on variable correlation as follows. John et al. presented a definition of relevant features. Suppose F to be the full set of features, $i \in F$ be a feature, $S_i = F - \{F_i\}$ and $S' \subseteq S_i$. Let s'_i be a value-assignment of all features in S'_i , f_i a value-assignment of feature F_i , and c a value-assignment of the target concept C . The definition can be formalized as follows. Definition: (Relevant feature) F_i is relevant to the target concept C if and only if there exists some s'_i , f_i and c , such that, for probability $(S'_i = s'_i, F_i = f_i) > 0, p(C=c | S'_i = s'_i, F_i = f_i) \neq p(C=c | S'_i = s'_i)$. Otherwise, feature F_i is an irrelevant feature. Definition 1 indicates that there are two kinds of relevant features due to different S'_i : (i) when $S'_i = S_i$, from the definition we can know that F_i is directly relevant to the target concept;

(ii) when $S' \not\subseteq S$, from the definition we may obtain that $p(C|S_i, Fi) = p(C|S_i)$. It seems that Fi is irrelevant to the target concept. However, the definition shows that feature Fi is relevant when using $S' \cup \{Fi\}$ to describe the target concept. The reason behind is that either Fi is interactive with $S' \cup \{Fi\}$ or Fi is redundant with $S' \cup \{Fi\}$. In this case, we say Fi is indirectly relevant to the target concept. Most of the information contained in redundant features is already present in other features. As a result, redundant features do not contribute to getting better interpreting ability to the target concept. It is formally defined by Yu and Liu based on Markov blanket. The definitions of Markov blanket and redundant feature are introduced as follows, respectively.

let $M_i \subseteq F$ ($Fi \in M_i$), M_i is said to be a Markov blanket for Fi if and only if $p(F - M_i - \{Fi\}, C | Fi, M_i) = p(F - M_i - \{Fi\}, C | M_i)$. Definition: (Redundant feature) Let S be a set of features, a feature in S is redundant if and only if it has a Markov Blanket within S . Relevant features have strong correlation with target concept so are always necessary for a best subset, while redundant features are not because their values are completely correlated with each other. Thus, notions of feature redundancy and feature relevance are normally in terms of feature correlation and feature-target concept correlation. Mutual information measures how much the distribution of the feature values and target classes differ from statistical independence. This is a nonlinear estimation of correlation between feature values or feature values and target classes. The symmetric uncertainty (SU) is derived from the mutual information by normalizing it to the entropies of feature values or feature values and target classes, and has been used to evaluate the goodness of features for classification by a number of researchers (e.g., Hall], Hall and Smith , Yu and Liu , , Zhao and Liu ,). Therefore, we choose symmetric uncertainty as the measure of correlation between either two features or a feature and the target concept

The symmetric uncertainty is defined as follows $S(X, Y) = 2 \times \text{Gain}(X|Y) / (H(X) + H(Y))$.

Where,

1) $H(X)$ is the entropy of a discrete random variable X . Suppose $p(x)$ is the prior probabilities for all values of X , $H(X)$ is defined by $H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$.

2) $\text{Gain}(X|Y)$ is the amount by which the entropy of Y decreases. It reflects the additional information about Y provided by X and is called the information gain which is given by $\text{Gain}(X|Y) = H(Y) - H(Y|X)$.

Where $H(Y|X)$ is the conditional entropy which

Quantifies the remaining entropy (i.e. uncertainty) of a random variable X given that the value of another random variable Y is known. Suppose $p(x)$ is the prior probabilities for all values of X and $p(x|y)$ is the posterior probabilities of X given the values of Y , $H(X|Y)$ is defined by $H(X|Y) = - \sum_{y \in Y} p(y) \sum_{x \in X} p(x|y) \log_2 p(x|y)$.

(4) Information gain is a symmetrical measure. That is the amount of information gained about X after observing Y is equal to the amount of information gained about Y after observing X . This ensures that the order of two variables (e.g., (X, Y) or (Y, X)) will not affect the value of the measure.

Symmetric uncertainty treats a pair of variables symmetrically, it compensates for information gain's bias toward variables with more values and normalizes its value to the range $[0, 1]$. A value 1 of $S(X, Y)$ indicates that knowledge of the value of either one completely predicts the value of the other and the value 0 reveals that X and Y are independent. Although the entropy-based measure handles nominal or discrete variables, they can deal with continuous features as well, if the values are discretized properly in advance

Given $SU(X, Y)$ the symmetric uncertainty of variables X and Y , the relevance T-Relevance between a feature and the target concept C , the correlation F-Correlation between a pair of features, the feature Redundancy F-Redundancy and the representative feature R-Feature of a feature cluster can be defined as follows.

Definition: (T-Relevance) The relevance between the feature $Fi \in F$ and the target concept C is referred to as The T-Relevance of Fi and C , and denoted by $S(Fi, C)$. If $S(Fi, C)$ is greater than a predetermined threshold θ , we say that Fi is a strong T-Relevance feature.

Definition: (F-Correlation) The correlation between any pair of features Fi and Fj ($Fi, j \in F \wedge i \neq j$) is called the F-Correlation of Fi and Fj , and denoted by $SU(Fi, Fj)$.

Let $F_k \subseteq F$ be a cluster of features. if $\exists F_j \in S, SU(F_j, C) \geq SU(F_i, C) \wedge SU(F_i, F_j) > SU(F_i, C)$ is always corrected for each $Fi \in S (i \neq j)$, then Fi are redundant features with respect to the given F_j (i.e. each Fi is a F-Redundancy).

Definition: (R-Feature) A feature $f_i \in F = \{f_1, f_2, \dots, f_n\} (n < |F|)$ is a representative feature of the cluster C (i.e. f_i is a R-Feature) if and only if, $f_i = \operatorname{argmax}_{f_j \in F} (T(f_j, C))$. This means the feature, which has the strongest T-Relevance, can act as a R-Feature for all the features in the cluster. According to the above definitions, feature subset selection can be the process that identifies and retains the strong

- T-Relevance features and selects R-Features from feature clusters. The behind heuristics are that
- 1) Irrelevant features have no/weak correlation with Target concept;
 - 2) Redundant features are assembled in a cluster and a representative feature can be taken out of the Cluster.

IV. ALGORITHM AND ANALYSIS

The proposed FAST algorithm logically consists of three steps:

- (i) removing irrelevant features,
- (ii) constructing a MST from relative ones,
- (iii) Partitioning the MST and selecting Representative features.

For a data set D with n features $F = \{f_1, f_2, \dots, f_n\}$ and class C , we compute the T-Relevance $T(f_i, C)$ value for each feature $f_i (1 \leq i \leq n)$ in the first step.

The features whose $T(f_i, C)$ values are greater than a predefined threshold θ comprise the target-relevant feature subset $F' = \{f'_1, f'_2, \dots, f'_m\} (m \leq n)$.

In the second step, we first calculate the F-Correlation $F(f'_i, f'_j)$ value for each pair of features f'_i and $f'_j (f'_i, f'_j \in F' \wedge i \neq j)$. Then, viewing features f'_i and f'_j as vertices and $F(f'_i, f'_j) (i \neq j)$ as the weight of the edge between vertices f'_i and f'_j , a weighted complete graph $G = (V, E)$ is constructed where $V = \{f'_i \mid f'_i \in F' \wedge i \in [1, m]\}$ and $E = \{(f'_i, f'_j) \mid (f'_i, f'_j \in F' \wedge i, j \in [1, m] \wedge i \neq j)\}$. As symmetric uncertainty is symmetric further the F-Correlation $F(f'_i, f'_j)$ is symmetric as well, thus G is an undirected graph.

The complete graph G reflects the correlations among all the target-relevant features. Unfortunately, graph G has m vertices and $m(m-1)/2$ edges. For high dimensional data, it is heavily dense and the edges with different weights are strongly interwoven. Moreover, the decomposition of complete graph is NP-hard. Thus for graph G , we build a MST, which connects all vertices such that the sum of the weights of the edges is the minimum, using the well-known Prim algorithm. The weight of edge (f'_i, f'_j) is F-Correlation $F(f'_i, f'_j)$. After building the MST, in the third step, we first remove the edges $E = \{(f'_i, f'_j) \mid f'_i, f'_j \in F' \wedge i, j \in [1, m] \wedge i \neq j\}$, whose weights are smaller than both of the T-Relevance $T(f'_i, C)$ and $T(f'_j, C)$, from the MST. Each deletion results in two disconnected trees T_1 and T_2 .

Assuming the set of vertices in any one of the final trees to be $V(T)$, we have the property that for each pair of vertices $(f'_i, f'_j \in V(T))$, $T(f'_i, C) \geq T(f'_j, C) \vee T(f'_j, C) \geq T(f'_i, C)$ always holds. From Definition 6 we know that this property guarantees the features in $V(T)$ are redundant. This can be illustrated by an example. Suppose the MST shown in Fig.2 is generated from a complete graph G . In order to cluster the features, we first traverse all the six edges, and then decide to remove the edge (f_0, f_4) . Take $V(T_1)$ as an example.

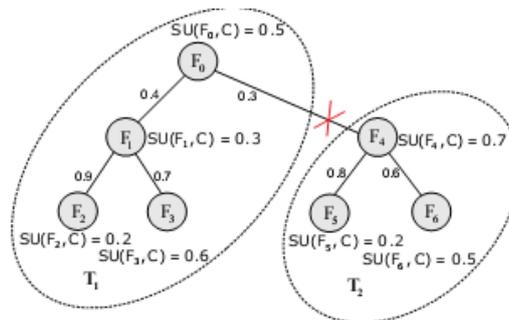


Fig 2. Example of Clustering Step

After removing all the unnecessary edges, a forest Forest is obtained. Each tree $T_i \in \text{Forest}$ represents a cluster that is denoted as (T_i) , which is the vertex set of T_i as well. As illustrated above, the features in each cluster are redundant.

The details of the FAST algorithm is shown in Algorithm 1.

```

Algorithm 1: FAST


---


inputs:  $D(F_1, F_2, \dots, F_m, C)$  - the given data set
            $\theta$  - the T-Relevance threshold.
output:  $S$  - selected feature subset.
//==== Part 1: Irrelevant Feature Removal ====
1 for  $i = 1$  to  $m$  do
2   T-Relevance =  $SU(F_i, C)$ 
3   if T-Relevance  $> \theta$  then
4      $S = S \cup \{F_i\}$ ;
//==== Part 2: Minimum Spanning Tree Construction ====
5  $G = \text{NULL}$ ; //G is a complete graph
6 for each pair of features  $\{F'_i, F'_j\} \subset S$  do
7   F-Correlation =  $SU(F'_i, F'_j)$ 
8   Add  $F'_i$  and/or  $F'_j$  to  $G$  with F-Correlation as the weight of
   the corresponding edge;
9  $\text{minSpanTree} = \text{Prim}(G)$ ; //Using Prim Algorithm to generate the
   minimum spanning tree
//==== Part 3: Tree Partition and Representative Feature Selection ====
10  $\text{Forest} = \text{minSpanTree}$ 
11 for each edge  $E_{ij} \in \text{Forest}$  do
12   if  $SU(F'_i, F'_j) < SU(F'_i, C) \wedge SU(F'_i, F'_j) < SU(F'_j, C)$  then
13      $\text{Forest} = \text{Forest} - E_{ij}$ 
14  $S = \emptyset$ 
15 for each tree  $T_i \in \text{Forest}$  do
16    $F'_k = \text{argmax}_{F'_k \in T_i} SU(F'_k, C)$ 
17    $S = S \cup \{F'_k\}$ ;
18 return  $S$ 


---



```

Time complexity analysis. The major amount of

work for Algorithm 1 involves the computation of SU values for T-Relevance and F-Correlation, which has linear complexity in terms of the number of instances in a given data set. The first part of the algorithm has a linear time complexity ($O(m)$) in terms of the number of features m . Assuming $(1 \leq \alpha \leq m)$ features are selected as relevant ones in the first part, when $\alpha = 1$, only one feature is selected. Thus, there is no need to continue the rest parts of the algorithm, and the complexity is $O(m)$. When $1 < \alpha \leq m$, the second part of the algorithm firstly constructs a complete graph from relevant features and the complexity is $O(\alpha^2)$, and then generates a MST from

The graph using Prim algorithm whose time complexity is $O(\alpha^2)$. The third part partitions the MST and chooses the representative features with the complexity of $O(\alpha)$. Thus when $1 < \alpha \leq m$, the complexity of the algorithm is $O(\alpha + \alpha^2)$. This means when $\alpha \leq \sqrt{m}$, FAST has linear complexity ($O(m)$), while obtains the worst complexity $O(m^2)$ when $\alpha = m$. However, α is heuristically set to be $\lfloor \sqrt{m * \lg m} \rfloor$ in the implementation of FAST. So the complexity is $O(m * \lg m)$, which is typically less than $O(m^2)$ since $\alpha^2 \leq m^2 < m$. This can be explained as follows. Let $\alpha(m) = \sqrt{m * \lg m}$, so the derivative $\alpha'(m) = \frac{1}{2} * (1 + \lg m) / \sqrt{m * \lg m}$, which is greater than zero when $m > 1$. So $\alpha(m)$ is an increasing function and it is greater than $\alpha(1)$ which is equal to 1, i.e., $\alpha(m) > \lg m$, when $m > 1$. This means the bigger the m is, the farther the time complexity of FAST deviates from $O(m^2)$. Thus, on high dimensional data, the time complexity of FAST is far more less than $O(m^2)$.

V. DATA SOURCE

For the purposes of evaluating the performance and effectiveness of our proposed FAST algorithm, verifying whether or not the method is potentially useful in practice, and allowing other researchers to confirm our results, 35 publicly available data sets were used. The numbers of features of the 35 data sets vary from 37 to 49152 with a mean of 7874. The dimensionality of the 54.3% data sets exceed 5000, of which 28.6% data sets have more than 10000 features. The 35 data sets cover a range of application domains such as text, image and bio microarray data classification.

VI. EXPERIMENTAL SETUP

To evaluate the performance of our proposed FAST algorithm and compare it with other feature selection algorithms in a fair and reasonable way, we set up our experimental study as follows. 1) The proposed algorithm is compared with five different types of representative feature selection algorithms. They are (i) FCBF, (ii) Relief, (iii) CFS, (iv) Consist and (v) FOCUS-SF [2], respectively. FCBF and Relief evaluate features individually. For FCBF, in the experiments, we set the relevance threshold to be the $\frac{1}{\log h}$ value of the $\frac{1}{\log h}$ ranked feature for each data set (h is the number of features in a given data set) as suggested by Yu and Liu. Relief searches for nearest neighbors of instances of different classes and weights features according to how well they differentiate instances of different classes. The other three feature selection algorithms are based on subset evaluation. CFS exploits best-first search based on the evaluation of a subset that contains features highly correlated with the target concept, yet uncorrelated with each other. The Consist method searches for the minimal subset that separates classes as consistently as the full set can under best-first search strategy. FOCUS-SF is a variation of FOCUS [2]. FOCUS has the same evaluation strategy as Consist, but it examines all subsets of features. Considering the time efficiency, FOCUS-SF replaces exhaustive search in FOCUS with sequential forward selection. For our proposed FAST algorithm, we heuristically set α to be the $\frac{1}{\log h}$ value of the $\frac{1}{\log h}$ ranked feature for each data set.

Four different types of classification algorithms are employed to classify data sets before and after feature selection. They are (i) the probability-based Naive Bayes (NB), (ii) the tree-based C4.5, (iii) the instance-based lazy learning algorithm IB1, and (iv) the rule-based RIPPER, respectively. Naive Bayes utilizes a probabilistic method for classification by multiplying the individual probabilities of every feature-value pair. This algorithm assumes independence among the features and even then provides excellent classification results. Decision tree learning algorithm C4.5 is an extension of ID3 that accounts for unavailable values, continuous attribute value ranges, pruning of decision trees, rule derivation, and so on. The tree comprises of nodes (features) that are selected by information entropy. Instance-based learner IB1 is a single-nearest-neighbor algorithm, and it classifies entities taking the class of the closest associated vectors in the training set via distance metrics. It is the simplest among the algorithms used in our study. Inductive rule learner RIPPER (Repeated Incremental Pruning to Produce Error Reduction) is a propositional rule learner that defines a rule based detection model and seeks to improve it iteratively by using different heuristic techniques. The constructed rule set is then used to classify new instances.

3) When evaluating the performance of the feature subset selection algorithms, four metrics, (i) the proportion of selected features (ii) the time to obtain the feature subset, (iii) the classification accuracy, and (iv) the Win/Draw/Loss record, are used. The proportion of selected features is the ratio of the number of features selected by a feature selection algorithm to the original number of features of a data set. The Win/Draw/Loss record presents three values on a given measure, i.e. the numbers of data sets for which our proposed algorithm FAST obtains better, equal, and worse performance than other five feature selection algorithms, respectively. The measure can be the proportion of selected features, the runtime to obtain a feature subset, and the classification accuracy, respectively.

VII. RESULTS AND ANALYSIS

In this paper present the experimental results in terms of the proportion of selected features, the time to obtain the feature subset, the classification accuracy, and the Win/Draw/Loss record. For the purpose of exploring the statistical significance of the results, we performed a nonparametric Friedman test followed by Nemenyi post-hoc test, as advised by Demsar and Garcia and Herrerato to statistically compare algorithms on multiple data sets. Thus the Friedman and the Nemenyi test results are reported as well

7.1. Proportion of selected features

Records the proportion of selected features of the six feature selection algorithms for each data set. From it we observe that) generally all the six algorithms achieve significant reduction of dimensionality by selecting only a small portion of the original features. FAST on average obtains the best proportion of selected features of 1.82%. The Win/Draw/Loss records show FAST wins other algorithms as well. 2) For image data, the proportion of selected features of each algorithm has an increment compared with the corresponding average proportion of selected features on the given data sets except Consist has an improvement. This reveals that the five algorithms are not very suitable to choose features for image data compared with for microarray and text data. FAST ranks 3 with the proportion of selected features of 3.59% that has a tiny margin of 0.11% to the first and second best proportion of selected features 3.48% of Consist and FOCUS-SF, and a margin of 76.59% to the worst proportion of selected features 79.85% of Relief. 3) For microarray data,

the proportion of selected features has been improved by each of the six algorithms compared with that on the given data sets. This indicates that the six algorithms work well with microarray data. FAST ranks 1 again with the proportion of selected features of 0.71%. Of the six algorithms, only CFS cannot choose features for two data sets whose dimensionalities are 19994 and 49152, respectively. 4) For text data, FAST ranks 1 again with a margin of 0.48% to the second best algorithm FOCUS-SF. TABLE 2: Proportion of selected features of the six feature selection algorithms. The Friedman test can be used to compare k algorithms over N data sets by ranking each algorithm on each data set separately. The algorithm obtained the best performance gets the rank of 1, the second best ranks 2, and so on. In case of ties, average ranks are assigned. Then the average ranks of all algorithms on all data sets are calculated and compared. If the null hypothesis, which is all algorithms are performing equivalently, is rejected under the Friedman test statistic, post-hoc tests such as the Nemenyi test can be used to determine which algorithms perform statistically different. The Nemenyi test compares classifiers in a pairwise manner. According to this test, the performances of two classifiers are significantly different if the distance of the average ranks exceeds the critical distance $CD = \frac{1}{\sqrt{2}} \sqrt{(k+1)6}$, where the $\frac{1}{\sqrt{2}}$ is based on the studentized range statistic [48] divided by $\sqrt{2}$. In order to further explore whether the reduction rates are significantly different we performed a Friedman test followed by a Nemenyi post-hoc test. The null hypothesis of the Friedman test is that all the feature selection algorithms are equivalent in terms of proportion of selected features. The test result is $p=0$. This means that at $\alpha = 0.1$, there is evidence to reject the null hypothesis and all the six feature selection algorithms are different in terms of proportion of selected features

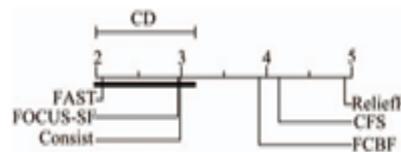


Fig. 3: Proportion of selected features

Comparison of all feature selection algorithms against each other with the Nemenyi test. In order to further explore feature selection algorithms whose reduction rates have statistically significant differences, we performed a Nemenyi test. Fig. 3 shows the results with $\alpha = 0.1$ on the 35 data sets. The results indicate that the proportion of selected features of FAST is statistically smaller than those of Relief, CFS and FCBF, and there is no consistent evidence to indicate statistical differences between FAST, Consist, and FOCUS-SF, respectively. The 10-fold cross-validation accuracies of the four different types of classifiers on the 35 data sets before and after each feature selection algorithm is performed, respectively. The classification accuracy of Naïve Bayes. From it we observe that 1) Compared with original data, the classification accuracy of Naive Bayes has been improved by FAST, CFS, and FCBF by 12.86%, 6.62%, and 4.32%, respectively. Unfortunately, Relief, Consist, and FOCUS-SF have decreased the classification accuracy by 0.32%, 1.35%, and 0.86%, respectively. FAST ranks 1 with a margin of 6.24% to the second best accuracy 80.60% of CFS. At the same time, the Win/Draw/Loss records show that FAST outperforms all other five algorithms. 2) For image data, the classification accuracy of Naïve Bayes has been improved by FCBF, CFS, FAST, and Relief by 6.13%, 5.39%, 4.29%, and 3.78%, respectively. However, Consist and FOCUS-SF have decreased the classification accuracy by 4.69% and 4.69%, respectively. This time FAST ranks 3 with a margin of 1.83% to the best accuracy 87.32% of FCBF. 3) For microarray data, the classification accuracy of Naive Bayes has been improved by all six algorithms FAST, CFS, FCBF, Relief, Consist, and FOCUS-SF by 16.24%, 12.09%, 9.16%, 4.08%, 4.45%, and 4.45%, respectively. FAST ranks 1 with a margin of 4.16% to the second best accuracy 87.22% of CFS. This indicates that FAST is more effective than others when using Naive Bayes to classify microarray data. 4) For text data, FAST and CFS have improved the classification accuracy of Naive Bayes by 13.83% and 1.33%, respectively. Other four algorithms Relief, Consist, FOCUS-SF, and FCBF have decreased the accuracy by 7.36%, 5.87%, 4.57%, and 1.96%, respectively. FAST ranks 1 with a margin of 12.50% to the second best accuracy 70.12% of CFS.

The classification accuracy of C4.5. From it we observe that 1) Compared with original data, the classification accuracy of C4.5 has been improved by FAST, FCBF, and FOCUS-SF by 4.69%, 3.43%, and 2.58%, respectively. Unfortunately, Relief, Consist, and CFS have decreased the classification accuracy by 3.49%, 3.05%, and 2.31%, respectively. FAST obtains the rank of 1 with a margin of 1.26% to the second best accuracy 81.17% of FCBF. 2) For image data, the classification accuracy of C4.5 has been improved by all the six feature selection algorithms FAST, FCBF, CFS, Relief, Consist, and FOCUS-SF by 5.31%, 4.54%, 7.20%, 0.73%, 0.60%, and 0.60%, respectively.

This time FAST ranks 2 with a margin of 1.89% to the best accuracy 83.6% of CFS and a margin of 4.71% to the worst accuracy 76.99% of Consist and FOCUS-SF. 3) For microarray data, the classification accuracy of C4.5 has been improved by all the six algorithms FAST, FCBF, CFS, Relief, Consist, and FOCUS-SF by 11.42%, 7.14%, 7.51%, 2.00%, 6.34%, and 6.34%, respectively. FAST ranks 1 with a margin of 3.92% to the second best accuracy 79.85% of CFS. 4) For text data, the classification accuracy of C4.5 has been decreased by algorithms FAST, FCBF, CFS, Relief, Consist and FOCUS-SF by 4.46%, 2.70%, 19.68%, 13.25%, 16.75%, and 1.90% respectively. FAST ranks 3 with a margin of 2.56% to the best accuracy 83.94% of FOCUS-SF and a margin of 15.22% to the worst accuracy 66.16% of CFS.

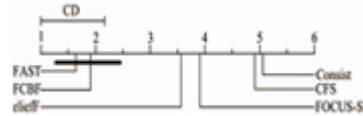


Fig.4: Runtime comparison of all feature selection algorithms against each other with the Nemenyi test.

The classification accuracy of RIPPER. From it we observe that 1) Compared with original data, the classification accuracy of RIPPER has been improved by the five feature selection algorithms FAST, FCBF, CFS, Consist, and FOCUS-SF by 7.64%, 4.51%, 4.08%, 5.48%, and 5.32%, respectively; and has been decreased by Relief by 2.04%. FAST ranks 1 with a margin of 2.16% to the second best accuracy 78.46% of Consist. The Win/Draw/Loss records show that FAST outperforms all other algorithms.

2) For image data, the classification accuracy of

RIPPER has been improved by all the six feature selection algorithms FAST, FCBF, CFS, Relief, Consist, and FOCUS-SF by 12.35%, 8.23 %, 4.67%, 3.86%, 4.90%, and 4.90%, respectively. FAST ranks 1 with a margin of 4.13% to the second best accuracy 76.52% of FCBF.

3) For microarray data, the classification accuracy of RIPPER has been improved by all the six algorithms FAST, FCBF, CFS, Relief, Consist, and FOCUS-SF by 13.35%, 6.13%, 5.54%, 3.23%, 9.02%, and 9.02%, respectively. FAST ranks 1 with a margin of 4.33% to the second best accuracy 77.33% of Consist and FOCUS-SF. 4) For text data, the classification accuracy of RIPPER has been decreased by FAST, FCBF, CFS, Relief, Consist, and FOCUS-SF by 10.35%, 8.48%, 7.02%, 20.21%, 7.25%, and 7.68%, respectively. FAST ranks 5 with a margin of 3.33% to the best accuracy 82.81% of CFS

This means that at $\alpha = 0.1$, there are evidences to reject the null hypotheses and the accuracies are different further differences exist in the six feature selection algorithms.

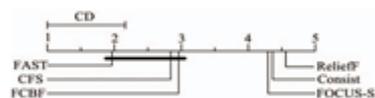


Fig. 5: Accuracy comparison of Naive Bayes with the six feature selection algorithms against each other with the Nemenyi test

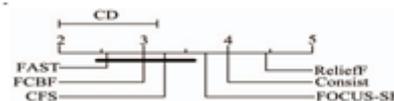


Fig. 6: Accuracy comparison of C4.5 with the six feature selection algorithms against each other with the Nemenyi test.

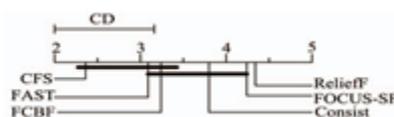


Fig. 7: Accuracy comparison of IB1 with the six feature selection algorithms against each other with the Nemenyi test

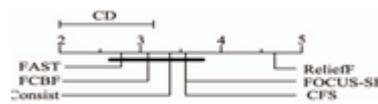


Fig. 8: Accuracy comparison of RIPPER with the six feature selection algorithms against each other with the Nemenyi test.

From Fig. 5 we observe that the accuracy of Naïve Bayes with FAST is statistically better than those with Relief, Consist, and FOCUS-SF. But there is no consistent evidence to indicate statistical accuracy differences between Naive Bayes with FAST and with CFS, which also holds for Naive Bayes with FAST and with FCBF. From Fig. 6 we observe that the accuracy of C4.5 with FAST is statistically better than those with Relief, Consist, and FOCUS-SF. But there is no consistent evidence to indicate statistical accuracy differences between C4.5 with FAST and with FCBF, which also holds for C4.5 with FAST and with CFS. From Fig. 7 we observe that the accuracy of IB1 with FAST is statistically better than those with Relief. But there is no consistent evidence to indicate statistical accuracy differences between IB1 with FAST and with FCBF, Consist, and FOCUS-SF, respectively, which also holds for IB1 with FAST and with CFS. From Fig. 8 we observe that the accuracy of RIPPER with FAST is statistically better than those with Relief. But there is no consistent evidence to indicate statistical accuracy differences between RIPPER with FAST and with FCBF, CFS, Consist, and FOCUS-SF, respectively. For the purpose of exploring the relationship between feature selection algorithms and data types, i.e. which algorithms are more suitable for which types of data, we rank the six feature selection algorithms according to the classification accuracy of a given classifier on a specific type of data after the feature selection algorithms are performed. Then we summarize the ranks of the feature selection algorithms under the four different classifiers, and give the final ranks of the feature selection algorithms on different types of data. Table 8 shows the results. From Table 8 we observe that (i) for image data, CFS obtains the rank of 1, and FAST ranks 3; (ii) for microarray data, FAST ranks 1 and should be the undisputed first choice, and CFS is a good alternative; (iii) for text data, CFS obtains the rank of 1, and FAST and FCBF are alternatives; and (iv) for all data, FAST ranks 1 and should be the undisputed first choice, and FCBF, CFS are good alternatives.

Sensitivity analysis : Like many other feature selection algorithms, our pro-posed FAST also requires a parameter α that is the threshold of feature relevance. Different α values might end with different classification results. In order to explore which parameter value results in the best classification accuracy for a specific classification problem with a given classifier, a 10 fold cross-validation strategy was employed to reveal how the classification accuracy is changing with value of the parameter.

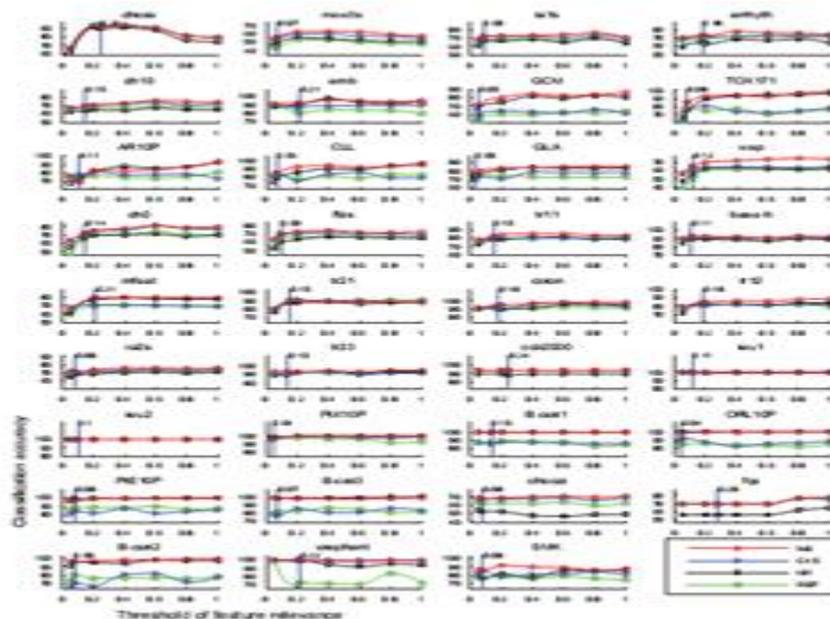


Fig. 9: Accuracies of the four classification algorithms with different α values.

Fig. 9 shows the results where the 35 subfigures rep-represent the 35 data sets, respectively. In each subfigure, the four curves denotes the classification accuracies of the four classifiers with the different θ values. The cross points of the vertical line with the horizontal axis repre-sent the default values of the parameter θ recommended by FAST, and the cross points of the vertical line with the four curves are the classification accuracies of the corresponding classifiers with the θ values. From it we observe that: Classification accuracies; (ii) there is a θ value where the corresponding classification accuracy is the best; and (iii) the θ values, in which the best classification accuracies are obtained, are different for both the different data sets and the different classification algorithms. Therefore, an appropriate θ value is desired for a specific classification problem and a given classification algorithm. 2) In most cases, the default θ values recommended by FAST are not the optimal. Especially, in a few cases (e. g., data sets GCM, CLL-SUB-11, and TOX-171), the corresponding classification accuracies are very small. This means the results presented in Section 4.4.3 are not the best, and the performance could be better. 3) For each of the four classification algorithms, al-though the θ values where the best classification accuracies are obtained are different for different data sets, The value of 0.2 is commonly accepted because the corresponding classification accuracies are among the best or nearly the best ones. When determining the value of θ , besides classification accuracy.

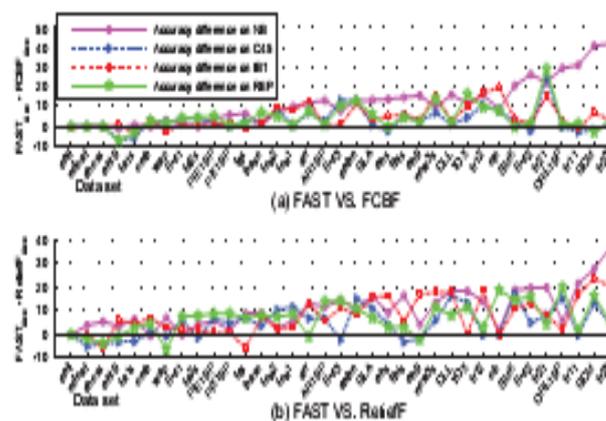


Fig. 10: Accuracy differences between FAST and the comparing algorithms

Just like the default θ values used for FAST in the experiments are often not the optimal in terms of classification accuracy, the default threshold values used for FCBF and Relief (CFS, Consist, and FOCUS-SF do not require any input parameter) could be so. In order to explore whether or not FAST still outperforms when optimal threshold values are used for the comparing algorithms, 10-fold cross-validation methods were firstly used to determine the optimal threshold values and then were employed to conduct classification for each of the four classification methods with the different feature subset selection algorithms upon the 35 data sets. The results reveal that FAST still outperforms both FCBF and Relief for all the four classification methods, Fig. 10 shows the full details. signed ranks tests with $\alpha = 0.05$ were performed to confirm the results as advised by Demsar. All the α values are smaller than 0.05, this indicates that the FAST is significantly better than both FCBF and Relief.

VIII. CONCLUSION

In this paper, we have presented a novel clustering-based feature subset selection algorithm for high dimensional data. The algorithm involves (i) removing irrelevant features, (ii) constructing a minimum spanning tree from relative ones, and (iii) partitioning the MST and selecting representative features. In the proposed algorithm, a cluster consists of features. Each cluster is treated as a single feature and thus dimensionality is drastically reduced. We have compared the performance of the proposed algorithm with those of the five well-known feature selection algorithms FCBF, Relief, CFS, Consist, and FOCUS-SF on the 35 publicly available image, microar-ray, and text data from the four different aspects of the proportion of selected features, runtime, classification accuracy of a given classifier, and the Win/Draw/Loss record. Generally, the proposed algorithm obtained the best proportion of selected features, the best runtime, and the best classification accuracy for Naive Bayes, C4.5, and RIPPER, and the second best classification accuracy for IB1. The Win/Draw/Loss records confirmed the conclusions. We also found that FAST obtains the rank of 1 for microarray data, the rank of 2 for text data, and the rank of 3 for image data in terms of classification accuracy of the four different types of classifiers, and CFS is a good alternative.

At the same time, FCBF is a good alternative for image and text data. Moreover, Consist and FOCUS-SF are alternatives for text data. For the future work, we plan to explore different types of correlation measures, and study some formal properties of feature space.

REFERENCES

- [1] Almuallim H. and Dietterich T.G., Algorithms for Identifying Relevant Features, In Proceedings of the 9th Canadian Conference on AI, pp 38-45, 1992.
- [2] Almuallim H. and Dietterich T.G., Learning boolean concepts in the presence of many irrelevant features, *Artificial Intelligence*, 69(1-2), pp 279-305, 1994.
- [3] Arauzo-Azofra A., Benitez J.M. and Castro J.L., A feature set measure based on relief, In Proceedings of the fifth international conference on Recent Advances in Soft Computing, pp 104-109, 2004.
- [4] Baker L.D. and McCallum A.K., Distributional clustering of words for text classification, In Proceedings of the 21st Annual international ACM SIGIR Conference on Research and Development in information Retrieval, pp 96-103, 1998.
- [5] Battiti R., Using mutual information for selecting features in supervised neural net learning, *IEEE Transactions on Neural Networks*, 5(4), pp 537-550, 1994.
- [6] Bell D.A. and Wang, H., A formalism for relevance and its application in feature subset selection, *Machine Learning*, 41(2), pp 175-195, 2000.
- [7] Biesiada J. and Duch W., Features election for high-dimensional data: Pearson redundancy based filter, *Advances in Soft Computing*, 45, pp 242-249, 2008.
- [8] Butterworth R., Piatetsky-Shapiro G. and Simovici D.A., On Feature Selection through Clustering, In Proceedings of the Fifth IEEE international Conference on Data Mining, pp 581-584, 2005.
- [9] Cardie, C., Using decision trees to improve case-based learning, In Proceedings of Tenth International Conference on Machine Learning, pp 25-32, 1993.
- [10] Chanda P., Cho Y., Zhang A. and Ramanathan M., Mining of Attribute Interactions Using Information Theoretic Metrics, In Proceedings of IEEE international Conference on Data Mining Workshops, pp 350-355, 2009.
- [11] Chikhi S. and Benhammada S., ReliefMSS: a variation on a feature ranking Relief algorithm. *Int. J. Bus. Intell. Data Min.* 4(3/4), pp 375-390, 2009.
- [12] Cohen W., Fast Effective Rule Induction, In Proc. 12th international Conf. Machine Learning (ICML'95), pp 115-123, 1995.
- [13] Dash M. and Liu H., Feature Selection for Classification, *Intelligent Data Analysis*, 1(3), pp 131-156, 1997.
- [14] Dash M., Liu H. and Motoda H., Consistency based feature Selection, In Proceedings of the Fourth Pacific Asia Conference on Knowledge Discovery and Data Mining, pp 98-109, 2000.
- [15] Das S., Filters, wrappers and a boosting-based hybrid for feature Selection, In Proceedings of the Eighteenth International Conference on Machine Learning, pp 74-81, 2001.
- [16] Dash M. and Liu H., Consistency-based search in feature selection. *Artificial Intelligence*, 151(1-2), pp 155-176, 2003.
- [17] Demsar J., Statistical comparison of classifiers over multiple data sets, *J. Mach. Learn. Res.*, 7, pp 1-30, 2006.
- [18] Dhillon I.S., Mallela S. and Kumar R., A divisive information theoretic feature clustering algorithm for text classification, *J. Mach. Learn. Res.*, 3, pp 1265-1287, 2003.
- [19] Dougherty, E. R., Small sample issues for microarray-based classification. *Comparative and Functional Genomics*, 2(1), pp 28-34, 2001.
- [20] Fayyad U. and Irani K., Multi-interval discretization of continuous-valued attributes for classification learning, In Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence, pp 1022-1027, 1993.