

Fast Data Collection for High Dimensional Data in Data Mining

Raju¹& Pooja Srivastvas²

¹M-Tech Computer Science and Engineering, Aurora's Scientific Technology & Research Academy, Mail Id: - raju_m593@yahoo.com

²Assistant Professor Computer Science and Engineering, Aurora's Scientific Technology & Research Academy, Mail Id: - poojasrivastavas@gmail.com

Abstract—

In machine learning, feature selection is preprocessing step and can be effectively reduce high dimensional data, remove irrelevant data, increase learning accuracy, and improve result comprehensibility. High dimensionality of data takes over efficiency and effectiveness points of view in feature selection algorithm. Efficiency stands required time to find a subset of features, and the effectiveness belongs to good quality of the subset of features. In feature selection technique high dimensional data contains many irrelevant and redundant features. Irrelevant features make available no useful information in any context, and redundant features provide no more information than the selected features. Good feature subsets contain features highly predictive of (correlated with) the class, yet not predictive of (uncorrelated with) each other. A subset of useful features to produce compatible results as the original set of features is identified from feature selection.

Keywords-Feature subset selection; graph-theoretic clustering; filter method.

INTRODUCTION

In machine learning, feature selection, also known as variable subset selection, is the process of selecting a subset of relevant features for use in model construction. Feature selection techniques have benefits when constructing correlated models: improved model to interpret the hidden meaning, shorter (small) training times, and enhanced generalization by reducing over fitting. Feature selection is helpful as part of the data analysis process, as it identifies important features for prediction. Choosing a subset of good features according to target concepts, feature subset selection has been effective to reduce dimensionality, removing irrelevant data, increasing learning accuracy, and improving comprehensibility. Feature subset selection algorithms for machine learning applications can be divided into four main categories: Wrapper, Filter, Hybrid, and Embedded methods.

Wrapper methods use a predetermined learning model to score feature subsets. A wrapper methods train a fresh model for new subset, they have high accuracy but are expensive to compute

and also limited in generality of selected features. Filter methods are faster than wrapper methods but produces a features set which is independent from learning algorithms with better generality. Filter methods measures include the correlation coefficient, Mutual Information, distance and consistency measurements to sort a good subset. Filtering approach to feature selection involves a greater degree of search through the feature space but the accuracy of the algorithms is not guaranteed.

Embedded algorithms integrate feature subset selection as a training process and they are fixed to learning methods, hence more efficient than Wrapper and Filter methods.

Decision tree algorithms are best example of embedded methods. A combination of filter methods and wrapper methods form hybrid methods which achieves best possible performance with a specific learning algorithm with similar time complexity like the filter methods. The wrapper methods tend to over fit on small training sets. The main benefits of filter

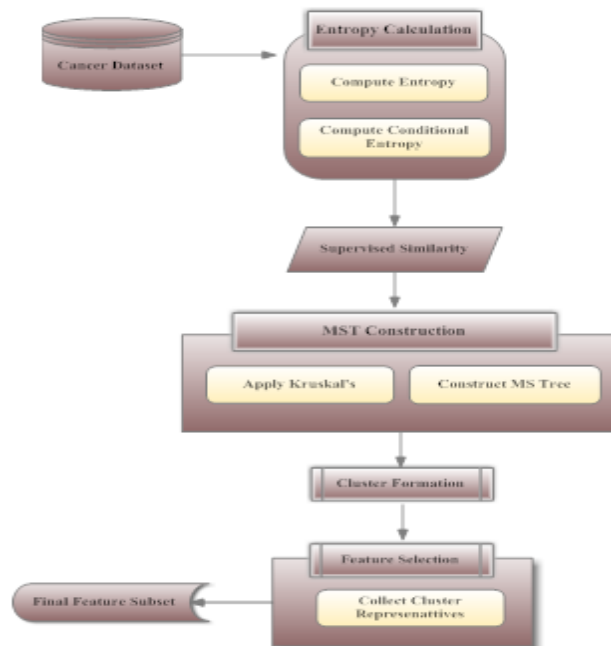
methods are they are faster and they have ability to scale to large datasets. With respect to the filter feature selection methods, the application of cluster analysis clearly gives practical demonstration and explanation to be more effective than traditional feature selection algorithms. The distributional clustering of words is agglomerative in nature and reduces the high dimensionality of text data since each word cluster can be treated as single feature but are expensive compute.

In cluster analysis, most of the applications use graph theoretic methods because they produce good results. The graph-theoretic clustering is simple since it computes a neighborhood graph of instances, and then deletes any edge in graph that is much short or long than its neighbors. The graph theoretic clustering results in forest and trees in forest represents a cluster. In this survey graph-theoretic clustering algorithms are used to **CLUSTERING**

Clustering and segmentation are the processes of creating a partition so that all the members of each set of the partition are similar according to some metric. A cluster is a set of objects grouped together because of their similarity or proximity. Objects are often decomposed into an exhaustive and/or mutually exclusive set of clusters. Clustering according to similarity is a very powerful technique, the key to it being to translate some intuitive measure of similarity into a quantitative measure. When learning is unsupervised then the system has to discover its

It is widely recognized that a large number of features can adversely affect the performance of inductive learning algorithms, and clustering is not an exception. However, while there exists a large body of literature devoted to this problem for supervised learning task, feature selection for Supervision available. In machine learning and statistics, feature selection, also known as variable selection, attribute selection or variable subset selection, is the process of selecting a subset of relevant features for use in model construction.

features; particularly minimum spanning tree based clustering algorithms.



own classes i.e. the system clusters the data in the database. The system has to discover subsets of related objects in the training set and then it has to find descriptions that describe each of these subsets. There are a number of approaches for forming clusters. One approach is to form rules which dictate membership in the same group based on the level of similarity between members. Another approach is to build set functions that measure some property of partitions as functions of some parameter of the partition.

FEATURE SELECTION

clustering has been rarely addressed. The problem appears to be a difficult one given that it inherits all the uncertainties that surround this type of inductive learning. Particularly, that there is not a single performance measure widely accepted for this task and the lack of supervision available.

The central assumption when using a feature selection technique is that the data contains many redundant or irrelevant features. Redundant features are those which provide no more information than the currently selected features,

and irrelevant features provide no useful information in any context. Feature selection techniques are a subset of the more general field of feature extraction. Feature extraction creates new features from functions of the original features, whereas feature selection returns a subset of the features. Feature selection techniques are often used in domains where there are many features and comparatively few samples (or data points). The archetypal case is the use of feature selection in analyzing DNA microarrays, where there are many thousands of features, and a few tens to hundreds of samples. Feature selection techniques provide three main benefits when constructing predictive models Improved model interpretability,

Shorter training times, Enhanced generalization by reducing over fitting.

Feature selection is also useful as part of the data analysis process, as shows which features are important for prediction, and how these features Wrapper methods are widely recognized as a superior alternative in supervised learning problems, since by employing the inductive algorithm to evaluate alternatives they have into account the particular biases of the algorithm. However, even for algorithms that exhibit a moderate complexity, the number of executions that the search process requires results in a high computational cost, especially as we shift to more exhaustive search strategies. The wrapper

Hybrid Approach

The hybrid methods are a combination of filter and wrapper methods by using a filter method to reduce search space that will be considered by the subsequent wrapper. They mainly focus on combining filter and wrapper methods to achieve the best possible performance with a particular learning algorithm with similar time complexity of the filter methods.

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 instances, then delete any

are related. With such an 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. 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.” 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

Wrapper Filter

methods use the predictive accuracy of a predetermined learning algorithm to determine the goodness of the selected subsets, 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. Their computational complexity is low, but the accuracy of the learning algorithms is not guaranteed

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 selection algorithm FAST was tested various numerical. Investigates the application of the mutual information as a criterion to evaluate a set of candidate features and to select an informative subset to be used as input data for a neural network classifier. Because the mutual information measures arbitrary dependencies between random variables, it is suitable for assessing the “information content” of features in complex classification tasks, where methods based on linear relations (like the correlation) are prone to mistakes.

The fact that the mutual information is independent of the coordinates chosen permits a robust estimation. Nonetheless, the use of the mutual information for tasks characterized by high input dimensionality requires suitable approximations because of the prohibitive demands on computation and samples. An algorithm is proposed that is based on a “greedy” selection of the features and that takes both the mutual information with respect to the output class and with respect to the already-selected features into account. Finally the results of a series of experiments are discussed.

During “preprocessing” stage, where an appropriate number of relevant features are extracted from the raw data, has a crucial impact both on the complexity of the learning phase and on the achievable generalization performance. While it is essential that the information contained in the input vector is sufficient to determine the output class, the presence of too many input features can burden the training process and can produce a neural network with more connection weights than those required by the problem.

A major weakness of these methods is that they are not invariant under a transformation of the

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 classification accuracy.

Using Mutual Information for Selecting Features in Supervised Neural Net Learning

variables. For example a linear scaling of the input variables (that may be caused by a change of units for the measurements) is sufficient to modify the PCA results. Feature selection methods that are sufficient for simple distributions of the patterns belonging to different classes can fail in classification tasks with complex decision boundaries. In addition, methods based on a linear dependence (like the correlation) cannot take care of arbitrary relations between the pattern coordinates and the different classes. On the contrary, the mutual information can measure arbitrary relations between variables and it does not depend on transformations acting on the different variables.

Our objective was less ambitious, because only the first of the above options was considered (leaving the second for the capabilities of the neural net to build complex features from simple ones). We assumed that a set of candidate features with globally sufficient information is available and that the problem is that of extracting from this set a suitable subset that is sufficient for the task, thereby reducing the processing times in the operational phase and, possibly, the training times and the cardinality of the example set needed for a good generalization.

In particular we were interested in the applicability of the mutual information measure. For this reason we considered the estimation of the MI from a finite set of samples, showing that the MI for different features is over-estimated in approximately the same way. This estimation is the building block of the MIFS algorithm, where the features are selected in a “greedy” manner, ranking them according to their MI with respect to

the class discounted by a term that takes the mutual dependencies into account. 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 shows that good classifiers can be built by using a small number of attributes located at the centers of the clusters identified in the dendrogram. This type of data compression can be achieved with little or no penalty in terms of the accuracy of the classifier produced and highlights the relative importance of attributes.

Clustering's were extracted from the tree produced by the algorithm by cutting the tree at various heights starting with the maximum height of the tree created above (corresponding to a single cluster) and working down to a height of 0 (which consists of single-attribute clusters). A „representative“ attribute was created for each cluster as the attribute that has the minimum total

IRRELEVANT FEATURES REMOVAL

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

Load Data and Classify

Load the data into the process. The data has to be preprocessed for removing missing values, noise and outliers. Then the given dataset must be converted into the arff format which is the standard format for WEKA toolkit. From the arff

distance to the other members of the cluster, again using the Barth 'elemyMontjardet distance. A similar study was undertaken for the zoo database, after eliminating the attribute animal which determines uniquely the type of the animal. These results suggest that this method has comparable accuracy to the wrapper method and CSF. However, the tree of attributes helps to understand the relationships between attributes and their relative importance.

Attribute clustering help to build classifiers in a semi-supervised manner allowing analysts a certain degree of choice in the selection of the features that may be considered by classifiers, and illuminating relationships between attributes and their relative importance for classification. With the increased interest of data miners in bio-computing in general, and in microarray data in particular, 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

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 1) the construction of the minimum spanning tree from a weighted complete graph; 2) the partitioning of the MST into a forest with each tree representing a cluster; and 3) the selection of representative features from the clusters.

format, only the attributes and the values are extracted and stored into the database. By considering the last column of the dataset as the class attribute and select the distinct class labels

from that and classify the entire dataset with **Information Gain Computation**

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. To find the relevance of each attribute with the class label, Information gain is computed in this module. This is also said to be Mutual Information measure. 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 The relevance between the feature $F_i \in F$ and the target concept C is referred to as the T-Relevance of F_i and C , and denoted by $SU(F_i, C)$. If $SU(F_i, C)$ is greater than a predetermined threshold, we say that F_i is a strong T-Relevance feature.

The correlation between any pair of features F_i and F_j ($F_i, F_j \in F \wedge i \neq j$) is called the F-Correlation of F_i and F_j , and denoted by $SU(F_i, F_j)$. The equation symmetric uncertainty which is

MST Construction

With the F-Correlation value computed above, the Minimum Spanning tree is constructed. For that, we use Kruskal's algorithm which form MST effectively. Kruskal's algorithm is a greedy algorithm in graph theory that finds a minimum spanning tree for a connected weighted graph.

Description:

1. Create a forest F (a set of trees), where each vertex in the graph is a separate tree.
 2. Create a set S containing all the edges in the graph
 3. While S is nonempty and F is not yet spanning
- At the termination of the algorithm, the forest forms a minimum spanning forest of the graph. If the graph is connected, the forest has a single component and forms a minimum spanning tree. The sample tree is as follows,

respect to class labels.

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

The symmetric uncertainty is defined as follows:
 $Gain(X|Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$

To calculate gain, we need to find the entropy and conditional entropy values. The equations for that are given below:

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$

$$H(X|Y) = - \sum_{y \in Y} p(y) \sum_{x \in X} p(x|y) \log_2 p(x|y)$$

Where $p(x)$ is the probability density function and $p(x|y)$ is the conditional probability density function.

T-Relevance Calculation

$$SU(X, Y) = 2 \times Gain(X|Y) / (H(X) + H(Y))$$

After finding the relevance value, the redundant attributes will be removed with respect to the threshold value

F-Correlation Calculation

used for finding the relevance between the attribute and the class is again applied to find the similarity between two attributes with respect to each label.

This means it finds a subset of the edges that forms a tree that includes every vertex, where the total weight of all the edges in the tree is minimized. If the graph is not connected, then it finds a minimum spanning forest (a minimum spanning tree for each connected component).

- Remove an edge with minimum weight from S
 If that edge connects two different trees, then add it to the forest, combining two trees into a single tree
 Otherwise discard that edge.

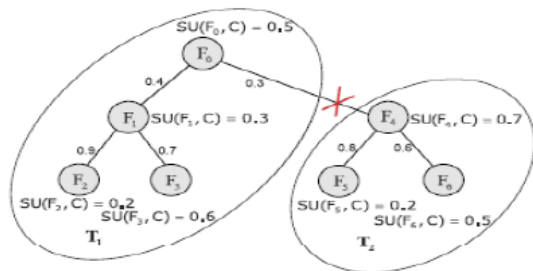


Fig 2. Correlations

ALGORITHM:-

Inputs: $D (F_1, F_2 \dots F_m, C)$ - the given data set θ - 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 =====

In this tree, the vertices represent the relevance value and the edges represent the F-Correlation value. The complete graph G reflects the correlations among all the target-relevant features. Unfortunately, graph G has k vertices and $k(k-1)/2$ edges. For high-dimensional data, it is heavily dense and the edges with different weights

Cluster Formation

After building the MST, in the third step, we first remove the edges whose weights are smaller than both of the T-Relevance $SU(F_i, C)$ and $SU(F_j, C)$, from the MST. After removing all the unnecessary edges, a forest is obtained. Each tree

In this Project present a FAST clustering-based feature subset selection algorithm for high dimensional data. The algorithm involves 1) removing irrelevant features, 2) constructing a minimum spanning tree from relative ones, and 3) 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

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** **weight** **of** **the** **corresponding** **edge**;

9 minSpanTree = $\text{KRUSKALS} (G);$ //Using KRUSKALS Algorithm to generate the minimum spanning tree

//==== Part 3: Tree Partition and Representative Feature Selection =====

10 Forest = minSpanTree

11 for each edge \in 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 Forest = $\text{Forest} - E_{ij}$

14 $S = \emptyset$

15 for each tree \in Forest **do**

16 $F_{jr} = \text{argmax}_{F \in S_{U_i}} S_{U_i}$

17 $S = S \cup \{F_{jr}\};$

18 returns S

are strongly interwoven. Moreover, the decomposition of complete graph is NP-hard. Thus for graph G , we build an MST, which connects all vertices such that the sum of the weights of the edges is the minimum, using the well known Kruskal algorithm. The weight of edge (F^i, F^j) is F-Correlation $SU(F^i, F^j)$.

$T_j \in \text{Forest}$ represents a cluster that is denoted as $V (T_j)$, which is the vertex set of T_j as well. As illustrated above, the features in each cluster are redundant, so for each cluster $V (T_j)$ we choose a representative feature $F_j R$ who's T-Relevance $SU(F_j R, C)$ is the greatest.

CONCLUSION

drastically reduced. The text data from the four different aspects of the proportion of selected features, run time, classification accuracy of a given classifier. Clustering-based feature subset selection algorithm for high dimensional data. For the future work, we plan to explore different types of correlation measures, and study some formal properties of feature space. In feature we are going to classify the high dimensional data.

FUTURE WORK

Minimum-spanning tree with the help of visual features to identify the main list in a page.

The key contributions of this concept are:

We defined a novel top-k list extract problem which is useful in knowledge discovery and fact answering;

We designed an unsupervised general-purpose algorithm along with a number of key optimizations that is capable of extracting top-k lists from any web pages

Our evaluation shows that our algorithm scales with the data size and achieves significantly better accuracy than competing methods.

Our basic algorithm runs in four steps. First, we compute the tag path for every node in the

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] 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.

[3] Biesiada J. and Duch W., Features election for high-dimensional data: a Pearson redundancy based filter, Advances in Soft Computing, 45, pp 242-249, 2008.

[4] 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.

[5] 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.

[6] Dash M. and Liu H., Consistency-based search in feature selection. Artificial Intelligence, 151(1-2), pp 155-176, 2003.

[7] Demsar J., Statistical comparison of classifiers over multiple data sets, J. Mach. Learn. Res., 7, pp 1-30, 2006.

Minimum Spanning tree of the input page. Second, we group nodes with an identical tag path into one equivalence class, and we select those equivalence classes which have exactly k members as our candidate classes. In effect, an equivalence class represents a list of item components. Third, for each of these candidate classes, we employ an GrowUp operation to merge some of the equivalence classes together, which essentially form a number of candidate lists. Now the item components that belong to the same list item are grouped together. Finally, we rank the candidate list by their importance to the page, and return the top ranking list as the result.

[8] Fleuret F., Fast binary feature selection with conditional mutual Information, Journal of Machine Learning Research, 5, pp 1531-1555, 2004.

[9] Forman G., An extensive empirical study of feature selection metrics for text classification, Journal of Machine Learning Research, 3, pp 1289-1305, 2003.

[10] Garcia S and Herrera F., An extension on "Statistical Comparisons of Classifiers over Multiple Data Sets" for all pairwise comparisons, J. Mach. Learn. Res., 9, pp 2677-2694, 2008.

[11] Guyon I. and Elisseeff A., An introduction to variable and feature selection, Journal of Machine Learning Research, 3, pp 1157-1182, 2003.

[12] Hall M.A., Correlation-Based Feature Selection for Discrete and Numeric Class Machine Learning, In Proceedings of 17th International Conference on Machine Learning, pp 359-366, 2000.