



INTERNATIONAL JOURNAL OF ENGINEERING SCIENCES & RESEARCH TECHNOLOGY

Subset Selection in High Dimensional Data by Using Fast Technique

Ms.R.Devipriya (M.E)^{*1}, Ms.J.Jacquin Margret,M.E(AP)²

^{*1,2} Department of CSE, Muthayammal Engineering College, India

devipriyaatrs@gmail.com

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 using MST construction. 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. 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 algorithm not only reduces the number of features, but also improves the performances of the four well-known different types of classifiers such as the probability-based Naive Bayes, the tree-based C4.5, the instance-based IB1, and the rule-based RIPPER before and after feature selection. We can construct FAST algorithm with prim's algorithm based on MST construction. Our experimental results show that improves the performances of the four types of classifiers.

Keywords: Feature Subset Selection, Feature Clustering, Graph Based Clustering

Introduction

Data mining (the analysis step of the "Knowledge Discovery in Databases" process, or KDD), an interdisciplinary subfield of computer science, is the computational process of discovering patterns in large data sets involving methods at the intersection of artificial intelligence, learning, statistics, and database systems. The overall goal of the data mining process is to extract information from a data set and transform it into an understandable structure for further use. A side from the raw analysis step, it involves database and data management in aspects, model and inference considerations, metrics, complexity considerations, post-processing.

Besides the term clustering, there are a number of terms with similar meanings, including automatic classification, numerical taxonomy, botryology and typological analysis. The subtle differences are often in the usage of the results: while in data mining, the resulting groups are the matter of interest, in automatic classification primarily their discriminative power is of interest. This often leads to misunderstandings between researchers coming from the fields of data mining and machine learning, since they use the same terms and often the same algorithms, but have different goals

Data Clustering

Due to the increase in data size, human manual labeling has become extremely difficult and expensive. Therefore, automatic labeling has become indispensable step in data mining. Data clustering is one of the most popular data labeling techniques.

In data clustering, we are given unlabeled data and we are to put similar samples in one pile, called a cluster, and the dissimilar samples should be in different clusters. Usually, neither cluster's description nor its quantification is given in advance unless a domain knowledge exists, which poses a great challenge in data clustering.

Partitioning Methods

The partitioning methods generally result in a set of M clusters, each object belonging to one cluster. Each cluster may be represented by a centroid or a cluster representative; this is some sort of summary description of all the objects contained in a cluster. The precise form of this description will depend on the type of the object which is being clustered. In case where real-valued data is available, the arithmetic mean of the attribute vectors for all objects within a cluster provides an appropriate representative; alternative types of centroid may be required in other cases, e.g., a cluster of documents can be represented by a list of those keywords that

occur in some minimum number of documents within a cluster. If the number of the clusters is large, the centroids can be further clustered to produce hierarchy within a dataset.

Preprocessing

Data pre-processing is an important step in the data mining process. The phrase "garbage in, garbage out" is particularly applicable to data mining and machine learning projects. Data-gathering methods are often loosely controlled, resulting in out-of-range values, impossible data combinations missing values, etc. Analyzing data that has not been carefully screened for such problems can produce misleading results. Thus, the representation and quality of data is first and foremost before running an analysis. If there is much irrelevant and redundant information present or noisy and unreliable data, then knowledge discovery during the training phase is more difficult. Data preparation and filtering steps can take considerable amount of processing time. Data preprocessing that includes cleaning, normalization, transformation, feature extraction and selection, etc. The product of data pre-processing is the final training set.

MST Construction

A minimum spanning tree (MST) or minimum weight spanning tree is then a spanning tree with weight less than or equal to the weight of every other spanning tree. More generally, any undirected graph (not necessarily connected) has a minimum spanning forest, which is a union of minimum spanning trees for its connected components. Finding the smallest edge can be done at the same time as updating minimum spanning tree.

When building a minimum spanning tree on a complete graph, an algorithm which has a complexity based on the number of edges must have a complexity better than $O(M)$ to beat Prim's algorithm.

Tree Partition

Each tree in the MST represents a cluster. In this module, 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.

Feature Selection

Feature subset selection was split up into two parts, subset searching and criterion functions. For both parts, the common algorithms were introduced and analyzed. 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. 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.

Applications

Data clustering has immense number of applications in every field of life. One has to cluster a lot of thing on the basis of similarity either consciously or unconsciously. So the history of data clustering is old as the history of mankind.

In computer field also, use of data clustering has its own value. Specially in the field of information retrieval data clustering plays an important role. Some of the applications are listed below.

Proposed Work

The feature subset selection can be viewed as the process of identifying and removing as many irrelevant and redundant features is proposed. This is because irrelevant features do not contribute to the predictive accuracy and redundant features do not add to getting a better predictor for that they provide mostly information which is already present in other feature(s) and construct minimum spanning trees to evaluate whether two sets of n-dimensional data are from the same distribution. Irrelevant features. The irrelevant feature removal is straightforward once the right relevance measure is defined or selected, while the redundant feature elimination A minimum spanning tree is built across the data points, and edges which connect data from one distribution to the other are removed. If many edges are removed, then the data from the distributions are mixed up together, and so they must come from the same distribution. A minimum-spanning tree is a sub-graph of a weighted, connected and undirected graph. It is acyclic, connects all the nodes in the graph, and the sum of all of the weight of all of its edges is minimum. That is, there is no other spanning tree, or sub-graph which connects all the nodes and has a smaller sum.

Conclusion

Feature subset selection is an effective way for reducing dimensionality, removing irrelevant

data, increasing learning accuracy. In this cluster based feature subset selection algorithm is used to select the features in efficient and accuracy. The cluster based feature subset selection algorithm involves (i) removing irrelevant features, (ii) constructing a minimum spanning tree from relative ones, (iii) partitioning the MST and selecting representative features. In the proposed system FAST algorithm is used. A cluster consists of features. Each

cluster is treated as a single feature and thus dimensionality is drastically reduced. The FAST algorithm obtained the best proportion of selected features, the best runtime, and the best classification accuracy for Naïve Bayes, C4.5, and RIPPER, and the second best classification accuracy for IB1. The clustering-based strategy of FAST has a high probability of producing a subset of useful and independent features.

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: a 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.