

Improving Efficiency of Feature Selection by Using Global Redundancy Minimization

Misss.Shweta Satish Shringarputale

Department of Computer Science and Engineering
MarthawadaShikshanPrasarkMandal"sDeogiri
Institute of Engineering & Management Studies,
Aurangabad Maharashtra state, India.

Mr. P. R. Rathod

Associate Professor Department of Computer Science and
Engineering MarthawadaShikshanPrasarkMandal"sDeogiri
Institute of Engineering & Management Studies, Aurangabad
Maharashtra state, India.

Abstract—The amount of information over internet has been growing last few years. And it has caused risk of information problem of accessing related data to the users. The information demand of the online users can be figured out by evaluating user's web navigation behavior. Web Usage Mining (WUM) is used to extract knowledge from Web users access logs by using Data Mining Techniques. One of the applications of WUM is Web Sites Recommendation system which is personalized information filtering technique used to either determine whether a certain user will approve a given item or to identify a list of items which can be of significant importance to the user. In this paper the modified architecture that integrates item information with user's access log data and then find pattern and make pattern clustering. There after generates a set of recommendations for the user. So execution time and fetching time is reduced.

Other experiments compared the CFS to a wrapper - a well-known approach to feature selection that uses the target learning algorithm to evaluate sets of features. In many cases CFS has given results comparable to the envelope, and in general, surpassed the envelope on small sets of data. CFS runs much faster than the wrapper, enabling it to extend to larger sets of data.

Keywords-*Feature selection, Feature ranking, Redundancy minimization, Radial Basis Function Kernel, Direct kernel*

I. INTRODUCTION

Web mining is the utilization of information mining systems to concentrate learning from web information, including web archives, hyperlinks between reports, use logs of sites, and so forth. Web mining can be comprehensively separated into three unmistakable classifications, as indicated by the sorts of information to be mined. Web Mining has three types that we mention below. Web Content Mining Web content mining is the way toward separating valuable data from the substance of web archives. Content information is the accumulation of certainties a website page is intended to contain. It might comprise of content, pictures, sound, video, or organized records, for example, records and tables. Use of content mining to web content has been the most generally examined. Issues tended to in content mining incorporate point revelation and following, removing affiliation designs, grouping of web archives and arrangement of website pages. Web Structure Mining The structure of a run of the mill Web chart comprises of Web pages as hubs, and hyperlinks as edges interfacing related pages. Web Structure Mining is the way toward finding structure data from the Web. This can be further isolated into two sorts in light of the sort of structure data utilized.

II. LITERATURE SURVEY

In general, there are three models of characteristic selection methods in the literature: (1) filtering methods [14] where selection is independent of classifiers, (2) wrapping methods [12] where the method of Prediction is used as a black box to score subsets of features, and (3) integrated methods where the feature selection procedure is integrated directly into the training process. In bioinformatics applications, many methods for selecting the characteristics of these categories have been proposed and applied.

Methods for selecting widely used filter characteristics include statistical F [4], relief [11, 13], mRMR, t-test and information gain [12] which calculate sensitivity Correlation, or relevance) of a characteristic with respect to (wrt) the class label distribution of the data. These methods can be characterized by the use of global statistical information. Wrapper type selection methods are tightly coupled to a specific classifier, such as the correlation-based feature selection (CFS) [9], the support vector. Recursive elimination machine (SVM-RFE) [8]. They often perform well, but their computational cost is very expensive. Recently, the regularity of sparsity in the reduction of dimensionality has been widely studied and also applied in characteristic selection studies. 1-SVM was proposed to perform characteristic selection using 1-normal

regularization which tends to give a scattered solution [3]. Because the number of selected functionalities using SVM-1 is greater than the sample size, a Huberized Hybrid MVS (HHSVM) was proposed combining both Standard 1 and Standard 2 to form a more structured regularization. But it was designed only for binary classification. In multi-task learning, in parallel work, Obozinsky Et al and Argyriou et al. Al. [1] developed a similar model for the regularization of the 2.1 standard to couple the selection of characteristics between tasks. Such regularization has close ties with the group lasso [28]. In this article we propose a new efficient and robust method of characteristic selection to use the joint minimization of the norm 2.1 on the loss function and the regularization. Instead of using a loss function based on standard 2 that is sensitive to outliers, a loss function based on the 2.1 standard is adopted in our work to suppress outliers. Motivated by previous research [1, 18], a '2.1 normal' regularization is performed to select characteristics across all data points with common sparsity, ie each characteristic (expression Gene or mass-Scores value for all data points or has large scores on all data points To solve this new objective of robust characteristic selection we propose an efficient algorithm to solve this problem of minimization of the norm 2.1 We also provide algorithmic analysis and prove the convergence of our algorithm. We have extensive experiments on six sets of bioinformatics data and our method outperforms five other commonly used methods of character selection in statistical and bioinformatics learning.

III. MATHEMATICAL EQUATION

The data matrix has been preprocessed and discredited with respect to the mean of each gene's expression (column). The number of output features (genes) say n is provided from outside by the user. The data matrix with classes $c = \{1, 2, \dots, C\}$ are the inputs. At the beginning, the first objective (obj1) i.e., the relevance of each gene is calculated by mutual information as per Equation 6. From the relevance score, the highest scorer gene id is extracted and added.

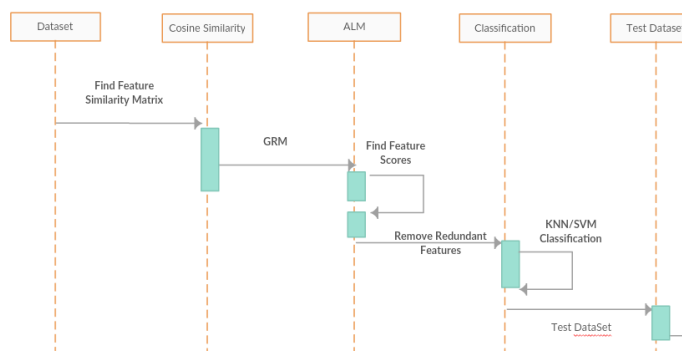


Figure 1.0 Sequence for proposed architecture

Algorithm 1 Proposed Feature Selection

Input: The feature id $idle\ ft$, first objective $ob\ j1$, second objective $ob\ j2$, $|ob\ j1| = |ob\ j2| = |idle\ ft|$.
Output: Non-dominated feature id $idns$, the second objective $ob\ j2ns$ of non-dominated features.

```

1: k = 1;
2: for i = 1 :|idle ft| do
3: t = 0;
4: for j = 1 :|idle ft| do
5: if then(i! = j)
6: if then(ob j1(i) ≤ ob j1(j)&ob j2(i) ≤ ob j2(j));
7: else if then(ob j1(i) < ob j1(j)&ob j2(i) > ob j2(j)||ob j1(i) > ob j1(j)&ob j2(i) < ob j2(j));
8: else
9: t = 1;
10: break;
11: end if
12: end if
13: end for
14: if then(t == 0&j == |idle ft|)
15: idns(k) = i;
16: obj2ns(k) = ob j2(i);
17: k = k + 1;
18: end if
19: end for
    
```

in the final solution set. Next a looping is performed for the remaining output features. Now the redundancy between the output feature and the remaining features ($idle\ ft$) is calculated as per Equation 5. If the output feature set contains more than one feature then the mean is considered as the redundancy score as in Equation 1

$$\text{mean-redundancy}(i) = \frac{\sum_{k=1}^F (\text{mutual-info}[x_k, x_i])}{|F|},$$

..... EQ. 1

where F is output feature set, X_k is output feature vector and x_i is the i th feature vector. Then the second objective (obj2) is modeled as the ratio of relevance to the redundancy and it is to be maximized. After calculating the two objectives for each feature the non-dominated features are identified. A reference feature is called the non-dominated feature if it satisfies the following conditions: 1) if the obj1 of the reference feature is greater than or equal to all the other features' obj1 and the obj2 of the reference feature is greater than or equal to all the other features' obj2 2) if the obj1 of the reference feature is greater than all the other features' obj1 and the obj2 of the reference feature is less than all the other features' obj2 and vice-versa. Afterwards, from the non-dominated features, the feature having maximum obj2 is included in the output feature set.

IV. RESULT

One real life data sets is used for the comparative study. The Prostate Cancer dataset is collected from the website:

www.biomedpubs.org/supp/bi-cancer/projections/info/. The dataset contain two classes of samples.

1. Prostate: Gene expression measurements for samples of prostate tumors and adjacent prostate tissue not containing tumor were used to build this classification model. It contains 50 normal tissue and 52 prostate tumor sample. The expression matrix consists of 12533 numbers of genes and 102 numbers of samples.

Data Set	Method	Sensitivity	Specificity	Accuracy	Fscore	AUC
Prostate Cancer	Proposed Method	0.96	0.92	0.94	0.92	0.98
	Existing Method	0.92	0.86	0.89	0.90	0.94

Table 1.0 Results with existing methods

The actual data sets described above are first standardized with the Min-Max normalization technique. Then, with respect to the mean of each characteristic (gene) or column, the data are discretized. In this article, the number of output functions is taken as 100 for all algorithms. Using 10-fold cross-validation, sensitivity, specificity, precision and fscore score are calculated. Then, the mean correlation for evaluating the redundancy of the selected characteristics is also calculated. A smaller correlation value indicates that the selected functions are less redundant. In addition, the area under ROC curve (AUC) is also reported.

The metric performance values of the proposed method, existing method on the different real datasets are shown in Table 1.0 It is evident from the table that for the data series on cancer Prostate sensitivity, specificity, and AUC are respectively 0.96, 0.92, 0.94, 0.92 and 0.98, which are better than the existing method patterns in all cases.

Summary Results

Iterations	3595
Total basis functions	57
Number Correct	26
Number Incorrect	8
Percentage Correct	76

Table (a)

Iterations	7498
Total basis functions	701
Number Correct	26
Number Incorrect	8
Percentage Correct	91.17

Table (b)

	Normal	Tumor
Normal	9	0
Tumor	8	17

	Normal	Tumor
Normal	9	0
Tumor	3	22

Confusion Matrix (a) Confusion Matrix (b)

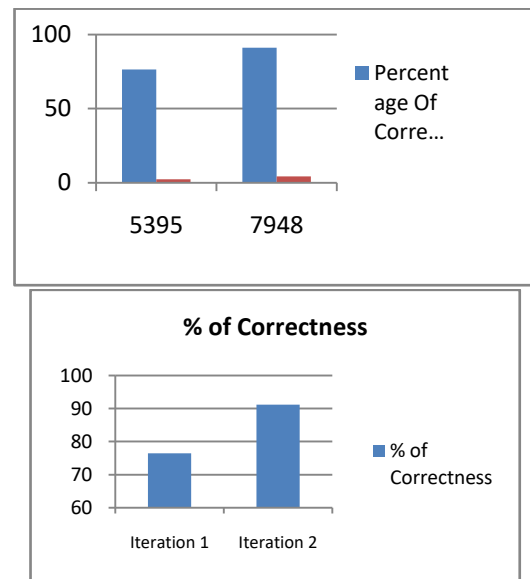


Figure 1.3 Graphs for Percentage of correctness based upon confusion matrix

CONCLUSION

There are different types of feature selection methods available in the existing literature. But in most cases, we have seen that the fundamental objective of the method is either relevance or redundancy. In this paper, we have proposed a method where relevance and redundancy are supported in parallel. To measure the relevance and redundancy of a characteristic or a gene, mutual information was considered. Relevance is defined as the mutual information between a feature vector and class labels. Redundancy is described as mutual information among the characteristics. The number of resulting functions is provided by the user. The performance of the proposed technique is evaluated on the basis of some sets of real life microarray gene expression data to select non-redundant and relevant genes.

REFERENCES

- [1]. Pena, J.M., Lozano, J.A., Larranaga, P., Inza, I. Dimensionality reduction in unsupervised learning of conditional gaussian networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 2001;23(6):590–603.
- [2]. Kurun, O., Akar, C.O., Favorov, O., Aydin, N., Urgen, F. Using covariates for improving the minimum redundancy maximum relevance feature selection method. *Turkish Journal of Electrical Engineering and Computer Sciences* 2010;18(6):975–987.
- [3]. Kamandar, M., Ghassemian, H. Maximum relevance, minimum redundancy band selection for hyperspectral images. In: *19th Iranian Conference on Electrical Engineering (ICEE)*, 2011.

- [4]. Dy, J.G., Brodley, C.E., Kak, A., Broderick, L.S., Aisen, A.M.. Unsupervised feature selection applied to content-based retrieval of lung images. *IEEE Transaction on Pattern Analysis and Machine Intelligence* 2003;25(3):373–378.
- [5]. Zhang, Z., R.Hancock, E..A graph-based approach to feature selection. In: *International Workshop on Graph-Based Representations in Pattern Recognition*. 2011,.
- [6]. Cai, D., Zhang, C., He, X.. Unsupervised feature selection for multi-cluster data. In: *16th ACM SIGKDD International Conference on Knowledge Discovery and Data mining*. 2010,.
- [7]. Ruiza, R., Riquelmea, J.C., Aguilar-Ruizb, J.S..Incremental wrapper-based gene selection from microarray data for cancer classification. *Pattern Recognition* 2006;39(12):2383–2392.
- [8]. Mitra, P., Murthy, C., Pal, S.K.. Unsupervised feature selection using feature similarity. *IEEE Transaction on Pattern Analysis and Machine Intelligence* 2002;24(3):301–312.
- [9]. Sondberg-Madsen, N., Thomsen, C., Pena, J.M.. Unsupervised feature subset selection. In: *In Proc. of the Workshop on Probabilistic Graphical Models for Classification*. 2003,.
- [10].Ding, C.H.Q.. Unsupervised feature selection via two-way ordering in gene expression analysis. *Bioinformatics* 2003;19(10):1259–1266.
- [11].Kohavi, R., John., G.. Wrapper for feature subset selection. *Artificial Intelligence* 1997;97:273–324.
- [12].Jiang, S., Wang, L.. An unsupervised feature selection framework based on clustering. In: *New Frontiers in Applied Data Mining*. 2008,.
- [13].Morita, M., Oliveira, L.S., Sabourin, R.. Unsupervised feature selection for ensemble of classifiers. In: *Frontiers in Handwriting Recognition*. 2004,.
- [14].Handl, J., Knowles, J.. Feature subset selection in unsupervised learning via multiobjective optimization. *International Journal of Computational Intelligence Research* 2006;2(3):217–238.
- [15].Dash, M., Liu, H.. Unsupervised feature selection. In: *In Proc. Pacific Asia Conf. Knowledge Discovery and Data Mining*. 2000,.

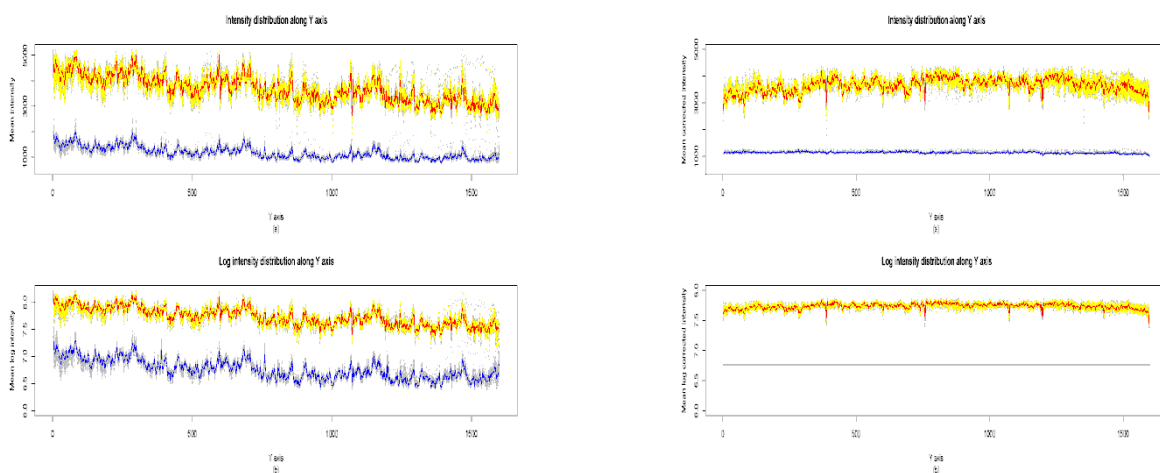


Figure 1. Example of a TWO-COLUMN figure caption: (a) this is the format for referencing parts of a figure.