

# Optimal Predictive Model for Large Scale Classification

Anna Joshy, Leya Elizabeth Sunny, Linda Sara Mathew

**Abstract:** Biosensors calculate the expression pattern of multiple genes in experimental work. A unique genomic chip is possible to produce levels of expression from multiple genes. The ability to interpret these high-dimensional samples fuels the creation of methods of automated analysis. Even though the existing methods undergo imbalanced problems and less classification accuracy over gene expression datasets. Therefore, a novel computational method has been developed in order to increase the classification performance of gene expression dataset and accurate disease prediction. By adding fuzzy memberships, we take into account the features of imbalanced data. Within our work, both the sample entropies and the expense for each class decide the fuzzy memberships in order to understand the different samples with various contributors to the judgment boundary. Thus, on imbalanced genomic datasets, the current proposed approach will result in more desirable classification outcomes. In addition, to build a new algorithm, we integrate the fuzzy memberships into current MKL. The results show that the proposed approach will tackle the imbalanced problem and achieve high accuracy rate.

**Keywords:** Biosensors, Imbalance problem, Fuzzy membership, Multiple kernel learning

## I. INTRODUCTION

Genomic arrays are a significant health treatment modality since they represent a cell's condition at the molecular scale. Usable genomic samples for categorizing tumor forms usually have a relatively limited sampling size due to the large number of genes concerned. This truth is referred to as a dimensionality curse, which is a hard issue. Gene selection is a fruitful method of solving this problem and plays an important role in determining a correct cancer classification, since only a small set of genes are related to the classification task. Gene selection resolves many challenges in microarray datasets, such as reducing the number of insignificant and disruptive genes and selecting the most genes associated. Over decades, the features of gene expression data have stayed almost the same. Small sample sizes along with class imbalance are daunting problems to be taken care of from these features of high dimensionality. Owing to its subjectivity in real life and scientific analysis, unstable issues have become one of the most serious matters

of machine learning and data mining. When we try to solve an uncommon but significant situation where the number of samples in one class is much lower than that of the others, there is an imbalanced problem. For eg, forgery, disease recognition and access management are common issues that are imbalanced. In forgery, only a limited percentage of regular activities are compensated for in cases of fraud. In addition, much of the time it manages the requests from relatives for the permission control scheme, although the records of outsiders are uncommon. In fact, misclassifying a stranger as a family member is far more complicated for the doorlocker scheme than misidentifying a family member as a stranger. Obviously, when coping with such issues, separate groups should be paid with varying attentions. To demonstrate exactly the imbalanced issues, the class with a large sample size is called the dominant class or negative class, while the class with less samples is the minor class or positive class. Typically, the positive class follows the dilemma of data scarcity, which includes actual scarcity and relative rarity. True shortage has to do with the indicates that the amount of positive cases is too low to fully reflect the sample creation. Furthermore, relative rarity shows that while there are a lot of samples in the positive class, their share is very tiny relative to that of the minority examples. While standard algorithms have optimal impact on balanced datasets, in imbalanced problems, they usually have a low positive class recognition accuracy. To counter this, there are two methods that are most widely used. One is the method to the database level, which relates to the techniques of sampling. However, there are certain clear disadvantages to the sampling process. The fuzzy memberships were used to support vector machines for much of the current work. But in this job, to change its meaning in the MKL context, we apply fuzzy memberships to all training set. On the initial training results, the fuzzy memberships are determined. The features of imbalanced data are calculated according to the distribution of samples and properly considered. Therefore, not only does the suggested approach pay more respect to the optimistic class, but it also contributes to a more rational judgment cap on imbalanced datasets. In fact, to deal with imbalanced issues, this approach is easily extended to other models. In addition, to form the algorithm, we apply the fuzzy memberships to current MKL. The following is the remainder of the article. In this area, Segment ii presents the relevant works. The proposed architecture is explained briefly in Segment iii. The findings obtained are defined in Segment iv, and the last portion concludes the article.

Manuscript received on January 07, 2021.

Revised Manuscript received on January 18, 2021.

Manuscript published on January 30, 2021.

\* Correspondence Author

**Anna Joshy\***, Department of Computer Science and Engineering, Mar Athanasius College of Engineering, Kothamangalam, Kerala, India. Email: [annazzzz21@gmail.com](mailto:annazzzz21@gmail.com)

**Leya Elizabeth Sunny**, Department of Computer Science and Engineering, Mar Athanasius College of Engineering, Kothamangalam, Kerala, India. Email: [leyabijoy@gmail.com](mailto:leyabijoy@gmail.com)

**Linda Sara Mathew**, Department of Computer Science and Engineering, Mar Athanasius College of Engineering, Kothamangalam, Kerala, India. Email: [lindasaramathew@gmail.com](mailto:lindasaramathew@gmail.com)

II. RELATED WORK

For Rukshan Batuwita and Vasile Palade[1], the Support vector machine was adaptive to both the issues of class imbalance and outliers/noise in 2010. Even though, Algorithm can overcome the class imbalance problem with the new CIL techniques available, they can also be vulnerable to the outliers problem noise and voice. In another side, although the FSVM solution may address the problem of outliers/noise, since there is no improvement in the FSVM method, it may suffer from the problem of class imbalance. To make it less sensitive compared to the original SVM algorithm to class imbalances. Presents an addition in order to resolve the issue of class imbalance for the FSVM learning process, which is named FSVM-CIL. A new fuzzy membership evaluation that specifies the fuzzy membership evaluation is proposed in the Entropy-based fuzzy support vector machine for imbalanced datasets[2]. Membership on the basis of class accuracy of the samples. That is, broader fuzzy memberships are distributed with greater class accuracy to the observations. Since entropy is used to measure class assurance, the fuzzy membership calculation is referred to as an entropy-based fuzzy membership test. The Entropybased FSVM (EFSVM) is proposed using the entropy-based fuzzy membership. EFSVM can pay more attention to the samples with higher class accuracy. In the sub-space weighting co-clustering of gene expression data[3], a gene feature space weight matrix is incorporated to characterize the contribution of gene objects to the differentiation of various sample clusters. We are creating a new co-clustering optimal solution to retrieve the co-clusters in the gene expression profiles in which the feature space weight matrix is used. Every observation is viewed as negative and positive groups in a new fuzzy SVM to measure credit risk[4]. The SVM fuzzifies the punishment word to reduce the vulnerability of less relevant data sets.

III. PROPOSED WORK

The Alizadeh samples explains the genomic expression information of the three most common adult lymphocytic tumors in the suggested method: diffuse large B-cell lymphoma (DLBCL), follicle lymphoma (FL) and chronic lymphocytic leukemia (CLL). For 59 samples, the dataset comprises 2234 expression profiles: 40 DL- BCL, 9 FL, and 10 CLL.

A. Fuzzy membership generation

Fuzzy memberships, in fuzzy logic, reflect the amount of truth as an expansion of value. In fact, in most cases, the results of training samples are different. That means that in the ranking, some of the training data are more significant than others. It is usually important that meaningful training data must be appropriately categorized. However, whether or not they are misclassified, certain training samples such as outliers and sounds wouldn't care. That means that certain training samples are not exactly part of a specific class anymore. 70% belong to one group and 30% may be insignificant. In coparison, 20% of them may belong to one category and 90% be insignificant. In other terms, each of the training  $a_i$  training samples accompines  $0 < w_i < 1$  with

a fuzzy membership. Therefore, in a classification, the fuzzy membership  $w_i$  can be regarded as the behaviour of the corresponding training data  $a_i$  against one group and the valuation  $(1-s_i)$  can be considered to be the stance of unnecessary.

Provided the training samples  $(a_i, b_i)_i = 1$  where N is the number of training samples,  $b_i = 1$  is the sample  $a_i$  of the positive class, and  $b_i = - 1$  is the sample  $a_i$  of the negative class. For each specimen, we measure the entropy at the beginning. The entropy is estimated for the sample  $a_i$  as in Equation 1:

$$T_i = \sum_{m=1}^c P^m(a_i) \ln P^m(a_i) \quad (1)$$

Where c represents the number of groups, and the likelihood that  $x_i$  belongs to class m is  $P^m(a_i)$ . The probability is generally determined by the adjacent neighborhood samples. But we use the nearest neighborhood search approach with no variables.

$$FRNN(a_i, a_{all}, R) = \{ a_{candi} \in a_{all}, d(a_i, a_{candi}) \leq R \}$$

An instance of the relationship between entropy and sample distribution is given, as seen in Figure 1.. The points in the red circle and the blue triangle represent samples in distinct categories. From the figure, we can get T1 for a 1; 0.43 for a1, and T2 for x2 is 0.57. The entropy of a1 is lower than the entropy of a2, which states that a 1 is as certain as a2. Based on the sample distributions, a2 is closer to the margin of the positive and negative groups than a1 . We would also find that samples closer to the margin between positive and negative classes that are prone to noise and have less certainty will have higher entropy. Suppose that  $E_1, E_2, \dots, E_N$ , for all training samples is the entropy. The first aim is to define the fuzzy membership for each sample, per the entropy. For a training sample,  $x_i$ , fuzzy membership  $0 < s_i < 1$  is given that is the behaviour of  $a_i$  towards a specific group. For each subclass, the fuzzy membership is determined as follows Equation 3

$$F_m = 1 - \alpha(m - 1), j = 1, 2, \dots, n \quad (3)$$

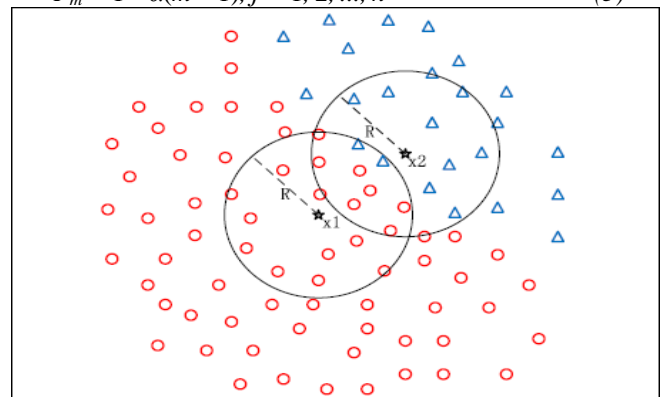


Figure 1: The FRNN process determination

B. Architecture of CFMKL

A more modular kernel approach called MKL, which utilizes the combination of various kernels, has been magnet in machine learning in recent years[5]. Two key sections are included in the objective function for algorithm and can be defined as in Equation 4.



$$J = \sum_{K=1} (U_{emp}^k(D) + c^k U_{reg}^k) + \gamma(D) \quad (4)$$

The kernel  $k$  represents the  $k$ th kernel where  $K$  is the number of kernels accepted. In the first section, two terms are used, including the analytical risk term  $U^k(D)$  and also the structure risk term  $U_{reg}$ . In this section, the optimization objective of each base classifier is represented and guarantees the classification accuracy of the classifier. The next term maintains the output of each kernel as close to the average output of all kernels as appropriate.

#### IV. RESULT AND DISCUSSION

In studies, microarrays simultaneously calculate the expression of multiple genes. A unique genomic chip is capable of producing levels of expression from thousands and thousands of genes and the information is commonly obtained from different tissue datas, with varying ecological factors, in multiple patients[6]. The need for such high-dimensional data to be processed is driving the advancement of automatic processing methods. A novel statistical technique has been developed to increase the classification performance of the dataset of gene expression and accurate disease prediction, in order to increase the classification performance of the dataset of gene expression and accurate disease prediction. To address this problem, by implementing fuzzy memberships, we take into account the characteristics of imbalanced data. In our work, in order to create distinct kinds have distinct contributions to the decision boundary, both the sample entropies and the cost for each class determine the fuzzy memberships. Therefore, the recently current technique could consist in a more favorable identification of performance on imbalanced datasets. In addition, we incorporate fuzzy memberships into current MKL in order to create a new algorithm called CFMKL in brief. Experimental findings in Figure 2 confirm the great efficacy of the proposed CFMKL on virtual, real-world binary and multiclass imbalanced dataset outcomes.

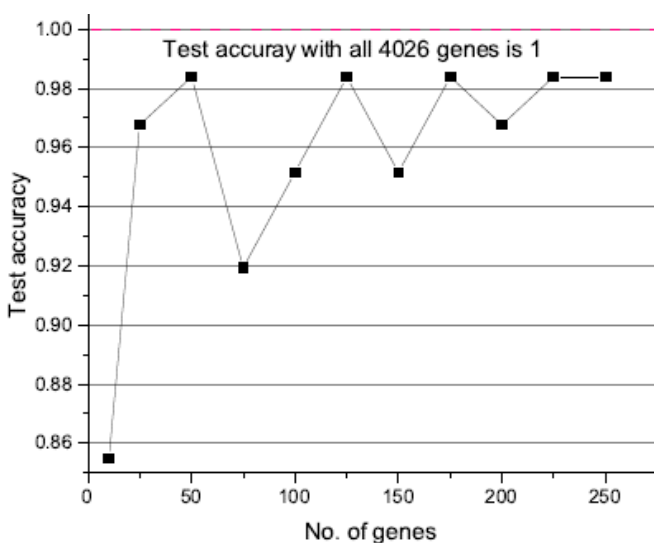


Figure 2. Test accuracy using CFMKL over gene expression dataset

#### V. CONCLUSION

MKL seeks to solve issues that are imbalanced. While current MKL algorithms have certain advantages over imbalanced problems[7][8], they are not able to take into account the characteristics of imbalanced data. We are proposing a new algorithm to solve this issue. Based on the current MKL algorithm, the approach employs the fuzzy memberships in our paper. It determines the first of all training samples for fuzzy memberships. The fuzzy memberships are being used to determine the relevance of the samples to the classification model in particular[9]. Furthermore, to measure the fuzzy memberships to match the imbalanced problems, we have built a newly function. The results suggest that the current algorithm simply[10][11] gives more intentions to the correctly classified and reduces the effects of unknown samples susceptible to noise[12][13]. Our newly built algorithms average classification performance is 7.43 percent higher than that of the initial MKL on imbalanced binary type datasets and the results of new algorithm rank first among all algorithms compared[9]. Compared with MKL and kernel dependent algorithms, the proposed algorithm still achieves the best results on multi-class imbalanced datasets[14]. Therefore, the findings of the experiment demonstrate that the algorithm outperforms the other algorithms, which implies that our proposal approach is a feasible and efficient way of coping with the imbalanced issues[15].

#### REFERENCES

1. R. Batuwita, V. Palade (2010) "FSVM-CIL: Fuzzy Support Vector Machines For Class Imbalance Learning", IEEE Transaction. Fuzzy syst. 18(3) 558-517
2. Q. Fan, Z. Wang, D. Li, D. Gao, H. Zha, (2017) "Entropy-based fuzzy support vector machine for imbalanced datasets", Knowl. Based Syst. Elsevier 11587-99.
3. Xiaojun Chen, Joshua Z. Huang 2017, Subspace Weighted Co-Clustering of gene expression data, IEEE/ACM transaction on computational biology and bioinformatics
4. Y. Wang, S. Wang, K. Lai, "A new fuzzy support vector machine to evaluate credit risk", IEEE Trans. Fuzzy Syst. 13 (6) (2012) 820-831
5. A. Rakotomamonjy, F. Bach, S. Canu, Y. Grandvalet, Simplemkl, J. Mach. Learn. Res. 9 (11) (2008) 2491-2521.
6. M. Gonen, E. Alpaydin, Multiple Kernel Learning algorithms, J. Mach. Learn. Res. 12(7) (2011) 2211-2268.
7. Z. Xu, R. jin, H. Yang, I. King, M. Lyu, Simple and efficient multiple kernel learning by group lasso, in: Proceedings of the 27<sup>th</sup> International Conference on Machine Learning, 2010, pp. 1175-1182
8. X. Xu, I. W. Tsang, D. Xu, Soft margin multiple kernel learning, IEEE Trans. neural Netw. learn. Syst. 24(5) (2013) 749-761.
9. Q. Mao, I. W. Tsang, S. Gao, L. Wang, Generalized multiple kernel learning with data dependent prior, IEEE Trans. Neural Netw. Learn. Syst. 26(6) (2015) 1134-1148
10. J. Tang, Y. Tian, A multi-kernel framework with nonparallel support vector machine, Neurocomputing 266(2017) 226-238
11. G. Govaert and M. Nadif, Co-clustering. Wiley-ISTE, 12 2013, 31
12. L. Breiman, "Random forest", machine learning, vol. 45, no: 1
13. S. Boyd and L. Vandenberghe, convex optimization. Cambridge university press, 2004
14. R. Diaz-Uriarte and S. A. De Andres, "Gene selection and classification of microarray data using random forest", BMC bioinformatics, vol. 7, no: 1, 2006
15. L. Lazzaroni and A. Owen, "Plaid models for gene expression data", 2002

## AUTHORS PROFILE



**Anna Joshy** received Bachelor of Technology in Computer Science and Engineering from Younus College of Engineering in 2018 and currently pursuing Master of Technology in Computer Science and Engineering from Mar Athanasius College of Engineering, Kothamangalam affiliated to APJ Abdul Kalam Technological University. Her research interest is in machine learning and

data mining



**Prof. Leya Elizabeth Sunny** is currently working as assistant professor in the Department of Computer Science and Engineering, Mar Athanasius College of Engineering, Kothamangalam, Kerala, India. She received her B-Tech Degree in 2004 in computer science and engineering from MG University and M-Tech in 2011 in Information System Security. She has around 12 years of teaching and industrial experience. She is interested in the areas of Data Security and Cryptography.



**Prof. Linda Sara Mathew** is currently working as assistant professor in the Department of Computer Science and Engineering, Mar Athanasius College of Engineering, Kothamangalam, Kerala, India. She received her B-Tech Degree in 2002 and M-Tech in 2011 in Computer Science and Engineering from Mahatma Gandhi university and Anna University. She has around 15 years of teaching and industrial experience. She is interested in the areas of Data Mining, Neural Network, Image Processing, Soft Computing