

# An Efficient Classification Approach for Examine Health Records of Cause of Death

Simma Narendra Prasad<sup>1</sup>, Konni Srinivasa Rao<sup>2</sup>

Final M.Tech Student<sup>1</sup>, Asst.professor<sup>2</sup>

<sup>1,2</sup>Dept of CSE, Sarada Institute of Science, Technology and Management (SISTAM), Srikakulam, Andhra Pradesh

## Abstract:

Now a day's examining the health of each person in the every country is an integral part of healthcare. After examining the health of each person we can identify type of risk to be occurred. The analysis of risk based unlabelled data can be done by using classification approach in the data mining. Particularly we are take unlabelled data contains information related to participants in the health examination whose health condition is vary from great health to very ill. In this study we formulated the task of risk prediction as a multi-class classification problem using the Cause of Death (COD) information as labels, regarding the health-related death as the "highest risk". The goal of risk prediction is to effectively classify 1) whether a health examination participant is at risk, and if yes, 2) predict what the key associated disease category is. In other words, a good risk prediction model should be able to exclude low-risk situations and clearly identify the high-risk situations that are related to some specific diseases. In the examination of health we are identifying different states of health without ground truth. So that by predicting risk of each participant by using classification approaches in the data mining. In this paper we proposed Mixed Probability Binary Rule Based Classification Algorithm is used to predict health risk of participate person. By implementing this algorithm we can get efficient classification result and also give better performance.

## Keywords:

data mining, Predict Analysis, Classification, Electronic Medical Records, Health Examination Records.

## I. INTRODUCTION

Nowadays, the data accumulated in medical databases are progressively growing up quickly, this makes extracting hidden knowledge from medical database complex and more time consuming. Analyzing these data is critical for medical decision makers and managers. The performance of patient management tasks will be improved by analyzing the medical data. Enormous Amounts of Electronic Health Records (EHRs) composed over the years

have provided a rich base for risk examination and forecast. An EHR contains numerically warehoused healthcare info about an individual, such as interpretations, laboratory tests, diagnostic reports, medications, procedures, patient identifying information, and allergies. A special type of EHR is the Health Examination Records (HER) from annual general health check-ups. For example, governments such as Australia, U.K., and Taiwan proposal periodic geriatric health examinations as an essential part of their matured care programs. Since clinical care frequently has a specific problem in mind, at a point in time, only a limited and often small set of measures considered necessary are collected and stored in a person's EHR. By contrast, HERs are gathered for consistent investigation and defensive purposes, covering a inclusive set of general health measures, all together at a point in time in a methodical way. Identifying contributors at risk based on their current and past HERs is important for early cautionary and preventive intervention. By "risk", we unsolicited outcomes such as mortality and morbidity. In this study we expressed the task of risk forecast as a multi-class classification problem using the Cause of Death (COD) information as labels, concerning the health-related death as the "highest risk". The goal of risk prediction is to effectively classify 1) whether a health examination participant is at risk, and if yes, 2) predict what the key related disease category is. In other words, a good risk prediction model should be able to exclude low-risk situations and clearly identify the high-risk conditions that are related to certain exact diseases.

Most developed countries have experienced dramatic growth in elderly populations from the beginning of this century. In recent years, this, together with the rising cost of healthcare has created an urgent need for improving predictions and efficient treatment. These programs enable continuous and comprehensive recording of a person's health status, as well as the tracking of his/her health changes. However, it is always a difficult task for healthcare professionals to provide an overall report on personal health after a comprehensive medical check-up has been performed because of hundreds of the parameters

available to be considered. One particular focus of preventive healthcare is risk assessment. The goal is to identify individuals at risk for further investigation or early treatment and intervention. Traditionally, risk assessments have been conducted manually by clinical professionals based on their expertise. These manual assessments have been constrained by the capacity of the human brain to process information within a limited time during the period of an appointment with a patient. Many risk-scoring systems have been developed in the field of medicine to assist clinical decision-making. As a general practice in medical research, these methods have been defined based on factors selected with expert knowledge and validated via population-based studies [1]. With the advances in computing technology and the availability of EHRs, an increasing number of data mining and machine learning applications have been developed to support healthcare decision making [2, 3]. In recent years methods for clinical risk classification have been developed [4, 5, 6, 7, and 8]. However, most existing studies have their focus on EHRs. GHE records and the unique challenges they pose have not yet been well explored. This gap has driven our research to advance risk prediction models for GHE records.

## II. RELATED WORK

In this section we review existing related studies, namely those on mining health examination data and those on classification with unlabeled data in healthcare applications. Although Electronic Health Records (EHRs) have attracted increasing research attention in the data mining and machine learning communities in recent years [9], [10], [11], [12], [13], [14], [15],[16] mining general health examination data is an area that has not yet been well-explored, except a few studies on risk prediction such as the chronic disease early warning system proposed in [17] and our previous work on health score classification framework [13], [18]. However, none of the them considered unlabeled data. In addition, the approach presented in [13] is limited to a binary classification problem (using alive/deceased labels) and consequently it is not informative about the specific disease area in which a person is at risk.

## III. PROPOSED SYSTEM

In this paper we are propose mixed probability binary rule based classification algorithm for predicting health risk of participant. To solve the problem of health risk prediction based on health examination of records of participant. Our algorithm takes health examination data and linked cause of death labels as inputs. Its key components process health examination records and predicting disease class. Before processing health examination records

we are take the training data set contains information related to test result with type of disease class. By taking those dataset as training data set and predicting examination records of participant. Take the more than one record of participant and processing those records for predicting risk. The predicting health records will consider as testing data set for finding type of disease class. To identify type of disease class we proposed mixed probability binary rule based classification algorithm. The algorithm combines the advantages of for class discovery and for handling heterogeneity to solve a specific problem induced by evidence-based risk prediction from health examination records. To train a disease risk prediction model that is capable of identifying high-risk individuals given no ground truth for “healthy” cases, we treated the “unknown” class as a class to be learned from data. We incorporated the class discovery mechanism of into our method to handle the “unknown” class. To handle unknown class we propose mixed probability binary rule based classification algorithm for predict type of disease class. The implementation procedure of mixed probability binary rule base classification algorithm is as follows.

### A) *Mixed Probability Binary Rule Based Classification Algorithm:*

In this module we are implementing mixed probability binary rule based classification algorithm for predicting type of disease classes. By implementing this algorithm we can get best predictive result and also improve the performance. The implementation of steps of mixed probability binary rule base classification algorithm is as follows.

1. Read the training data set contains information related to type of disease class with test results.
2. Read the testing data set contains information related to test result. By taking those testing result we can predict type of disease class.
3. Take the first attribute value from the testing dataset and compare to with training dataset same attribute. If the testing dataset attribute value is greater than equal to training dataset attribute values then put one for that attribute of record. Here we can also consider the testing dataset attribute value is greater than or equal to normal test result we can put one to the particular record.
4. In the comparison process testing data set attribute values less than training dataset attribute value or the testing data set attribute values less than normal test result the put zero as status to particular record.

5. Take the second attribute value from testing data set and perform the step 3, 4 put status of each record with one or zero.

6. This process repeated until the completion of all attributes in training data set and testing dataset values can be converted into in form of zero or one.

7. Take the each attribute value and calculate probability of each attribute related to testing data set attribute.

8. The calculation of probability of each attribute can be done by using following equation.

$Prob_{yes} = \text{Total number of once} / \text{total number of records.}$

$Prob_{no} = \text{Total number of zeroes} / \text{total number of records.}$

9. Calculate each attribute yes, no probability and find out final yes, no probability of each record. The calculation of final yea, no probability is as follows.

$P_{yes} = \text{multiplication of all attributes yes probability.}$

$P_{no} = \text{multiplication of all attributes no probability.}$

10. After completion of probability calculation we can perform the rule based classification process. The rule based classification process contains If then rules predicting disease class.

11. The rule based classification makes use of a set if then rules for classification. We can express the rule in following form.

If condition then conclusion.

12. In the rule based classification we can take if part of rule is called rule antecedent or precondition.

13. The then part of rule based classification is called rule consequent.

14. The antecedent part of condition consists of one or more attribute tests and these tests or logically And.

15. The antecedent part of our process will take condition as probabilities of yes or no and test dataset values of each attributes and normal test result of each attribute.

16. If all condition of each record satisfies particular disease class of training data set attribute values and take those disease class as predict of risk. Then the predict result is consequent to testing data set result and those participant face the type of disease class.

17. Step 16 will be repeated until total completion records in testing data set.

After completion of this process we can get type of disease that participant will face and also get efficient result. By implementing mixed probability binary rule based classification technique we can get best predict result and also improve performance of system. Because in this algorithm we can calculate probability of each attribute and also generate rule for each record. By performing those two operations we can retrieve more related predict result.

#### IV. CONCLUSIONS

Mining health examination data is challenging especially due to its heterogeneity, intrinsic noise, and particularly the large volume of unlabelled data. By examining the unlabelled dataset we are using classification technique the data mining. In this paper we are proposed mixed probability binary rule based classification process for predicting type of disease class. In the proposed system we are calculate each attribute probability related to training dataset and using that probability for predicting disease class. In this paper we can also implement the rule based classification approach for identifying type of risk based disease class. In this project we can take two type of dataset for identifying risk. The first data set is training dataset contains information related type of disease class with related test result of attributes. The second data set testing data set is used to identify predict type of disease class using training dataset. By implementing those two processes we can improve efficiency and also get best predict result with the type of disease class.

#### REFERENCES

- [1] M. Woodward. Epidemiology: study design and data analysis. CRC Press, 2013.
- [2] N. Esfandiary, M. R. Babavalian, A.-M. E. Moghadam, and V. K. Tabar. Knowledge discovery in medicine: Current issue and future trend. *Expert Systems with Applications*, 41(9), Jul. 2014.
- [3] J.-Y. Yeh, T.-H. Wu, and C.-W. Tsao. Using data mining techniques to predict hospitalization of haemodialysis patients. *Decision Support Systems*, 50(2):439–448, Jan. 2011.
- [4] F. Wang, P. Zhang, B. Qian, X. Wang, and I. Davidson. Clinical risk prediction with multilinear sparse logistic regression. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 145–154, New York, USA, 2014. ACM.
- [5] T. Tran, D. Phung, W. Luo, and S. Venkatesh. Stabilized sparse ordinal regression for medical risk stratification. *Knowledge and Information Systems*, 43(3):555–582, mar 2015.
- [6] J. Wiens, E. Horvitz, and J. V. Gutttag. Patient Risk Stratification for Hospital-Associated C. diff as a Time-Series Classification Task. In *Neural Information Processing Systems*, pages 476–484, 2012.
- [7] Y. Mao, W. Chen, Y. Chen, C. Lu, S. Louis, M. Kollef, and T. C. Bailey. An integrated data mining approach to real-time clinical monitoring and deterioration warning. In *Proceedings of the 18th ACM SIGKDD International*

Conference on Knowledge Discovery and Data Mining, pages 1140–1148, Beijing, China, 2012. ACM

- [11] H. Neuvirth, M. Ozery-Flato, J. Hu, J. Laserson, M. S. Kohn, S. Ebadollahi, and M. Rosen-Zvi. Toward personalized care management of patients at risk: the diabetes case study. In SIGKDD, pages 395–403, California, USA, 2011. ACM.
- [12] T. Tran, D. Phung, W. Luo, and S. Venkatesh, “Stabilized sparse ordinal regression for medical risk stratification,” Knowledge and Information Systems, pp. 1–28, Mar. 2014.
- [13] M. S. Mohktar, S. J. Redmond, N. C. Antoniadis, P. D. Rochford, J. J. Pretto, J. Basilakis, N. H. Lovell, and C. F. McDonald, “Predicting the risk of exacerbation in patients with chronic obstructive pulmonary disease using home telehealth measurement data,” Artificial Intelligence in Medicine, vol. 63, no. 1, pp. 51–59, 2015.
- [14] J. M. Wei, S. Q. Wang, and X. J. Yuan, “Ensemble rough hypercuboid approach for classifying cancers,” IEEE Transactions on Knowledge and Data Engineering, vol. 22, no. 3, pp. 381–391, 2010.
- [15] E. Kontio, A. Airola, T. Pahikkala, H. Lundgren-Laine, K. Junttila, H. Korvenranta, T. Salakoski, and S. Salanter’a, “Predicting patient acuity from electronic patient records,” Journal of Biomedical Informatics, vol. 51, pp. 8–13, 2014.
- [17] Q. Nguyen, H. Valizadegan, and M. Hauskrecht, “Learning classification models with soft-label information,” Journal of the American Medical Informatics Association : JAMIA, vol. 21, no. 3, pp. 501–8, 2014.
- [18] G. J. Simon, P. J. Caraballo, T. M. Therneau, S. S. Cha, M. R. Castro, and P. W. Li, “Extending Association Rule Summarization Techniques to Assess Risk of Diabetes Mellitus,” IEEE Transactions Knowledge and Data Engineering, vol. 27, no. 1, pp. 130–141, 2015.
- [19] L. Chen, X. Li, S. Wang, H.-Y. Hu, N. Huang, Q. Z. Sheng, and M. Sharaf, “Mining Personal Health Index from Annual Geriatric Medical Examinations,” in 2014 IEEE International Conference on Data Mining, 2014, pp. 761–766.
- [20] S. Pan, J. Wu, and X. Zhu, “CogBoost: Boosting for Fast Costsensitive Graph Classification,” IEEE Transactions on Knowledge and Data Engineering, vol. 6, no. 1, pp. 1–1, 2015.
- [21] Y. Zhao, G. Wang, X. Zhang, J. X. Yu, and Z. Wang, “Learning phenotype structure using sequence model,” IEEE Transactions on Knowledge and Data Engineering, vol. 26, no. 3, pp. 667–681, 2014.
- [22] L. Chen, X. Li, Y. Yang, H. Kurniawati, Q. Z. Sheng, H.-Y. Hu, and N. Huang, “Personal health indexing based on medical examinations: A data mining approach,” Decision Support Systems, vol. 81, pp. 54 – 65, 2016.