### Journal of Science, Computing and Engineering Research (JSCER)

### Volume-7, Issue-11, November 2024.

DOI: https://doi.org/10.46379/jscer.2023.071102

# Implementation of Feature Selection and Random Forest Machine Learning Classification Algorithm on a Medical Based Datas

Swanti Jodhi, Prof. S.K Dubey, Prof. Shasi kant Shrivastava

LDRP Institute of Technology and Research, Gandhinagar, Gujarat

### **Article Information**

Received

Revised : 06 Nov 2024 Accepted : 10 Nov 2024

02 Nov 2024

Published : 14 Nov 2024

Corresponding Author:

Swanti Jodhi

Abstract— For a given medical data set, there is a huge possibility that the data includes hundreds of features, each representing a symptom or a parameter based on which diagnosis can be carried out. While a lot of these features contribute towards the results, it is often the case that quite a few of these features turn out to be either irrelevant or have very little bearing in terms of their overall impact on the results and only end up crowding the data set. Feature selection provides a solution to this problem as the features that provide the highest contribution while predicting an output are retained and the irrelevant features are identified and subsequently eliminated. This helps in the model being trained faster and leads to a better interpretation of the model further allowing better diagnosis of the disease. Apart from feature selection, random forest classifier is being used as a means to predict the outcomes. Since random forest is made up of decision trees, it helps in better classification for a given problem.

Keywords: Feature Selection, Random Forest, Overfitting, Medical Dataset, Classification

Copyright © 2024: Swanti Jodhi, Prof. S.K Dubey, Prof. Shasi kant Shrivastava, This is an open access distribution, and reproduction in any medium, provided Access article distributed under the Creative Commons Attribution License the original work is properly cited License, which permits unrestricted use.

**Citation:** Swanti Jodhi, "Implementation of Feature Selection and Random Forest Machine Learning Classification Algorithm on a Medical Based Datas", Journal of Science, Computing and Engineering Research, 7(11), November 2024.

#### I. INTRODUCTION

[1] Machine learning can be applied to a wide variety of fields and there is always scope to come up with great insights and solutions for a given field with the help of machine learning. The approach to a solution naturally varies, based on the field of application and the type of problem at hand. The medical field is one such area that has greatly benefited from machine learning. Since medicine is a field of gargantuan proportions, the applications in the field also exist in a wide range.

This wide range of applications include enhancement in the ways in which health related data is managed, predicting illnesses at earlier stages by identifying various disease markers and patterns that pose a health risk, improving the accuracy of the diagnosis for a variety of diseases, etc.

[4]While the application may vary, the common goal for using machine learning in the field of medicine, most of the times, is to provide better health care at a lower cost or/and better performance and accuracy.

[4] This paper deals with the aspect of large medical datasets that contain a multitude of attributes and need a reduction in the number of attributes as some of them do not contribute much towards the resulting outcome. Since the output variable of the dataset in question concerns itself with prognosis, the problem comes under the umbrella of a

classification problem. This is because the diseases covered by the prognosis attribute are to be classified into different categories. With way too many columns dictating the prediction of a disease, the processing speed and efficiency of the model is bound to be less than optimum. [1] Feature selection helps in getting rid of attributes of less importance and makes the dataset more compact and precise. The seeming fixation with lesser number of attributes can be justified by the fact that the complexity of the model is reduced as a result of feature selection which makes it easier for the model to be interpreted and explained.

Moreover, feature selection provides a solution for the problem of overfitting. Overfitting is an issue that arises when the model learns about more data than is necessary, i.e., the noise or the unrelated data is identified by the model as necessary data that is to be incorporated in the learning along with the useful data. This leads to inaccuracy in predictions as the factors that do not have an impact on the outcome are believed to be of importance.

[3] There are various methods for applying feature selection but these methods can be classified into three categories: filter, wrapper and embedded. Filter feature selection methods such as chi-squared test and information gain assign a score to each feature through statistical means and based on a certain value or threshold the features are

# Implementation of Feature Selection and Random Forest Machine Learning Classification Algorithm on a Medical Based Datas

Available at https://jscer.org

either selected or discarded. Wrapper methods look for combinations of features and constantly evaluate and compare them with other feature combinations. Embedded methods carry out feature selection when the model is being trained.

We have implemented feature selection by calculating feature importance of all attributes using forest of trees. Finally, random forest classifier is being used for the prediction of diseases. A random forest is a collection of lots of individual trees. Each one of them varying slightly in comparison to the other as each one of them works on a slightly different set of observations.

The output of all the individual trees is combined to form a final prediction. This finalprediction is a result of the feature selection and the random forest algorithm that have helped create a model that is in quite a simplified form and helps in better diagnosis of diseases.

#### II. RELATED WORKS

The paper [1] lays down the very basics of feature selection and goes on to explain the concept in detail by covering various areas related to it. The papers primary focus lies in creating subsets of features that help in creating a good predictor. Variable Ranking is discussed as a principle selection mechanism and mathematical explanations are provided for various techniques through which it can be carried out such as Correlation Criteria, Single Variable Classifiers, Information theoretic ranking criteria, etc.

The paper served as an apt introduction to variable and feature selection and helped in acquiring a fundamental understanding about it. The paper [2] talks about the implementation of mRMR (Minimum Redundancy Maximum Relevance) feature selection method in marketing, focusing on the mathematical aspect of mRMR, while also covering real data examples. Implementation in Production is discussed in terms of Architecture of the Platform, Challenges and Optimization and Online Experiment Evaluation.

The implementation of a specific feature selection method on a particular domain helped in understanding the overall approach that is to be taken while implementing a feature selection method to the project. It highlights the parameters that are to be kept in mind during ideation and in its subsequent implementation.

The paper [3] discusses the use of machine learning techniques in analyzing data. It gives a brief idea about various feature selection methods namely principal component analysis, factor analysis and attribute ranker. It sheds light on multiple feature selection methods providing us a basic understanding of each one of them and how they can be useful for analyzing data in thefield of medicine. The

paper [4] talks about data mining in medical dataset and applying feature selection for classification. Similar to [3], the paper highlights various feature selection methods for analyzing medical data.

The difference being, that unlike [3], this paper also discusses the methodology apart from the concepts of those methods. Hence, it helps in understanding the approach that is to be taken, in order to implement feature selection in the domain of medicine. The paper [5] deals with Multilayer perceptronbased feature selection algorithm.

The backpropagation algorithm trains the multilayer perceptron to determine the attributes to be removed from the data set. The concept of prominence is used for the objective function for the feature selection algorithm. It basically indicates the real relevance of an attribute for a given task. While this paper takes a different approach by using an MLP based feature selection method but the basic objective of eliminating redundancy and irrelevance matches our objective.

The paper [6] discusses feature ranking and feature selection for a linear model. It provides a detailed procedure for ranking a feature, which is an essential step during the implementation of a feature selection method. The entire mathematical procedure is explained in detail, which helps in understanding not just the theoretical aspect but also its actual implementation in various applications.

#### III. METHODOLOGY

Medical Dataset Selection The primary goal while looking for a dataset was to select one that provided a great deal of information. At the same time, it was important to keep in mind that the information gained from the data needed to be comprehensible.

Initially, we leaned towards datasets that have information to predict a particular disease. Most of the datasets that we found using this approach led us to some pretty common datasets dealing with heart related diseases and breast cancer. These datasets had already seen way too many implementations and besides that, the number of columns in it didn't seem to be sufficient enough. We ended up choosing a dataset that focused on multiple diseases with enough data available on them.

The dataset that we worked on has 133 columns in total with 132 of those columns signifying various symptoms and the output variable prognosis mentioned the disease that was caused due to these symptoms. Every row signified a patient record and had values of either 1 or 0 which means symptom exists and symptom doesn't exist respectively for 132 columns signifying various symptoms and the last entry for every row stated the disease as mentioned earlier. The number of unique values, i.e., the number of different diseases in the prognosis attribute are 41. With 4920 rows of

# Implementation of Feature Selection and Random Forest Machine Learning Classification Algorithm on a Medical Based Datas

Available at https://jscer.org

information for training the model, the dataset is quite useful for the implementation of feature selection. The only bit of data pre-processing done was to classify the string values of the diseases to numbers, for the model to conveniently work on. Feature Selection With 132 different symptoms existing in the dataset, the objective was to find and eventually eliminate the features that did not contribute much to the final outcome.

The implementation of feature selection on the dataset used was done using random forests [7]. This is because random forests use tree-based strategies to naturally rank elements based on how well they improve the purity of a particular node. At the start of the trees, nodes with the greatest decrease in impurity are found.

Whereas, at the end of the trees, nodes with the least decrease in impurity are found. The impurity mentioned here is the gini impurity. Used for classification trees, gini impurity is a measure of the frequency of a randomly chosen element from the set being incorrectly labelled if it is being randomly labelled. Hence, the subset of important features is created by pruning trees below a particular node.

Table 1 shows the importance values of 10 features (symptoms) as an example. The selection criteria for a feature was set as the importance value for the given feature being more than or equal to the average importance of a feature.

Hence, all features whose importance value was less than the average value was discarded. Random Forest Classifier The classifier was used initially so as to apply feature selection[8] using forest of trees[9]. Once a subset of features was obtained from the original set of features through feature selection, the classifier was trained for a second time in order to run the model with the new dataset consisting of the retained features. We used a forest of 120 trees in our classifier. Individual trees do not provide accuracy for predictions.

A random forest consisting of largely uncorrelated trees come up with accurate predictions. The existence of uncorrelated trees is ensured through the concept of bagging which allows individual trees to randomly take samples for training from the dataset rather than allocating the same data to every tree. For example, let's assume that a training data has 5 rows. Instead of allowing a tree to train using those 5 rows, the tree is randomly allotted two instances of row 2, and 1 instance each of row 1, row 3 and row 5. Similarly, other trees are randomly allotted slightly different data resulting in each tree being slightly different than the others.

Since the overall prediction is based on the prediction of majority of trees, the errors of individual trees get overshadowed by the accurate prediction of other trees. These characteristics of random forest[10] classifier informed our decision to use it for our model.

#### IV. RESULTS AND DISCUSSIONS

Once the importance of all features was calculated using feature selection [11], the features having an importance less than the average importance were to be discarded. This resulted in close to 70 features being removed out of 133. A lot of those features hardly featured as symptoms for most of the diseases in the dataset. Once random forest [12] classifier was applied on the existing features, we fed certain symptoms to test our model. For example, on entering symptoms of diarrhoea and vomiting, the model predicted Gastroenteritis as the condition with an accuracy of 85%. On providing itching and blister as the symptoms, Drug Reaction was the prognosis made by the model with 95% accuracy. While testing for Hepatitis C an accuracy of 90% was obtained.

#### V. CONCLUSION

With data getting bigger by the day, it has become a necessity that people keep up with it and try to get useful insights from it as much as is possible. [4] Looking at the field of medicine, there is so much scope for the usage of machine learning, right from providing better care at low cost to making higher accuracy predictions for diseases. Feature Selection is just one of the many tools available, that helps achieve the aforementioned objectives. Yet it's relevance can be credited to the fact that it allows for faster processing of data, better interpretation of a model and better accuracy for predictions.

### REFERENCES

- [1]. P. Nirmala, T. Manimegalai, J. R. Arunkumar, S. Vimala, G. Vinoth Rajkumar, Raja Raju, "A Mechanism for Detecting the Intruder in the Network through a Stacking Dilated CNN Model", Wireless Communications and Mobile Computing, vol. 2022, Article ID 1955009, 13 pages, 2022. https://doi.org/10.1155/2022/1955009.
- [2]. D. Sathyanarayanan, T. S. Reddy, A. Sathish, P. Geetha, J. R. Arunkumar and S. P. K. Deepak, "American Sign Language Recognition System for Numerical and Alphabets," 2023 International Conference on Research Methodologies in Knowledge Management, Artificial Intelligence and Telecommunication Engineering (RMKMATE), Chennai, India, 2023, pp. 1-6, doi: 10.1109/RMKMATE59243.2023.10369455.
- [3]. J. R. Arunkumar, Tagele berihun Mengist, 2020" Developing Ethiopian Yirgacheffe Coffee Grading Model using a Deep Learning Classifier" International Journal of Innovative Technology and Exploring Engineering (IJITEE) ISSN: 2278-3075, Volume-9 Issue-4, February 2020. DOI: 10.35940/ijitee.D1823.029420.
- [4]. Ashwini, S., Arunkumar, J.R., Prabu, R.T. et al. Diagnosis and multi-classification of lung diseases in CXR images using optimized deep convolutional neural network. Soft

# Implementation of Feature Selection and Random Forest Machine Learning Classification Algorithm on a Medical Based Datas

Available at https://jscer.org

- Comput (2023). https://doi.org/10.1007/s00500-023-09480-3
- [5]. J.R.Arunkumar, Dr.E.Muthukumar," A Novel Method to Improve AODV Protocol for WSN" in Journal of Engineering Sciences" ISSN NO: 0377-9254Volume 3, Issue 1, Jul 2012.
- [6]. R. K, A. Shameem, P. Biswas, B. T. Geetha, J. R. Arunkumar and P. K. Lakineni, "Supply Chain Management Using Blockchain: Opportunities, Challenges, and Future Directions," 2023 Second International Conference on Informatics (ICI), Noida, India, 2023, pp. 1-6, doi: 10.1109/ICI60088.2023.10421633.
- [7]. Arunkumar, J. R. "Study Analysis of Cloud Security Chanllenges and Issues in Cloud Computing Technologies." Journal of Science, Computing and Engineering Research 6.8 (2023): 06-10.
- [8]. J. R. Arunkumar, R. Raman, S. Sivakumar and R. Pavithra, "Wearable Devices for Patient Monitoring System using IoT," 2023 8th International Conference on Communication and Electronics Systems (ICCES), Coimbatore, India, 2023, pp. 381-385, doi: 10.1109/ICCES57224.2023.10192741.
- [9]. S. Sugumaran, C. Geetha, S. S, P. C. Bharath Kumar, T. D. Subha and J. R. Arunkumar, "Energy Efficient Routing Algorithm with Mobile Sink Assistance in Wireless Sensor Networks," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ACCAI58221.2023.10201142.
- [10] R. S. Vignesh, V. Chinnammal, Gururaj.D, A. K. Kumar, K. V. Karthikeyan and J. R. Arunkumar, "Secured Data Access and Control Abilities Management over Cloud Environment using Novel Cryptographic Principles," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023, pp. 1-8, doi: 10.1109/ACCAI58221.2023.10199616.
- [11]. Syamala, M., Anusuya, R., Sonkar, S.K. et al. Big data analytics for dynamic network slicing in 5G and beyond with dynamic user preferences. Opt Quant Electron 56, 61 (2024). https://doi.org/10.1007/s11082-023-05663-2
- [12].Krishna Veni, S. R., and R. Anusuya. "Design and Study Analysis Automated Recognition system of Fake Currency Notes." Journal of Science, Computing and Engineering Research 6.6 (2023): 16-20.
- [13]. V. RamKumar, S. Shanthi, K. S. Kumar, S. Kanageswari, S. Mahalakshmi and R. Anusuya, "Internet of Things Assisted Remote Health and Safety Monitoring Scheme Using Intelligent Sensors," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023, pp. 1-8, doi: 10.1109/ACCAI58221.2023.10199766.
- [14].R. S. Vignesh, R. Sankar, A. Balaji, K. S. Kumar, V. Sharmila Bhargavi and R. Anusuya, "IoT Assisted Drunk and Drive People Identification to Avoid Accidents and Ensure Road Safety Measures," 2023 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ACCAI58221.2023.10200809.
- [15].I. Chandra, G. Sowmiya, G. Charulatha, S. D, S. Gomathi and R. Anusuya, "An efficient Intelligent Systems for Low-Power Consumption Zigbee-Based Wearable Device for Voice Data Transmission," 2023 International Conference on Artificial

- Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ICECONF57129.2023.10083856.
- [16].G. Karthikeyan, D. T. G, R. Anusuya, K. K. G, J. T and R. T. Prabu, "Real-Time Sidewalk Crack Identification and Classification based on Convolutional Neural Network using Thermal Images," 2022 International Conference on Automation, Computing and Renewable Systems (ICACRS), Pudukkottai, India, 2022, pp. 1266-1274, doi: 10.1109/ICACRS55517.2022.10029202.
- [17].R. Meena, T. Kavitha, A. K. S, D. M. Mathew, R. Anusuya and G. Karthik, "Extracting Behavioral Characteristics of College Students Using Data Mining on Big Data," 2023 International Conference on Artificial Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ICECONF57129.2023.10084276.
- [18].S. Bharathi, A. Balaji, D. Irene. J, C. Kalaivanan and R. Anusuya, "An Efficient Liver Disease Prediction based on Deep Convolutional Neural Network using Biopsy Images," 2022 3rd International Conference on Smart Electronics and Communication (ICOSEC), Trichy, India, 2022, pp. 1141-1147, doi: 10.1109/ICOSEC54921.2022.9951870.
- [19].I. Chandra, G. Sowmiya, G. Charulatha, S. D, S. Gomathi and R. Anusuya, "An efficient Intelligent Systems for Low-Power Consumption Zigbee-Based Wearable Device for Voice Data Transmission," 2023 International Conference on Artificial Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-7, doi: 10.1109/ICECONF57129.2023.10083856.
- [20].Revathi, S., et al. "Developing an Infant Monitoring System using IoT (INMOS)." International Scientific Journal of Contemporary Research in Engineering Science and Management 6.1 (2021): 111-115.

