# A Study of Employee Attrition and prediction using Logistic Regression in R programming

**Pramit Banerjee, Systems Analyst (IT MNC), Kolkata, India, pramiteejanai@gmail.com**

**Anchita Mishra, MBA Student (Full Time), IISWBM Kolkata, India, anchitamishra10@gmail.com**

**Abstract: Employee attrition is a major concern for the employer and most of the organizations are trying to finding ways to deal with the situation. It is very difficult to ascertain which set of employees will stay with the organization and who will leave the organization in the near future without understanding the pattern of employee behavior. This work aims to analyze the employee behavior based on the sample data and generate relationships among the data so as to build a model which may help any organization to predict the attrition of employees within the organization.**

*Keywords — Employee Attrition, Logistic Regression, Predictive Modelling, Machine Learning, R programming*
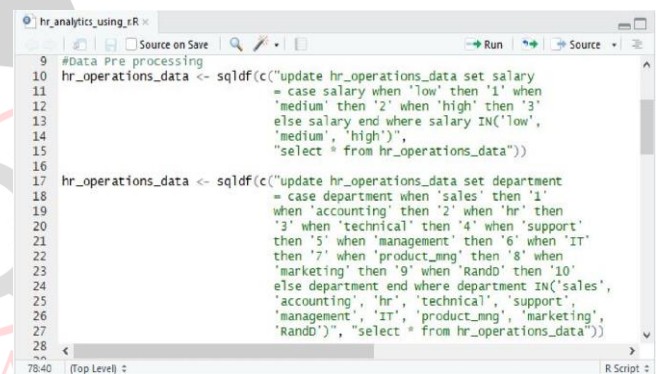
## I. INTRODUCTION

Employee attrition is a major concern for most of the organizations and it is tracked by organizations so as to maintain the optimal level of workforce for their day to day operations. If employees start leaving the organization, then without proper planning or a back up, it becomes difficult for the organization to maintain the daily operational activities, hampering the profitability in turn. In order to track the employee attrition, the different signals are being observed in the form of data which if properly analyzed, would lead to conclusions as to how many people and which set of people might exit the organization in the near future. One such attempt has been made through this research where historical data catering to the employees within an organization has been taken, such as the satisfaction level, last evaluation, project number, average monthly hours, time spent in company, work accident, promotion in last five years, department, salary and based on this the output has been predicted whether the employee is going to leave the organization or not. In this process, logistic regression model has been applied and the tool used is R programming. The data used here is open source data from Kaggle and github.

## II. DATA PREPARATION

The data consists of ten columns and out of these parameters, we consider the attribute employee left, as the output variable (dependent) which is a function of the other nine independent variables like satisfaction level, last evaluation, project number, average monthly hours, time spent in company, work accident, promotion in last five years, department, salary. At first, we have carried out the data preprocessing stage where the data has been cleaned and modified. For example, in the department column, there are ten distinct values and we have assigned categories from one to ten as numbers instead of text and also for the column salary, we have three distinct values and we have

assigned categories from one to three as numbers instead of text. The libraries used in R programming for operations here are sqldf, tidyverse, caret. The code snippet is attached herewith.



Figure - 1

In the above figure, the data set is first imported and then two update operations are performed using the sqldf package in R where the salary and department columns are revised by replacing the character values by numerical distinct values.

## III. METHODOLOGY AND IMPLEMENTATION

After the data pre processing stage, we have split the data into the test data set and the train data set according to the principles of machine learning. The data set contains around 15000 inputs and we have broken the dataset in such a way that we have trained around 80% of the dataset and based on the inputs and patterns of the train data set, we have applied it on the remaining 20% which is the test dataset. The code snippet is attached herewith.

Figure - 2

In the above figure, using the method createDataPartition using the caret library in R, the data partition has been done and two data frames train and test data have been created subsequently.

In order to ascertain better accuracy, we have done 10 times cross validation in this train and test data set before applying the model. This will in turn help the model to be more accurate. Finally, we have applied the model using logistic regression where employee left is taken as a function of the other parameters. The code snippet is attached herewith.



Figure - 3

In the above figure using the trainControl method in R, the ten-fold cross validation has been performed and afterwards, logistic regression is performed in train data set and the test data set is applied in the model.

## IV.  RESULTS AND CONCLUSION

The output of the model is found in the R console and it is given as below:



Figure - 4

From the above output, we see that there is strong negative correlation between employee left and satisfaction of the employee. This denotes that those employees whose satisfaction level is less, is more likely to leave the organization. Likewise, we see that those employees who have had a work-related accident, with less promotions in last five years are also more likely to leave the organization. From the above output, we see the p-value of the parameters. We find that for the parameters, satisfaction level, last evaluation, number project, average monthly hours, times pend in company, work accident, promotion in last five years, salary, the p value is less than 0.05 which means that they are statistically significant and hence these parameters play an important role in the employee attrition.

In order to identify the other metrics, we find out the accuracy, sensitivity and specificity of the model using the libraries ROCR and Metrics. The code snippet is given below:



Figure - 6

In the above figure, using the confusionMatrix function in R, the test data has been compared against the predicted data and thereby the true positive rate, i.e. the sensitivity, the true negative rate i.e. the specificity and the accuracy which is defined as ratio of sum of true positive and true negative to the sum of true positive, true negative, false positive and false negative data.

The other metrics are as follows:

Accuracy:  0.7902634

Sensitivity:  0.9243478

Specificity: 0.3490701

We find that the accuracy of the model is quite good and the sensitivity is also very high which denotes that the true positive rate is good and the correct prediction has been done on the test data.
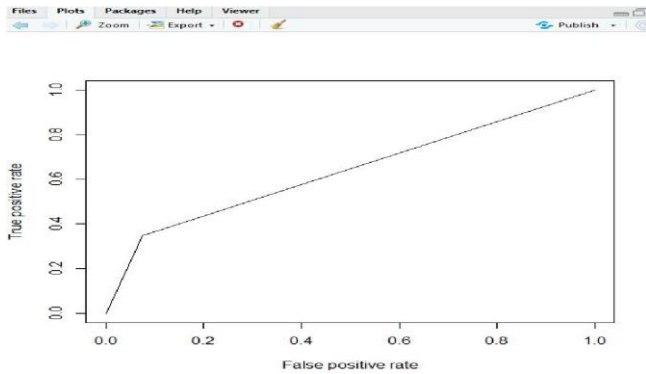


Figure - 7

In the above figure, we see the ROC (receiving operating characteristics) curve. Here, the X -axis is the false positive rate and the Y-axis is the true positive rate. The more the area under the curve, the better is the accuracy of the model.

We see that the area under the curve is: 0.636709, which indicates that the model is quite good.

## V.  FUTURE WORK

Through this model, it has been highlighted, how to predict the employee attrition of an organization using logistic regression. There is further scope of study using this model which can be carried out using higher volumes of data and across different data sets. We would also like to use different other algorithms in machine learning to optimize the results.

## VI.  REFERENCES

[1] Sung-Lin Chan, Xiang zhe Meng, S¨uha Kagan K¨ose, "EPFLMachineLearningCS-433-Project2Twitter Sentiment Analysis", 2018

[2] A logistic regression investigation of the relationship between the Learning Assistant model and failure rates in introductory STEM courses by Jessica L. Alzen, Laurie S. Langdon & Valerie K. Otero

[3] Logistic Regression Analysis of Graduate Student Retention by Sandra W. Pyke & Peter M. Sheridant, The Canadian Journal of Higher Education, volume XXIII-2, 1993

[4] Human Resource Analytics Kaggle Dataset by Ryan Nazareth and Hannes Draxl

[5] https://github.com/ryankarlos/Human-Resource-Analytics

[6] Twitter Sentiment Analysis using XG Boost and Random Forest Classification Algorithms: A Hybrid Approach by Ashwini M. Joshi, Sameer Prabhune, Nesara B R, IJREAM Vol – 05, 6th Sept 2019.

[7] Binary Logistic Regression Analysis in Assessment and Identifying Factors That Influence Students' Academic Achievement: The Case of College of Natural and Computational Science, Wolaita Sodo University, Ethiopia by Bereket Tessema Zewude (MSc) 1* Kidus Meskele Ashine (Ass. Professor)2 1. Wolaita Sodo University, College of Natural and Computational Sciences, Wolaita Sodo, Ethiopia, P.O.Box 138 2. Wolaita Sodo University, School of Law, Wolaita Sodo, Ethiopia. P.O. Box 138 - Journal of Education and Practice www.iiste.org ISSN 2222-1735 (Paper) ISSN 2222-288X (Online) Vol.7, No.25, 2016

[8] Application of Logistic Regression in the Study of Students' Performance Level (Case Study of Vlora University) by Miftar Ramosacaj Prof. Dr. Vjollca Hasani Prof. Dr. Alba Dumi - Journal of Educational and Social Research MCSER Publishing, Rome-Italy, Vol. 5 No.3 September 2015

[9] The importance of Logistic Regression implementation in the Turkish livestock sector and logistic regression implementations/fields by Murak Korkmaz, Selamy Guney, Sule Yuksel Yigiter – Journal Article, J.Agric. Fac. HR.U., 2012, 16(2): 25-36

[10] Noora Shrestha. Application of Binary Logistic Regression Model to Assess the Likelihood of Overweight. American Journal of Theoretical and Applied Statistics. Vol. 8, No. 1, 2019, pp. 18-25. doi: 10.11648/j.ajtas.20190801.13

[11] http://www.sthda.com/english/articles/36-classification-methods-essentials/151-logistic-regression-essentials-in-r/