# PREDICTING EMPLOYEE ATTRITION USING VARIOUS ML ALGORITHMS

## P. D. Sai Vardhan | A. Sai Kaushik | Balagopal. T. S | Dr. Sairabanu J | SCOPE

## Introduction

Employee attrition is the loss of employees due to personal reasons, low job satisfaction, low salary, and/or a bad business environment. Aim of our project is to develop a software that predicts employee attrition using various parameters utilizing various machine learning algorithms and also compare the results.

## Motivation

One of the major issues facing business leaders within companies is the loss of talented employees. This project studies employee job satisfaction using machine learning models.

## Scope of the Project

In the base paper the authors have used ADASYN as the oversampling technique and then fused it with other machine learning algorithms to predict whether the employee is going to leave the company or not. In our paper we are implementing SMOTE as the oversampling technique and then test it against other oversampling techniques. We use the database which is openly available from IBM's repository and then refine it to remove any anomalies that might skew our results later on. So in order to refine the database we use oversampling techniques as stated above and then make the refined database undergo prediction based machine learning algorithms to get results.
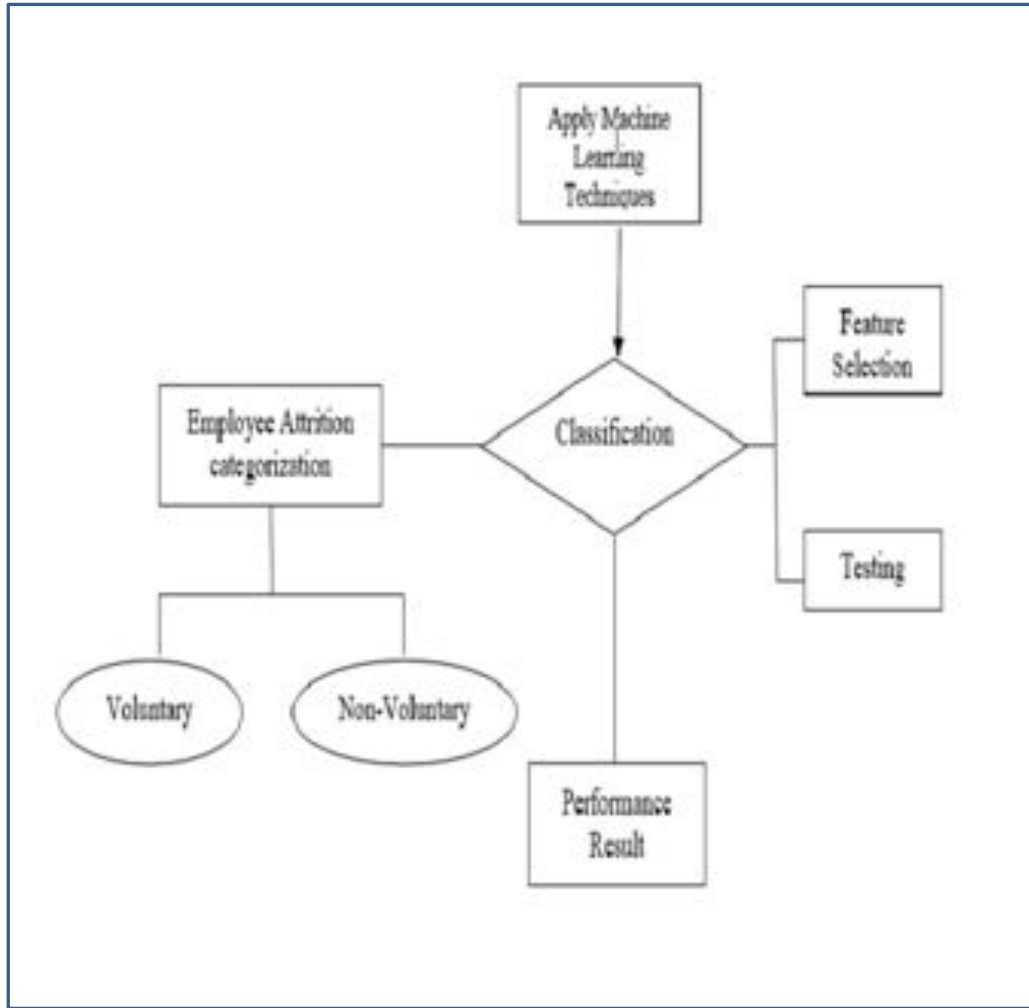
## Methodology

I. Related concepts

In this project the model is prepared in two steps. First the imbalanced data is converted to balanced data using SMOTE and Random Sampling. Then the data is trained using the Classifier models: KNN, Random Forest classifier, SVM and Artificial Neural Network(ANN).Then their performance is compared in terms of accuracy and precision. The existing paper uses ADASYN for sampling but we are using random sample and SMOTE instead. We also intend to compare the results thus obtained and infer which machine learning algorithm is best suitable for predicting the results.

II. Proposed System Methodology

Data Pre-processing:
The dataset is highly imbalanced and it is converted to balanced data by using over-sampling methods like ADASYN or SMOTE. After that the data is standardized or normalized to avoid over fitting also null values are replaced with zeros. Redundant and irrelevant columns are also removed in this step.



*Classification flowchart*

Data splitting:
Now 70 % of the data is used for training and 30% is used for testing, maintaining the class ratio. Training the model: As different models are to be used so k-fold cross validation is used for proper model selection. Then we fit our data to the models to evaluate the performance. Also we use different parameters for choosing the best model.

III. Requirement Analysis

Hardware Requirements
Processor – i3 / i5
RAM - 4GB / 8GB
System Free Space - Minimum 15GB
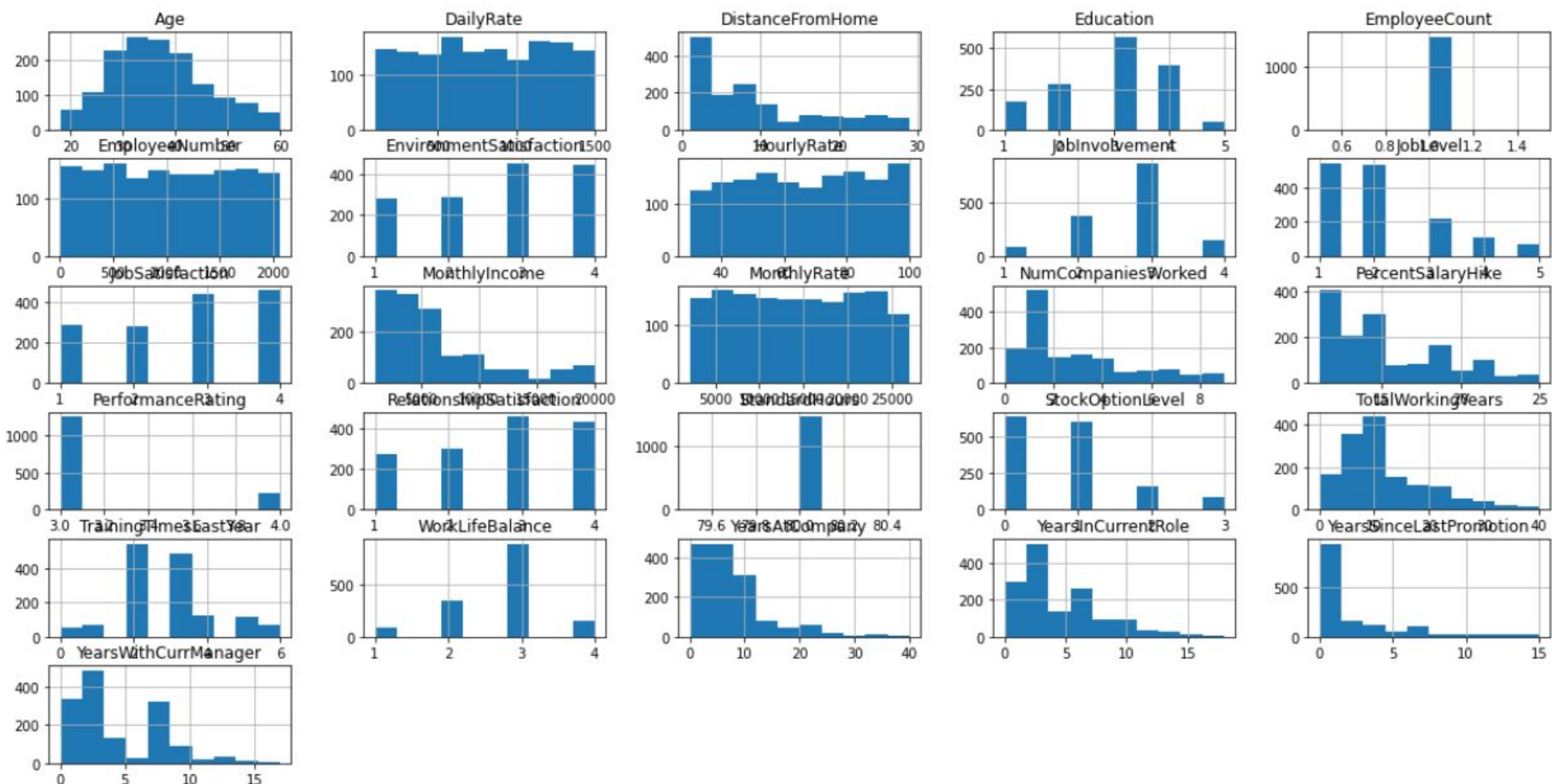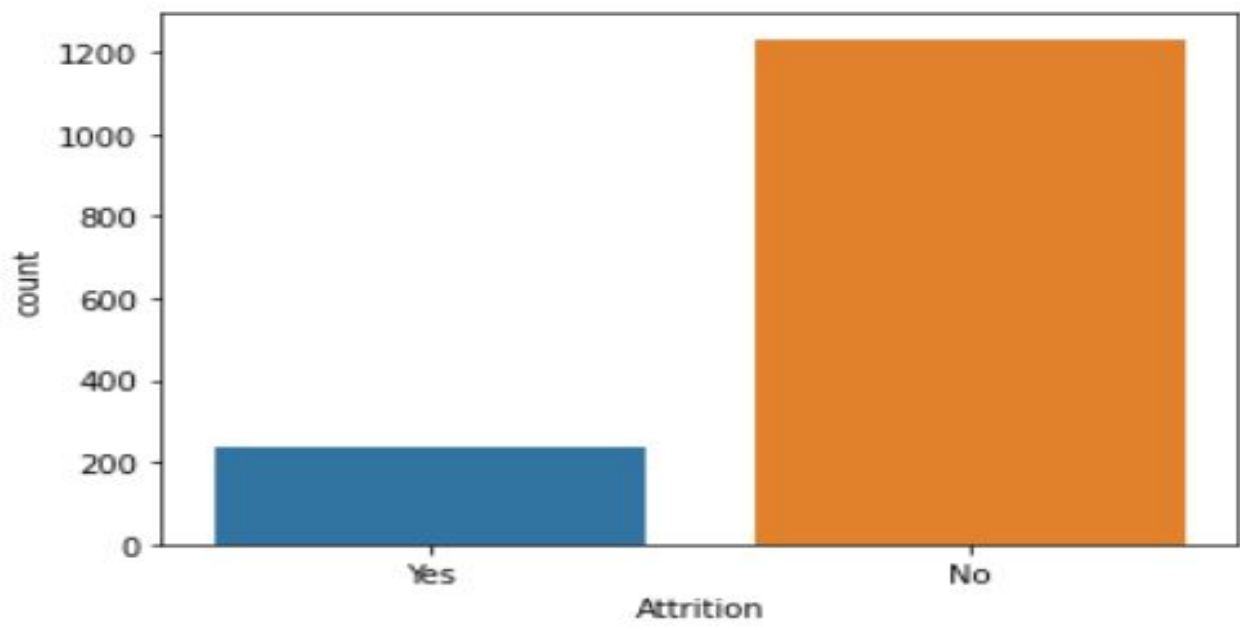Software Requirements
Programming Language - Python
IDE - Anaconda
User Interface - Visual Studio

## Results







Data processed using SMOTE tends to yield higher accuracy results after training with the above mentioned machine learning models when compared to random sampling. SMOTE gave accuracy of 95.2%, 77% and 75.8% for random forest, artificial neural network and k-NN, respectively.



|  | RF | ANN | KNN |
| --- | --- | --- | --- |
| Smote | 0.952381 | 0.770833 | 0.758929 |
| over sampling | 0.851852 | 0.833333 | 0.750000 |

## Conclusion

Based on the above results, we were able to conclude that of all the techniques we used, the combination of SMOTE and Random Forest has the highest accuracy. Thus we can say that SMOTE used with Random Forest gives a slightly better accuracy than ADASYN with Random Forest.

Compared with Random Oversampling, SMOTE was the better accurate oversampling technique even though it has less accuracy when combined with ANN. Random Forest algorithm provided better results compared with ANN and KNN.

## References

[1] S. Kaur and R. Vijay, "Job Satisfaction – A Major Factor Behind Attrition or Retention in Retail Industry", Imperial Journal of Interdisciplinary Research, vol. 2, no. 8, 2016.
[2] I. Paparrizos, B. B.Cambazoglu, A. Gionis, "Machine learned job recommendation", RecSys '11: Proceedings of the fifth ACM conference on Recommender systems, October 2011