

HR Analytics: Predictive analysis of employee promotion status

1. Introduction

HR analytics in companies plays a major role in restructuring the operanda of their HR department. A company of sizeable proportions deals with hundreds of employee records every day. Although HR analytics has been in operation in companies for years, some of these operations are still done manually. Automating such processes will aid in saving valuable time and increasing overall efficiency in the operation of the company.

Eligibility for promotions depends on many different criteria. These criteria vary between companies and even between departments within a company. Currently, there is no definitive method to determine why an employee receives a promotion since such decisions require logical reasoning and understanding the current environmental factors that computers are just not capable of. Problems such as this are where we utilize machine learning to our advantage.

Machine learning has been used to solve similar problems in different domains for several years now. In areas where some kind of human intervention is necessary, a well trained machine learning algorithm has been proven to be an acceptable substitute, if not an ideal solution. Here, we will be using machine learning techniques to not only predict the promotion status of future employees, but to also determine which of the provided attributes from the employees' data is most relevant to making this prediction. This paper proposes a solution involving machine learning techniques such as the Random Forest and XGBoost algorithms to learn from past employee records to aid this decision making process.

1.1 Random Forest

As useful as decision trees are, they fall short as problems become more and more complex. This is attributed to the fact that they are quite unstable, as even small changes in the data can affect the structure of the tree profoundly. This causes overfitting, which means that the tree cannot perform as well on real world data as it does on its training data.

The Random Forest algorithm aims to overcome that shortcoming. Rather than have a single tree making a definitive decision, it utilizes multiple trees, each making a prediction to arrive at its final conclusion. This algorithm uses the drawback of the decision tree to its advantage to create diverse trees by altering the data little by little. This helps keep the overfitting problem in check. Finally, the predictions from all the decision trees are polled and the prediction with the highest frequency is taken as the final prediction.

1.2 XGBoost

The XGBoost(short for Extreme Gradient Boosting) algorithm is essentially a decision tree which utilizes the concept of gradient boosting to enhance the performance of the decision tree. This algorithm performs best when being used to model small to medium sized structured data.

The XGBoost algorithm is generally preferred over other conventional algorithms as it combines many different traits from other algorithms, such as bagging, gradient boosting, random forest and decision trees and makes improves upon them through system optimization and algorithmic enhancements such as regularization, sparsity awareness, weighted quantile sketch, and so on.

1.3 Dataset

To fully explain the approach to solving this problem, we must first analyze and understand our dataset. Our dataset comprises of data collected by a company on previous candidates who were shortlisted for a promotion. Below, is a screenshot of the first 10 rows of our data. This dataset consists of 14 different attributes or columns, with 54808 total observations or rows.

employee_id	department	region	education	gender	recruitment_channel	no_of_trainings	age	previous_year_rating	length_of_service	KPIs_met>80%	awards_won?	avg_training_score	is_promoted
65438	Sales & Marketing	region_7	Master's& above	f	sourcing	1	35	5	8	1	0	49	0
65141	Operations	region_22	Bachelor's	m	other	1	30	5	4	0	0	60	0
7513	Sales & Marketing	region_19	Bachelor's	m	sourcing	1	34	3	7	0	0	50	0
2542	Sales & Marketing	region_23	Bachelor's	m	other	2	39	1	10	0	0	50	0
48945	Technology	region_26	Bachelor's	m	other	1	45	3	2	0	0	73	0
58896	Analytics	region_2	Bachelor's	m	sourcing	2	31	3	7	0	0	85	0
20379	Operations	region_20	Bachelor's	f	other	1	31	3	5	0	0	59	0
16290	Operations	region_34	Master's & above	m	sourcing	1	33	3	6	0	0	63	0
73202	Analytics	region_20	Bachelor's	m	other	1	28	4	5	0	0	83	0

Figure 1: Top 10 rows of the data

A breakdown of the attributes of this dataset is as follows:

employee_id (int): Unique id of the employee.

department (string): The department to which the employee belongs to. Possible values are: 'Analytics' .

'Finance' , 'HR' , 'Legal' , 'Operations', 'Procurement', 'R&D', 'Sales & Marketing', 'Technology'

region (string): Region of employment. Possible values are: 'region_10', 'region_11', 'region_12', 'region_13', 'region_14', 'region_15', 'region_16', 'region_17', 'region_18', 'region_19', 'region_2', 'region_20', 'region_21', 'region_22', 'region_23', 'region_24', 'region_25', 'region_26', 'region_27', 'region_28', 'region_29', 'region_3', 'region_30', 'region_31', 'region_32', 'region_33', 'region_34', 'region_4', 'region_5', 'region_6', 'region_7', 'region_8', 'region_9' .

education (string): Describes the level of education of the employee. Possible values are: 'Bachelor\'s', 'Below Secondary', 'Master\'s & above'.

gender (string): Gender of the employee. Possible values are: 'f', 'm'.

recruitment_channel (string): Channel of recruitment of the employee. Possible values are: 'referred', 'sourcing', 'other'.

no_of_trainings (int): Describes the number of training programs completed by the employee. Range: from 1 to 10.

age (int): Describes the age of the employee.

previous_year_rating (int): Employee rating from previous year. Range from 1 to 5.

length_of_service (int): The service length of the employee in years.

KPIs_met>80% (int): Describes whether the employee's Key Performance Indicators score are greater than 80%. Value is 1 if yes, else 0.

awards_won? (int): Whether the employee won any awards last year. Value is 1 if yes, else 0.

avg_training_score (int): Employee's average training score in training evaluations.

is_promoted (int): Whether the employee was promoted or not. Value is 1 if yes, else 0.

Here, the is promoted variable is our target variable. This is what we have to predict to produce our final result.

2. Method / Analysis:

2.1 Feature engineering

First off, we remove the employee_id column as it is just a column to distinguish records by and will not realistically impact the decision of whether an employee is to be promoted or not. This brings down our variable count to 13.

2.1.1 Duplicate values

Even though the size of our dataset is huge, we cannot be sure that all the observations in the data are unique. Looking at our data, we find that there are 118 duplicate records that are present in our dataset.

Even though it does not seem like a significantly large number, duplicate records will negatively impact the training process of machine learning algorithms and may cause them to overfit.

After the removal of the duplicate values, we have 54690 observations left.

2.1.2 Null values

Looking at our dataset, we see there are missing values. Observations with such values are not suitable for machine learning algorithms.

Luckily, out of 54690 observations, we have about 4042 observations with missing values. That accounts to about 7% of our total data.

A plot of this can be seen in figure 2.

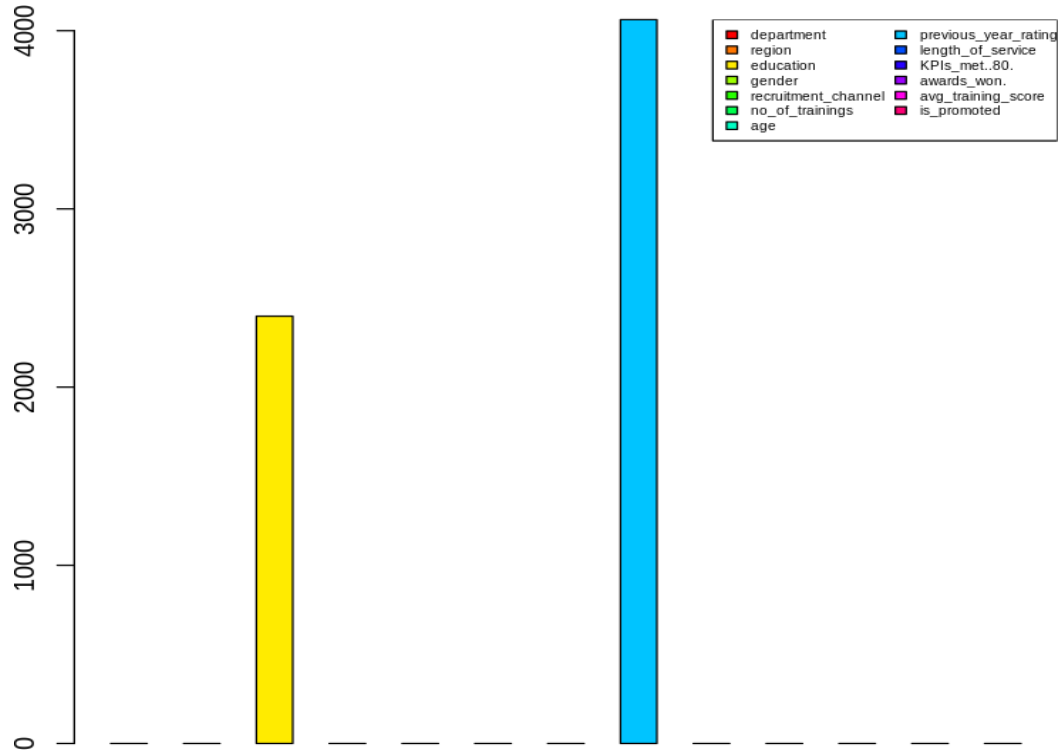


Figure 2: Barplot of columns along with the frequency of null values in them

This plot shows the frequency of null values in each column. There are only two columns with null values in our data, education and previous_year_rating, with 2398 and 4062 records with missing values respectively. These can be dealt with in one of two ways.

Either, these observations could be completely removed or they could be systematically generated using existing data. We choose the latter, as with that approach, we still maintain the size of our dataset.

Since both of these attributes are categorical, we can fill in missing values by finding the mode of the values that belong to that attribute. We also take into account the value of the target class for each attribute.

2.2 Pre-processing

The data is pre-processed to make it suitable for our algorithms. First, the target column values are renamed from “0” and “1” to “no” and “yes” so that the algorithm recognizes as categorical variables. Then, dummy variables are generated from the values of all the categorical attributes, or perform one-hot encoding for all the attributes.

2.3 Class imbalance problem

In this dataset, there is a clear imbalance between the values of our target variables. Out of 54690 observations, only 4665 observations are of class “no”, while 50025 observations are of class “yes”. This is a significant issue, as the difference exceeds more than 50% of the data.

For a machine learning algorithm to be properly able to parse and understand the given data, there should ideally be an equal distribution of the number of examples with the different classes it is meant to classify other data into. When one class takes precedence over the other class in the dataset, the algorithm is less likely to learn what the properties of each class are and tends to forget the less frequent class’s properties all together during training.

To ensure that this does not happen, we must make sure that there are equal number of examples for both the cases. Here, we simply randomly sample a subset of the data where the class is “yes”, as it is the class with higher frequency and append it to the observations with the class “no” to generate a new, minified training set with equal number of “yes” and “no” observations. This process is called undersampling.

This brings down the size of our dataset to 9330 observations total. Even though it is only a fraction of the original 54690, it is still a significant amount and should be enough to train our algorithms along with being balanced.

Figure 3(a): Distribution of class percentage before undersampling

Propn of people promoted/not promoted before undersampling

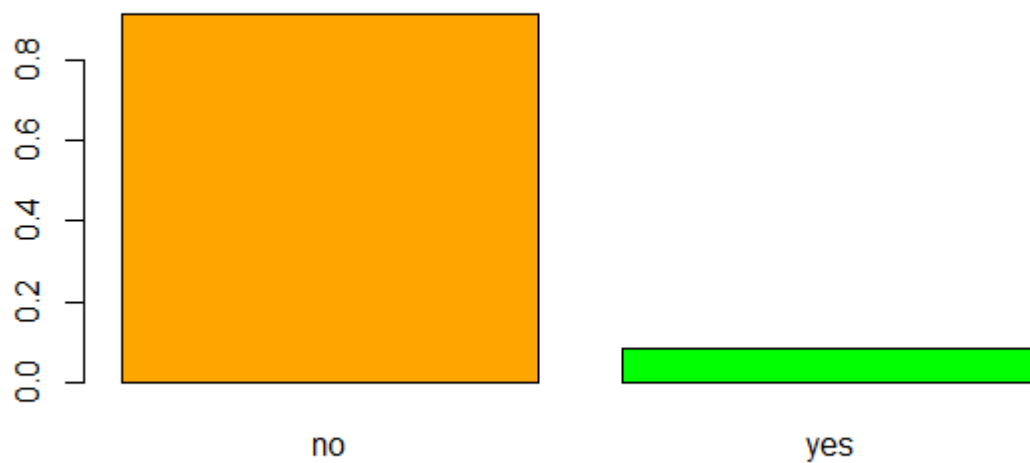
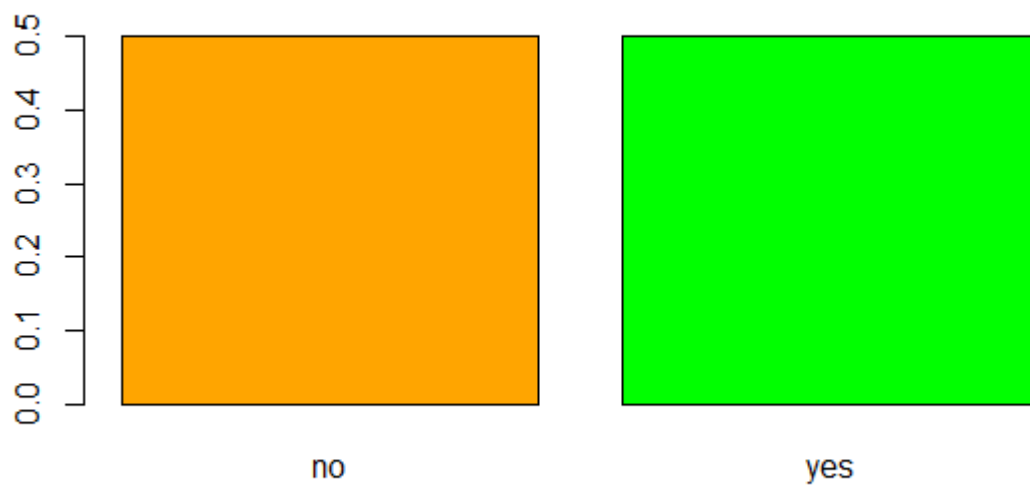


Figure 3(b): Distribution of class percentage after undersampling

Propn of people promoted/not promoted after undersampling



2.4 Training and Evaluation

The dataset is split into a training set and a testing set where the testing set contains 1/5ths of the records in the dataset and the rest belongs to the training set.

The chosen method of evaluation to measure the performance of our models is K-fold cross validation, where K is taken to be 5.

3.Results

The training of both the models resulted the following:

a. Random forest

The confusion matrix of the model evaluated on the training set, with 5 fold cross validation is as follows:

Cross-Validated (5 fold) Confusion Matrix

Prediction	Reference	
	no	yes
no	37.5	5.8
yes	12.5	44.2

Accuracy (average) : 0.8162

And the confusion matrix of the model evaluated on the test set:

Confusion Matrix and statistics

rfpreds no yes

no 719 122
yes 214 811

Accuracy : 0.8199
95% CI : (0.8017, 0.8371)
No Information Rate : 0.5
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.6399 McNemar's

Test P-Value : 6.889e-07

Sensitivity : 0.7706
Specificity : 0.8692 Pos
Pred Value : 0.8549 Neg
Pred Value : 0.7912
Prevalence : 0.5000
Detection Rate : 0.3853
Detection Prevalence : 0.4507
Balanced Accuracy : 0.8199

'Positive' Class : no

We can see that the performance of the model is consistent over both the training and testing set. So, we can ensure that our model has not overfit.

The random forest model also helps us calculate the importance of each attribute. The attribute importance scores, as predicted by our trained random forest model is as follows:

Attribute	Overall Score
avg_training_score	832.70
KPIs_met..80.	533.48
previous_year_rating	381.64
length_of_service	313.44
age	301.19
departmentSales & Marketing	186.40
awards_won.	114.01
departmentOperations	102.15
departmentProcurement	67.42
no_of_trainings	67.28
recruitment_channelSourcing	60.31
genderm	55.79
regionregion_2	41.94
departmentTechnology	39.70
departmentFinance	38.99
regionregion_22	38.76
regionregion_7	36.84
educationMaster's & above	31.27
departmentHR	30.92
educationBachelor's	30.88

Table 1: Top 20 most important features by Random Forest with no scaling of scores.

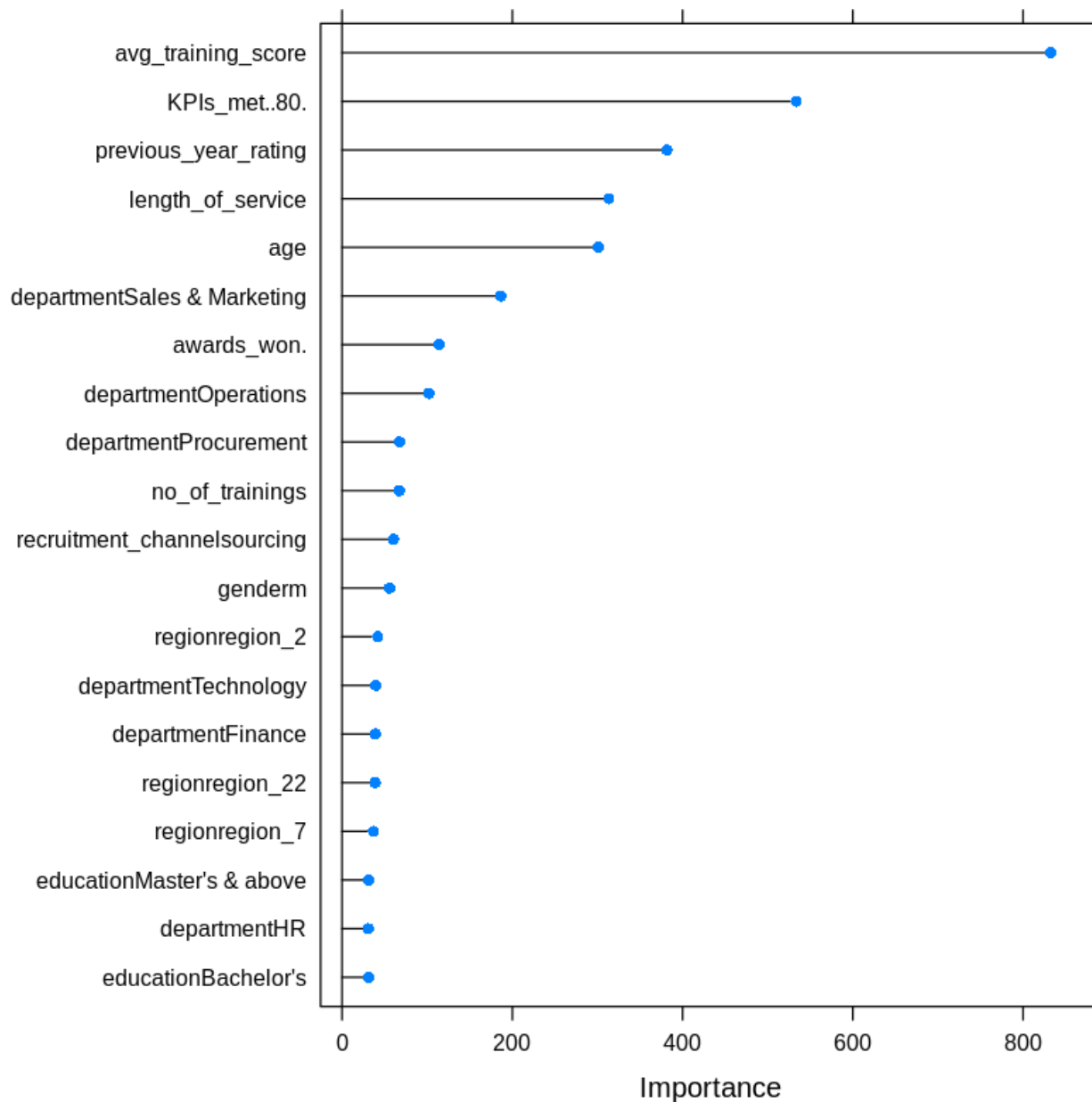


Figure 4: Plot of feature importance according to Random Forest

We can see here that the model has rated the average_training_score, KPI_mets>80 and the previous_year_rating as the 3 most important attributes for this decision.

b. XGBoost

Confusion matrix for the training set, evaluated with 5 fold cross validation is as follows:

Cross-validated (5 fold) Confusion Matrix

	Reference	
Prediction	no	yes
no	38.0	4.9
yes	12.0	45.1

Accuracy (average) : 0.8307

And the confusion matrix for the test set:

Confusion Matrix and Statistics

xgb_predict	no	yes
no	724	94
yes	209	839

Accuracy : 0.8376
95% C.I : (0.8201, 0.8541)
No Information Rate : 0.5
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.6752

Mcnemar's Test P-Value : 5.787e-11

Sensitivity : 0.7760
Specificity : 0.8992
Pos Pred Value : 0.8851
Neg Pred Value : 0.8006
Prevalence : 0.5000
Detection Rate : 0.3880
Detection Prevalence : 0.4384
Balanced Accuracy : 0.8376

'Positive' Class : no

The XGBoost model gives us a final validation accuracy of 0.8376(or 83.76%), which is marginally higher when compared to the random forest(81.99%). Other statistics such as the specificity and the sensitivity are also higher, which tells us that the XGBoost model will give us a better overall performance.

The feature importance table as calculated by the XGBoost tree is as follows:

xgbTree variable importance	
only 20 most important variables shown (out of 54)	
	Overall
avg_training_score	0.293120
KPIs_met..80.	0.218195
previous_year_rating	0.131789
departmentSales & Marketing	0.066747
length_of_service	0.057429
departmentOperations	0.041774
awards_won.	0.039652
departmentProcurement	0.028399
age	0.020241
departmentFinance	0.014365
departmentHR	0.012699
departmentTechnology	0.011097
no_of_trainings	0.008190
regionregion_22	0.005040
departmentR&D	0.003892
regionregion_4	0.003562
regionregion_7	0.002778
genderm	0.002639
regionregion_28	0.002507
educationBachelor's	0.002437

Table 2: Top 20 most important features by XGBoost model without scaling scores

20 most important variables according to XGBoost model

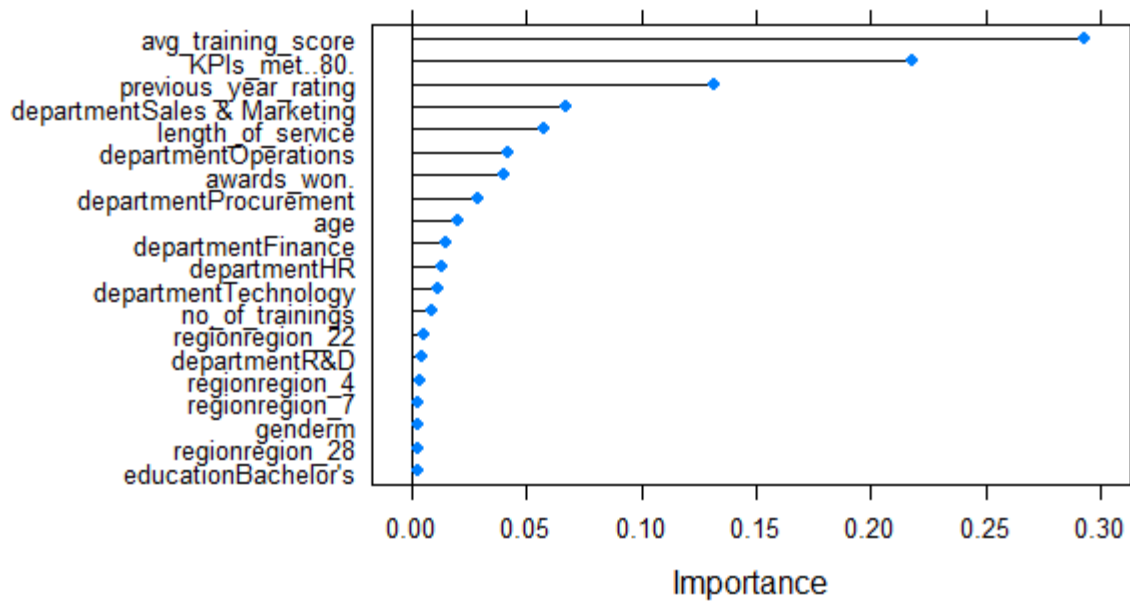


Figure 5: Feature importance plot by XGBoost

We notice that the list of the 10 most important features according to both the Random Forest and XGBoost models are nearly identical with the exception of `no_of_trainings` and `departmentHR`. This tells us that the other attributes can be discarded to get improved performance from the models.

4. Conclusion

We have downloaded the dataset, cleansed the dataset (removed duplicate values, balanced out classes using under-sampling, replaced null values with mode values), divided the data into testing and training datasets, trained the data using two machine learning algorithms (Random Forest and XGBoost) and tested the accuracy of the models using the testing data.

As we have demonstrated above, the use of machine learning as a predictive decision making tool or at least a suggestive tool is a completely viable solution for the presented problem. With fairly limited data, we were able to train algorithms to perform with significantly good accuracy. We can improve upon this with more data and more optimized solutions. Thus, we conclude that machine learning in the field of HR analytics by reducing the amount of time that goes into decision making, thereby increasing work efficiency.