



A REPORT ON HUMAN RESOURCE ANALYTICS

BY

Nital Vasvada

Rohini Bhandare

Mohit Joshi

Nitin Mali

Under the guidance of :

Dr. Jaideep Vaidya

FOR ACADEMIC YEAR 2017



RUTGERS
THE STATE UNIVERSITY
OF NEW JERSEY

ABSTRACT

When employees walk out the door, they take substantial value with them. Not surprisingly then, the broader application of predictive modeling across the enterprise along with the emergence of HR Analytics is leading organizations to ask how HR can start using data to predict and ultimately reduce employee turnover.

The specific goal here is to predict whether an employee will stay or voluntary leave within the next year. In the present data, this means predicting the variable “vol_leave” (0 = stay, 1 = leave) using the other columns of data. You can think of this data as historical data which tells us who did and who did not leave within the last year.

Our initial step is to describe and visualize our data. Then, we will develop two different kinds of predictive models. The first of these is a logistic regression model. Logistic regression models predict the likelihood of a categorical outcome, here staying or leaving.

The second kind of model is known as a decision tree (or a classification tree). A decision tree is essentially a set of rules for splitting the data into buckets to help us predict whether the employees in those buckets will end up in one group (staying) or another group (leaving).

In both cases, we are classifying into just two possible groups. This is known as “binary classification”.



TABLE OF CONTENTS

CHAPTER 1 : INTRODUCTION

CHAPTER 2 : DATASET DESCRIPTION AND RSTUDIO

CHAPTER 3 : DATA EXPLORATION AND
VISUALIZATION

CHAPTER 4: DATA MODELING

CHAPTER 5 : ROC CURVES

CHAPTER 6: DECISION TREE MODEL

CHAPTER 7 : CONCLUSION

INTRODUCTION



RUTGERS
THE STATE UNIVERSITY
OF NEW JERSEY

PROBLEM STATEMENT

Each and every organization needs to diminish costs, increment income, amplify operational proficiency, and concentrate on vital activities to remain productive. Whether in developed or developing markets, HR pioneers regularly battle to bolster the business with the gifted workforce it needs, because of budget and time requirements. One of the biggest challenges a company faces when it plans to launch a new line of services or products is recruiting the right people for the job in time for execution. So also, organizations bear unequivocal and implicit costs when talent exits the organization. It is more regrettable when representatives quit not long, but after in the wake of partaking in a costly preparing program supported by the association. Is there an approach to foresee such dangers and decrease the expenses connected with them?

SOLUTION

The attrition risk score of individual workers can be evaluated with predictive models of attrition. This can aid organizations to keep the potential attrition of high performing workers, factors contributing to attrition guarantee business continuation, and recognize loyal representatives. Employee's salary, designation, gender, performance, business area, age can be used for data analysis. Supervisors can recognize the key purposes for attrition and along these lines decrease its occurrence.



GOAL AND OVERVIEW

The particular objective here is to anticipate whether an employee will stay or voluntarily leave the following year. In the present information, this implies anticipating the variable "vol_leave" (0 = stay, 1 = leave) using the other columns of data. You can think about this information as historical data which lets us know who did and who did not leave the organization in the most recent year.

Our underlying stride is to depict and visualize our data. Then we will develop two kinds of predictive models. The first of these is a logistic regression model. Logistic regression models predict the likelihood of a categorical outcome, here staying or leaving.

The second sort of model is known as a decision tree. A decision tree is basically an arrangement of standards to part the data into bins to help us foresee whether the workers in those bins will end up in some group(staying) or another group (leaving). In both cases, we are characterizing into only two conceivable groups. This is known as "binary classification".

After we establish the framework for model development, we then disclose how to assess model quality utilizing overall model accuracy and the Receiver Operator Curve (ROC). The ROC lets us know what the best "limit" or "cutoff" while figuring out if somebody will leave

While there are absolutely more intricate techniques for foreseeing turnover, logistic regression and decision trees both work exceptionally well. Besides, they are similarly simple to execute and, above all, easier to interpret and explain. This is critical when translating modeling insights into action.



DATASET DESCRIPTION

AND R STUDIO



- DATA ANALYSIS TECHNOLOGIES USED

R is the main instrument for data analysis, statistics and machine learning. It is a programming dialect, so you can make your own particular objects and packages.

Like all projects, R programs unequivocally record the means of your analysis and make it simple to repeat as well as overhaul examination, which implies we can rapidly attempt numerous ideas or potentially revise issues. It is also platform independent and it's free, so you can utilize it at any business.

- DATASET AND R LIBRARIES

Our dataset is downloaded from Kaggle. This dataset has 8 unique attributes like those that you may have in your own HR information. These attributes and its description is mentioned below:

Role: This specifies the designation of the employee. We have 5 types of designation in our dataset namely CEO, VPs, Directors, Managers, and Individual Contributors.

Performance: The performance scale of an employee varies from 1-3. With 1 being lowest and 3 being highest.

Area: The area refers to the business department of an organization. We have areas such as Sales, Finance, Accounting, Marketing and others

Sex : Refers to the gender of the employee. It has two categories namely male female.

Id : Refers to the employee id

Age : Refers to the age of the employee

Salary : Refers to the salary of the employees

Vol_leave : This attribute is based on historical data. It depicts whether an employee has remained or voluntarily left an organization. '0' means to stay and '1' referring to leaving the organization.

Further more the R libraries used during the project implementation are given below:

- `library(plyr)` : The plyr package is a set of clean and consistent tools that implement the split-apply-combine pattern in R.
- `library(ggplot2)` : A system for 'declaratively' creating graphics, based on "The Grammar of Graphics".
- `library(caTools)` : Tools: moving window statistics, GIF, Base64, ROC AUC, etc.
- `library(RColorBrewer)` : Sequential, diverging and qualitative colour scales from colorbrewer.org
- `library(rpart.plot)` : Plot an rpart model, automatically tailoring the plot for the model's response type

DATA EXPLORATION AND VISUALIZATION



RUTGERS
THE STATE UNIVERSITY
OF NEW JERSEY

Let's start with the structure and summary command.

```
mydata = read.csv("humanresource.csv", header = TRUE)
str(mydata)
```

```
## 'data.frame': 11111 obs. of 8 variables:
## $ role : Factor w/ 5 levels "CEO","Director",...: 1 2 2 2 2 2 2 2 2 ...
## $ perf : int 3 3 1 2 3 1 2 3 2 1 ...
## $ area : Factor w/ 5 levels "Accounting","Finance",...: 5 3 2 5 3 4 1 2 5 3 .
## $ sex : Factor w/ 2 levels "Female","Male": 2 2 2 2 2 1 1 1 1 1 ...
## $ id : int 1 32 76 69 28 77 70 103 71 25 ...
## $ age : num 62 53.4 53.5 49.2 49.8 ...
## $ salary : num 1000000 258935 189828 207492 188205 ...
## $ vol_leave: int 0 0 1 0 0 0 0 0 1 0 ...
```

```
summary(mydata)
```

```
##      role      perf      area      sex
## CEO      : 1   Min.   :1.000   Accounting:1609   Female:6068
## Director: 100 1st Qu.:2.000   Finance   :1677   Male  :5043
## Ind      :10000 Median :2.000   Marketing :2258
## Manager  :1000 Mean    :2.198   Other     :2198
## VP       : 10  3rd Qu.:3.000   Sales     :3369
##          Max.    :3.000
##      id      age      salary      vol_leave
## Min.   : 1   Min.   :22.02   Min.    : 42168   Min.    :0.0000
## 1st Qu.:2778 1st Qu.:24.07   1st Qu. : 57081   1st Qu. :0.0000
## Median :5556 Median :25.70   Median  : 60798   Median :0.0000
## Mean   :5556 Mean   :27.79   Mean    : 65358   Mean   :0.3812
## 3rd Qu.:8334 3rd Qu.:28.49   3rd Qu. : 64945   3rd Qu. :1.0000
## Max.   :11111 Max.   :62.00   Max.    :1000000   Max.    :1.0000
```

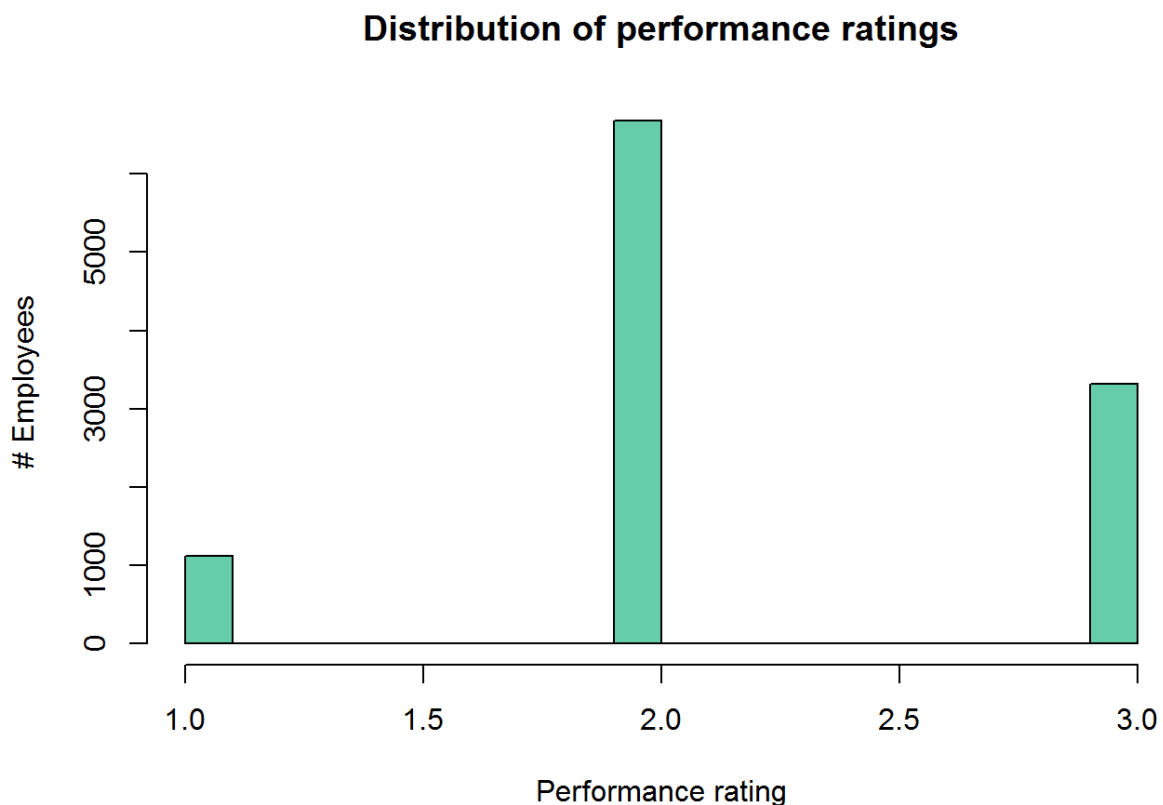
The summary information lets us know that we have 5 fundamental roles: CEO, VPs, Directors, Managers, and Individual Contributors. Since CEOs and VPs encounter an altogether different labormarket than Directors, Managers, and Individuals, hence incorporating them in our modeling effort doesn't bode well.

PERFORMANCE

We'll start with performance, here graded on a simple 1 (low) to 3 (high) scale.

```
table(mydata$perf)
##
##      1      2      3
## 1117 6675 3319

hist(mydata$perf,col = "aquamarine3", main = "Distribution of performance ratings",
      xlab = "Performance rating",ylab = "# Employees")
```

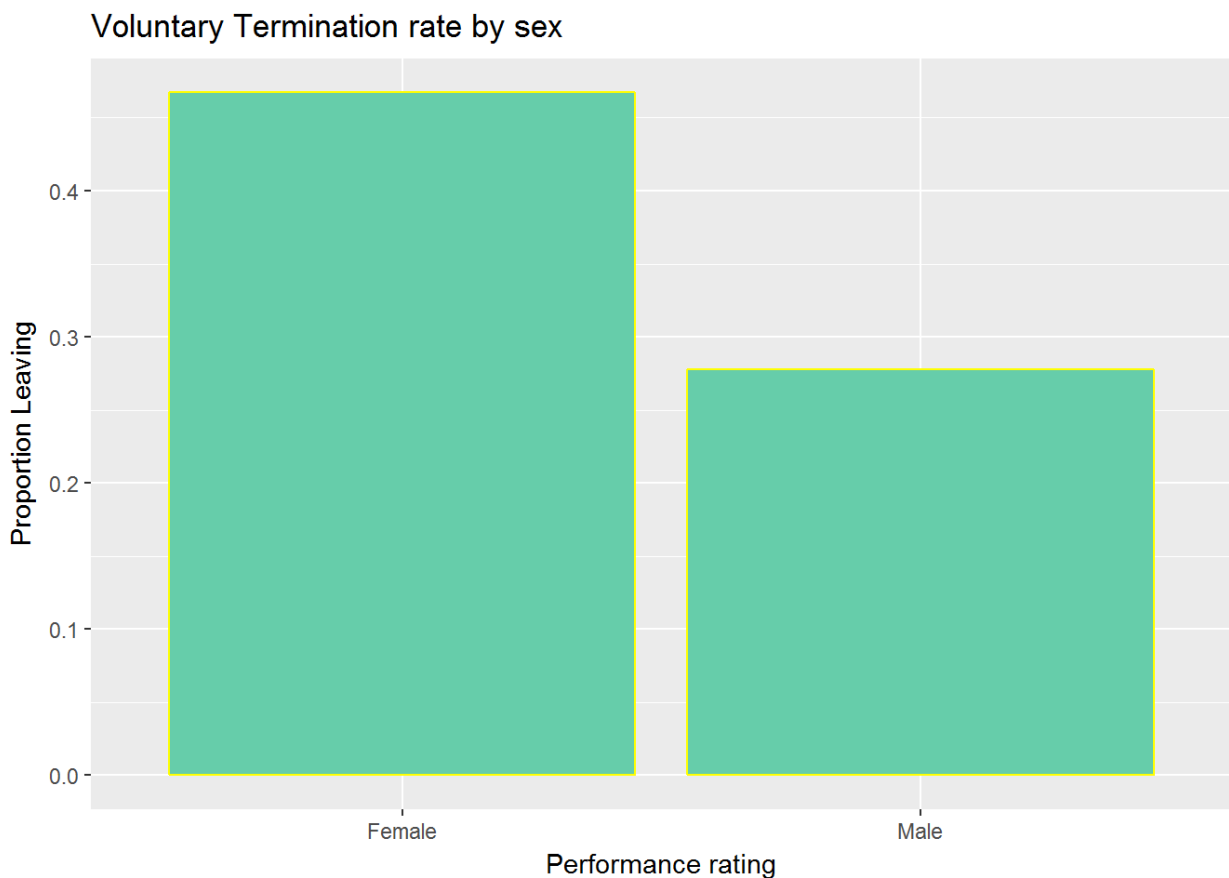


The table and histograms give us a glance at the distributions. With just three conceivable qualities, it doesn't bode well to discuss it, as being normal or skewed, however we can see that there are reasonable number of individuals in each group.

When we utilize the aggregate function to separate the probabilities by performance groupings, we can see turnover is much higher for the elite group. That is a major issue and certainly something we decided to keep for the model.

ANALYZING THE DISTRIBUTION OF MALES AND FEMALES.

```
performanceratingvsproportionleaving  
agg_sex=aggregate(vol_leave~sex,data = mydata,mean)  
print(agg_sex)  
##      sex vol_leave  
## 1 Female 0.4672050  
## 2  Male 0.2778108  
  
ggplot(agg_sex, aes(x=sex,y=vol_leave)) + geom_bar(stat = "identity",fill =  
"aquamarine3", colour = "yellow")+ ggtitle("Voluntary Termination rate by s  
ex")+  
labs(y="Proportion Leaving",x="Performance rating")
```



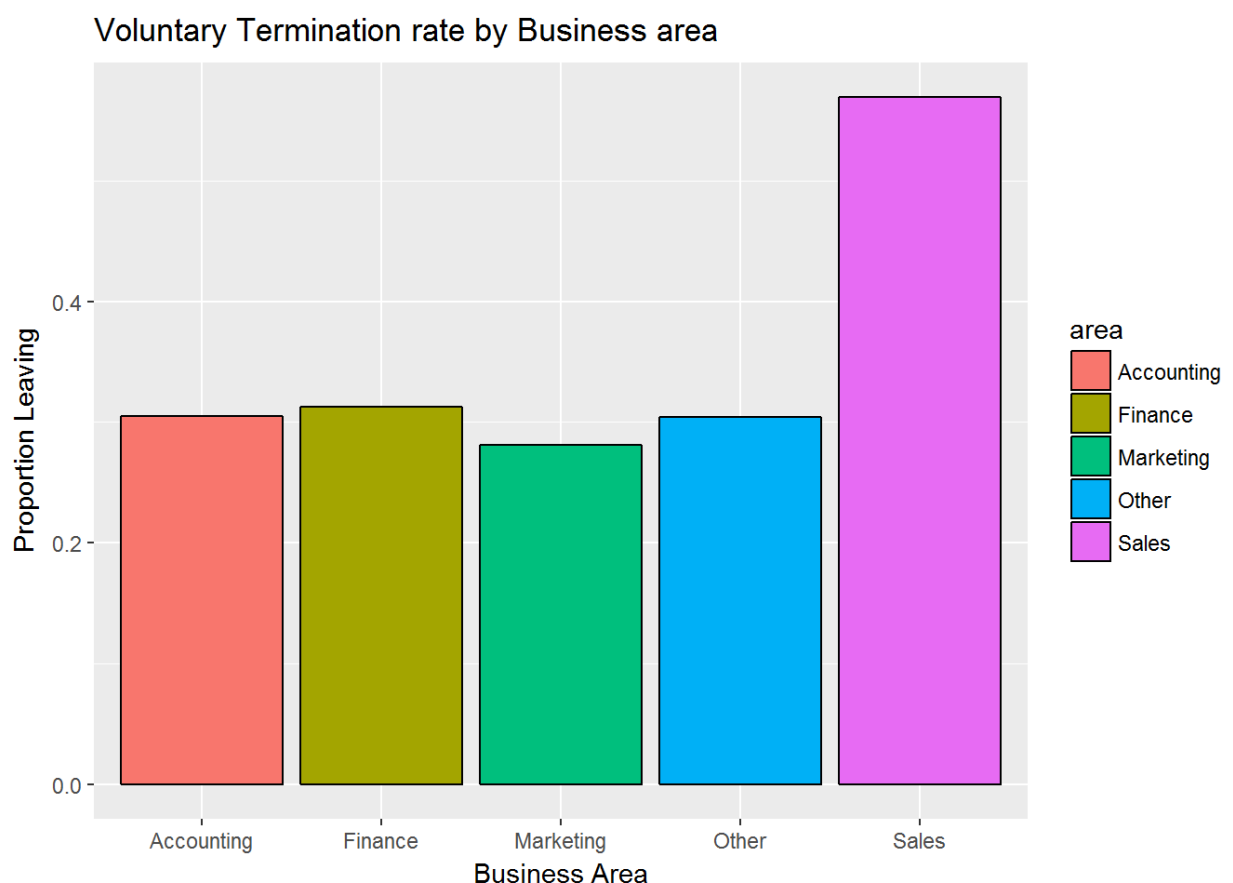
Unquestionably significant issues for female workers v/s. the male workers are observed. In such an event, alerts ought to go off and the investigative needs would begin with more distinct analyses and some extra focused modeling.

PLOTTING THE VOLUNTARY TERMINATION ON THE BASIS OF BUSINESS AREA

```
#businessareavsproportionleaving
agg_area=aggregate(vol_leave~area,data = mydata,mean)
print(agg_area)

##           area vol_leave
## 1 Accounting 0.3051585
## 2   Finance 0.3124627
## 3 Marketing 0.2812223
## 4    Other 0.3039126
## 5    Sales 0.5693084

ggplot(agg_area, aes(x=area,y=vol_leave,fill = area)) + geom_bar(stat = "id
entity", colour = "black")+ggtitle("Voluntary Termination rate by Business
area")+
  labs(y="Proportion Leaving",x="Business Area")
```



Those in Sales are much more likely to leave. Again, we'll definitely want to keep this for our models too.

SEGMENTING BUSINESS AREA BASED ON GENDER AND ANALYZING THE BUSINESS SECTORS AGAIN

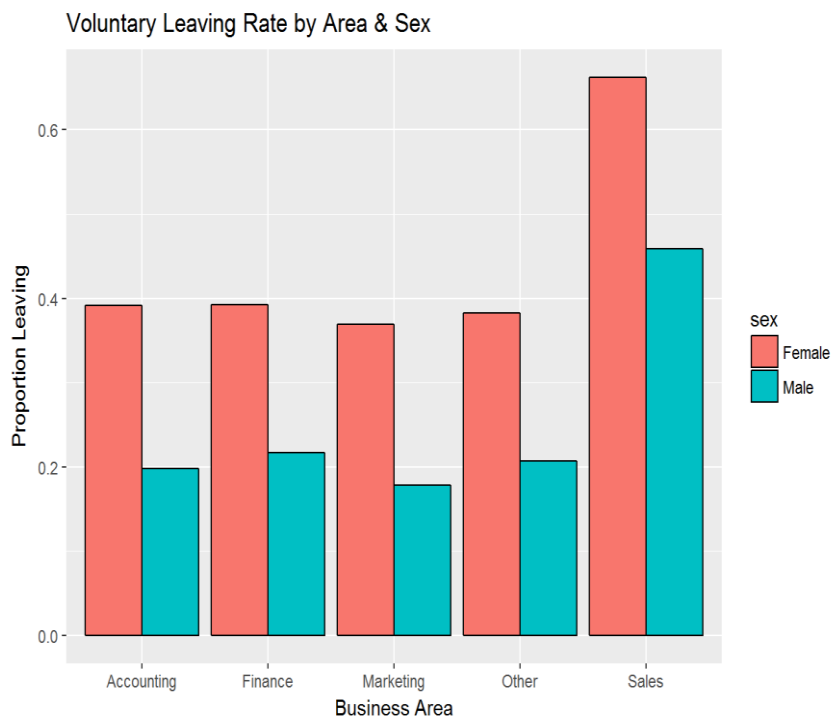


```
agg_as=aggregate(vol_leave~area+sex,data = mydata,mean)
print(agg_as)
```

```
##           area      sex vol_leave
## 1  Accounting Female 0.3918919
## 2    Finance Female 0.3923497
## 3  Marketing Female 0.3688525
## 4      Other Female 0.3828383
## 5     Sales Female 0.6623022
## 6  Accounting  Male 0.1983356
## 7    Finance  Male 0.2165354
## 8  Marketing  Male 0.1782274
## 9      Other  Male 0.2068966
## 10     Sales  Male 0.4583333
```

```
ggplot(agg_as,aes(x=area,sex,y=vol_leave))+geom_bar(aes(fill = sex), stat =
"identity", colour = "black", position = position_dodge()) +
  ggtitle("Voluntary Leaving Rate by Area & Sex") + labs(y = "Proportion Le
aving", x = "Business Area")
```

Plotting the Findings Below:



While observing the plot, we can see the difference is for the most part 20% across all areas. There might a marginally hoisted contrast for females in the business assemble, yet nothing gigantic.

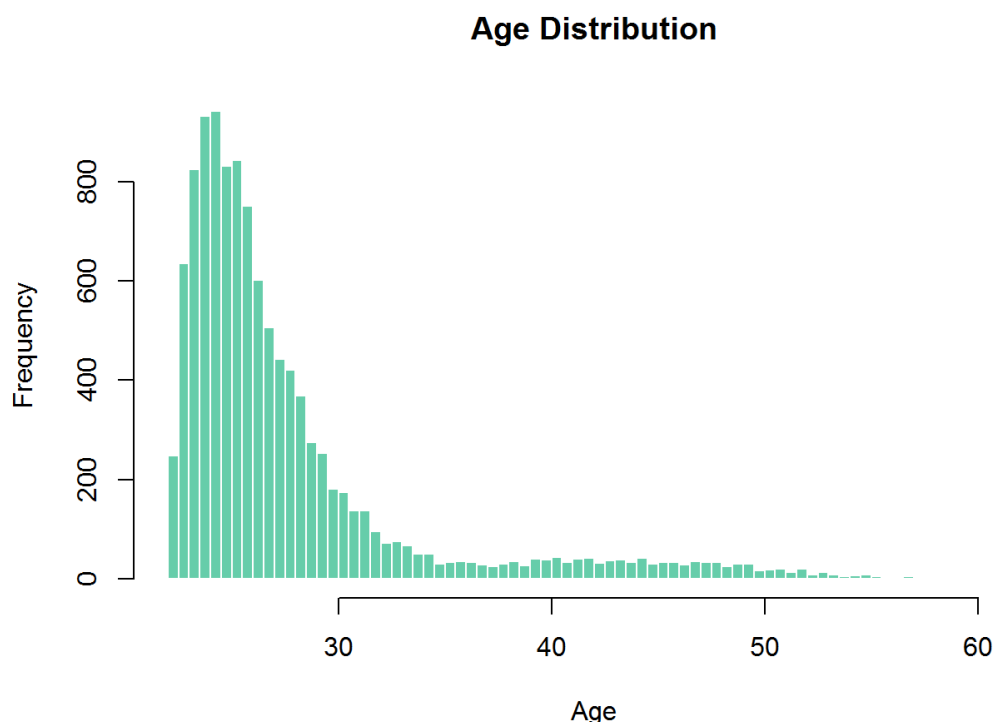
ANALYZING THE AGE OF EMPLOYEES

```
#analyzing the age of employees
```

```
hist(mydata$age,breaks=100,col="aquamarine3",main="Age Distribution",border  
=F,xlab="Age")
```

There is a solid skew here, with half of our workforce somewhere around 22 and 26 years old. We should recall however that we have three distinct levels: people, supervisors, and executives.

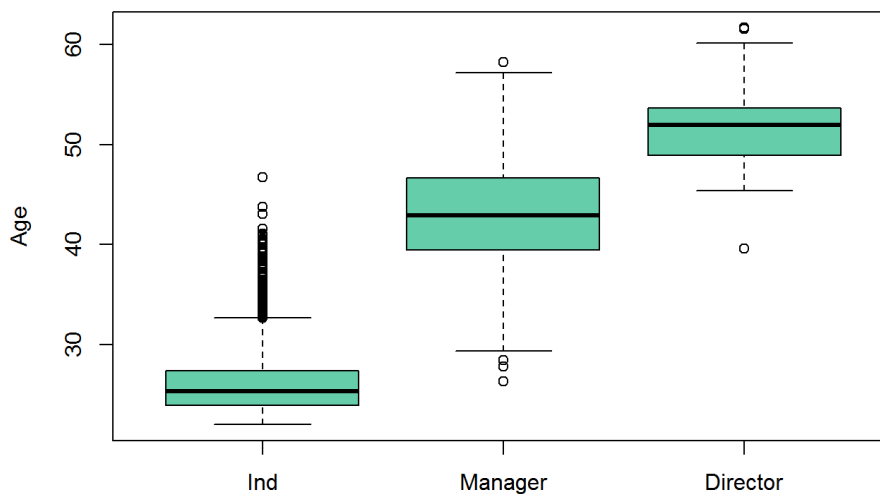
It will be more informative to see how those ages breakdown when we take that into account. Therefore box plots have been utilized for this purpose.



ANALYZING THE ROLES OF THE EMPLOYEES TO GET A BETTER IDEA OF THE AGE DISTRIBUTION

Our box plot indicates solid contrasts in age by the role with no meaningful overlap. This implies there is a solid relationship between Age and Role.

```
mydata$role<-factor(mydata$role,levels=c("Ind","Manager","Director"))  
boxplot(age~role,data=mydata,col="aquamarine3",ylab="Age")
```



If we just held one of these factors in the model, we might mistakenly attribute turnover differences to age when they are instead more related to role, or vice versa. Still, given that age is to a great degree skewed, we'll address that by making another variable that takes the log of age. This is known as a “data transformation” and is often done to make the distributions of a continuous variable more normally distributed.

This will give us a less skewed value and reduce any problems for inclusion in our later model.



ANALYZING VOLUNTARY TERMINATION BY AGE AND ROLE

```
#Logarithmicage
mydata$log_age = log(mydata$age)
summary(mydata$log_age)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max. 
##  3.092   3.181   3.246   3.305   3.349   4.127 

#ANALYZING VOLUNTARY TERMINATION BY AGE AND ROLE

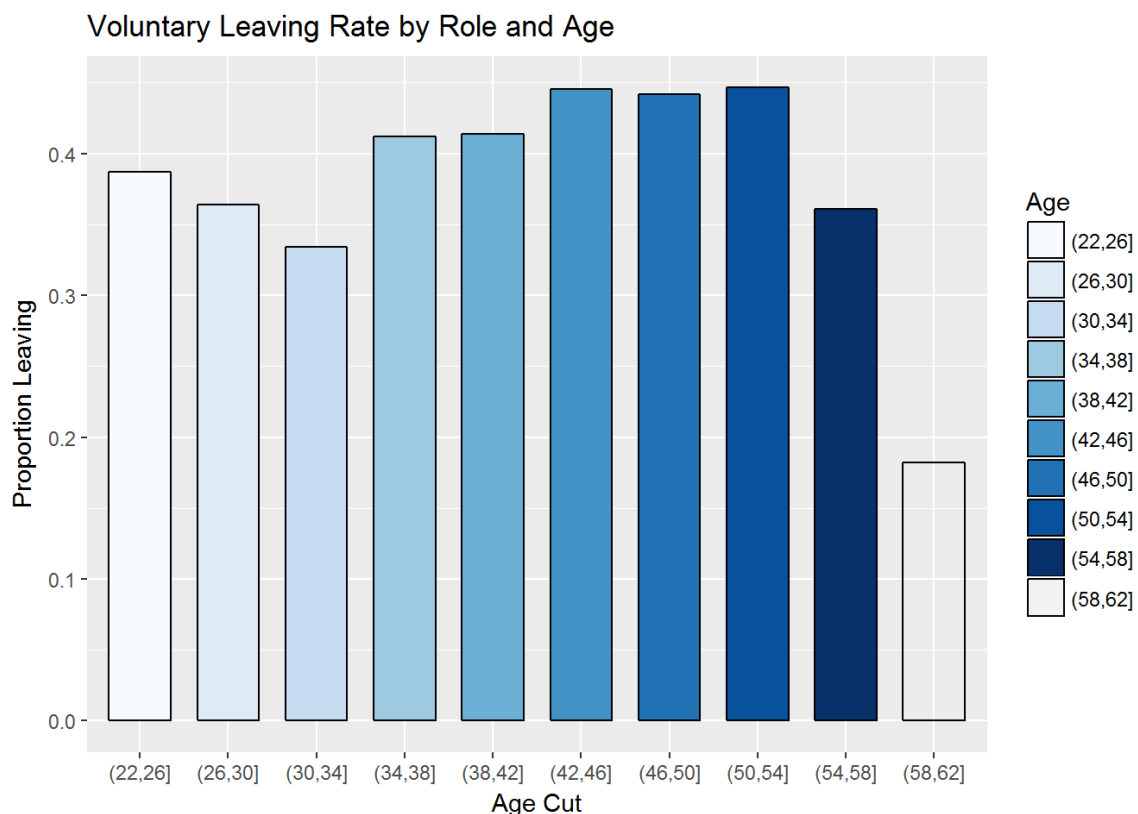
#before that segmenting age

age_agg = aggregate(x = mydata$vol_leave, by = list(cut(mydata$age, 10)), mean)

age_agg
names(age_agg) = c("Age", "Probability")

ggplot(age_agg, aes(x = Age, y = Probability, fill = Age)) + geom_bar(stat = "identity", width = .7, colour = 'black') +

  scale_fill_brewer() + ggtitle("Voluntary Leaving Rate by Role and Age") +
  labs(y = "Proportion Leaving", x = "Age Cut")
```



plot proposes a steep increment in turnover rate between age 34 and 46, for the most part hoisted rate until 58, and followed by a precarious drop-off.



ANALYZING THE SALARY PATTERN

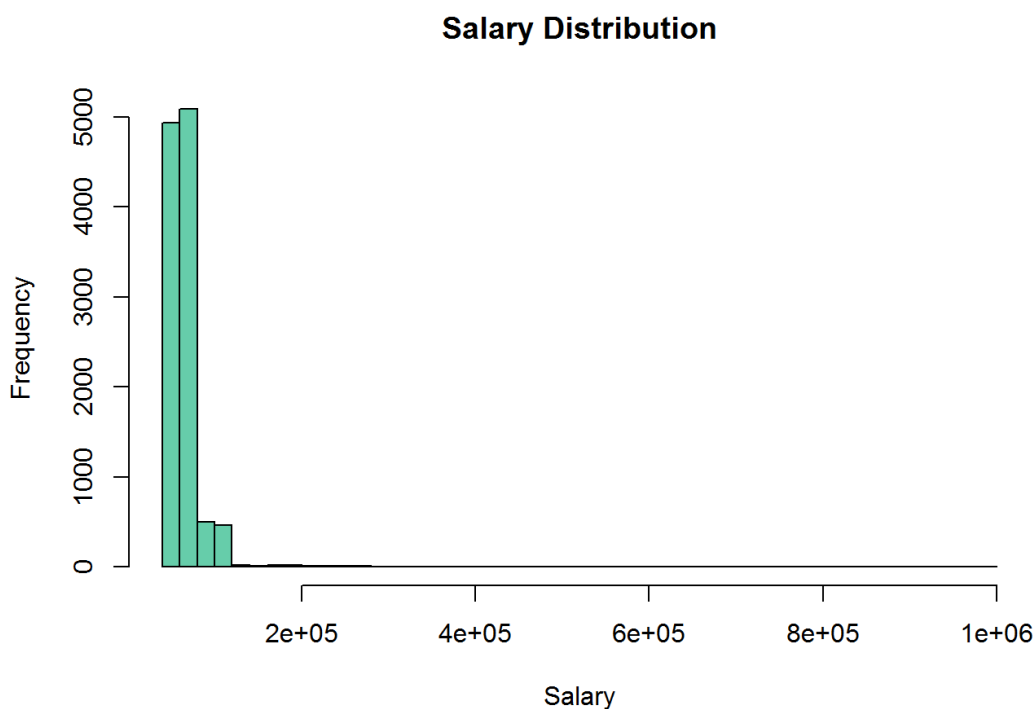
```
#analyzing the salary pattern
summary(mydata$salary)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  42170   57080   60800   65360   64950 1000000

hist(mydata$salary,breaks=50,col="aquamarine3",main="Salary Distribution",x
lab="Salary")
```

The median salary is \$60800 with the max being \$1000000 and the min being \$42170.

Plotting the Histogram for the same



FINDING THE SALARY DISTRIBUTION BASED ON QUANTILE

```
#FINDING THE SALARY DISTRIBUTION BASED ON QUANTILE
quantile(mydata$salary,probs=seq(0,1,.2))

##      0%      20%      40%      60%      80%     100%
## 42168.22 56191.54 59386.81 62317.40 66173.65 1000000.00
```

Salary variable is highly skewed with almost 80% of the people earning till \$66173.65

PLOTTING THE BOXPLOT TO SHOW SALARY DISTRIBUTION BASED ON ROLES



```
#PLOT THE BOXPLOT TO SHOW SALARY DISTRIBUTION BASED ON ROLES  
boxplot(salary~role,data=mydata,col="aquamarine3",main="Salary")
```



From the above graph we can observe three distinct pieces, relating firmly to the role. Given the tremendous contrasts in pay, we decided to find some way to represent salary that is more meaningful across the roles.

DATA MODELING



DATA SPLITTING

Before we start creating models, we need to split our data into a training set and a test set.

We will utilize two-thirds of the data for training and model development and one third of the data for testing the models.

The splitting strategy here is random with the limitation that both sets will have a similar extent of leavers and stayers. We set the random seed to a particular number so we can simply replicate our outcomes.

```
Data Modelling
set.seed(42)
spl<-sample.split(mydata$vol_leave,2/3)
train<-mydata[spl,]
test<-mydata[!spl,]

#1.Logistic regression

test_mean = mean(test$vol_leave)
train_mean = mean(train$vol_leave)
print(c(test_mean, train_mean))

## [1] 0.3812095 0.3812610
```

OUR MODELS: A BRIEF INTRODUCTION

Logistic regression builds a condition that as a result predicts the probability of a two-class result (staying or leaving) utilizing the chosen indicators. Each of the indicators are connected with a "significance" pointer that lets you know whether the indicator is helpful or not.

By contrast, decision trees work by utilizing the indicators to part the data into buckets using a set of decision rules.

In our present data for instance, the general likelihood of leaving is 38% and we are attempting to make distinctive buckets in which a few containers of workers are a great deal more inclined to leave (say, 70%) and other buckets have employees are much more likely to stay (say, 80%).

- **THE LOGISTIC REGRESSION MODEL**

With our dataset now split, we'll utilize the `glm` function to make our logistic regression model. The plan beneath demonstrates that we are foreseeing the `vol_leave` variable with the chose set of indicators.

In addition to including each of the variables we discussed in the descriptive analyses above, we are also including the “sex*area” interaction term.

This gives a test to whether the effect of the male/female variable on turnover varies as indicated by the business region (or equally, whether the effect of business region on turnover contrasts for males v/s females).

The “family” argument of the function is set to “binomial”, indicating to the model that we have a 0/1 response outcome.

INTERPRETING THE MODEL OUTPUT

fit the model

```
modell1<-glm(vol_leave~perf+role+log_age+sex+area+salary,data=train,family =  
"binomial")
```

```
summary(modell1)
```

```
## Call:  
## glm(formula = vol_leave ~ perf + role + log_age + sex + area +  
##      salary, family = "binomial", data = train)  
## Deviance Residuals:  
##      Min        1Q      Median        3Q        Max   
## -2.4613   -0.9193   -0.6081    1.1021    3.2016   
## Coefficients:  
##              Estimate Std. Error z value Pr(>|z|)      
## (Intercept)  5.772e+00  8.768e-01   6.584 4.59e-11 ***  
## perf         4.924e-01  4.383e-02  11.234 < 2e-16 ***  
## roleManager  3.491e+00  2.423e-01  14.409 < 2e-16 ***  
## roleDirector 8.925e+00  7.113e-01  12.547 < 2e-16 ***  
## log_age      -1.033e+00  2.512e-01  -4.111 3.94e-05 ***  
## sexMale      -9.476e-01  5.353e-02 -17.700 < 2e-16 ***  
## areaFinance  7.239e-02  9.642e-02   0.751   0.453   
## areaMarketing -2.652e-02  9.183e-02  -0.289   0.773   
## areaOther     9.349e-02  9.114e-02   1.026   0.305   
## areaSales     1.232e+00  8.345e-02  14.763 < 2e-16 ***  
## salary       -6.736e-05  4.643e-06 -14.509 < 2e-16 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## (Dispersion parameter for binomial family taken to be 1)  
##  
##      Null deviance: 9836.9  on 7397  degrees of freedom  
## Residual deviance: 8707.9  on 7387  degrees of freedom  
##      (9 observations deleted due to missingness)  
## AIC: 8729.9  
##  
## Number of Fisher Scoring iterations: 4
```

Not surprisingly, we see that those rated more highly on performance are significantly more likely to leave, consistent with our descriptive analyses.

The Role Variable

In the case of Managers, we see that the estimate is positive, meaning that Managers are significantly more likely to leave. Surprisingly though, the designation of Director did not significantly predict turnover.

Given the marked drop-off in turnover we saw for Directors in our descriptive analyses, how can this be?

The short answer is that when we include other variables such as age and performance, the impact of Director per se disappears.

The Age Variable

Our transformed age variable is significant and negative. This tells us that older employees are less likely to leave.

The Male/ Female Variable

We see a significant ($p < .05$) and negative parameter value males. This means males are less likely to leave than females.

The Area Variable

The interpretation for Area is a bit more complex. The bottom line is that those in the Sales are significantly more likely to leave (positive and statistically significant parameter value) while there are no significant differences associated with the other areas.

This is consistent with our descriptive analysis above. Again, to keep this primer focused, we will just leave the Area values in as they are. In principle, though, one might wish to simply code area as just two categories (Sales v. Non-Sales) to reduce the model complexity.

Some more Insights From the Above Summary:

1. First of all, we can see that areaFinance, areaMarketing and areaOther is not statistically significant.
2. As for the statistically significant variables, salary, areaSales and perf has the lowest p-value suggesting a strong association of these variable with the probability of leaving the company.
3. Now we can run the anova() function on the model to analyze the table of deviance.

Chi Square Test

```
4. anova(model1, test="Chisq")
5. ## Analysis of Deviance Table
6. ##
7. ## Model: binomial, link: logit
8. ##
9. ## Response: vol_leave
10. ##
11. ## Terms added sequentially (first to last)
12. ##
13. ##
14. ##           Df Deviance Resid. Df Resid. Dev    Pr(>Chi)
15. ## NULL                7397      9836.9
16. ## perf             1      111.41      7396      9725.5 < 2.2e-16 ***
17. ## role             2       18.80      7394      9706.7 8.277e-05 ***
18. ## log_age          1       11.14      7393      9695.6 0.0008449 ***
19. ## sex              1      290.47      7392      9405.1 < 2.2e-16 ***
```

```

20. ## area      4    458.46      7388      8946.7 < 2.2e-16 ***
21. ## salary    1    238.77      7387      8707.9 < 2.2e-16 ***
22. ## ---
23. ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Analysis:

1. The difference between the null deviance and the residual deviance shows how our model is doing against the null model (a model with only the intercept). The wider this gap, the better
2. A smaller p-value here indicates that all the variables in the model are significant

Assessing the predictive ability of the model

```

####Assessing the predictive ability of the model
fitted.results = predict(modell1, test, type = 'response')
fitted.results = ifelse(fitted.results > 0.5,1,0)

# Confusion Matrix
confMat=table(actual = test$vol_leave, prediction = fitted.results)

accuracy = sum(diag(confMat))/sum(confMat)
accuracy [1] 0.691248

```

The Accuracy of the Model is 0.691248



RUTGERS
THE STATE UNIVERSITY
OF NEW JERSEY

ROC CURVES



RUTGERS
THE STATE UNIVERSITY
OF NEW JERSEY

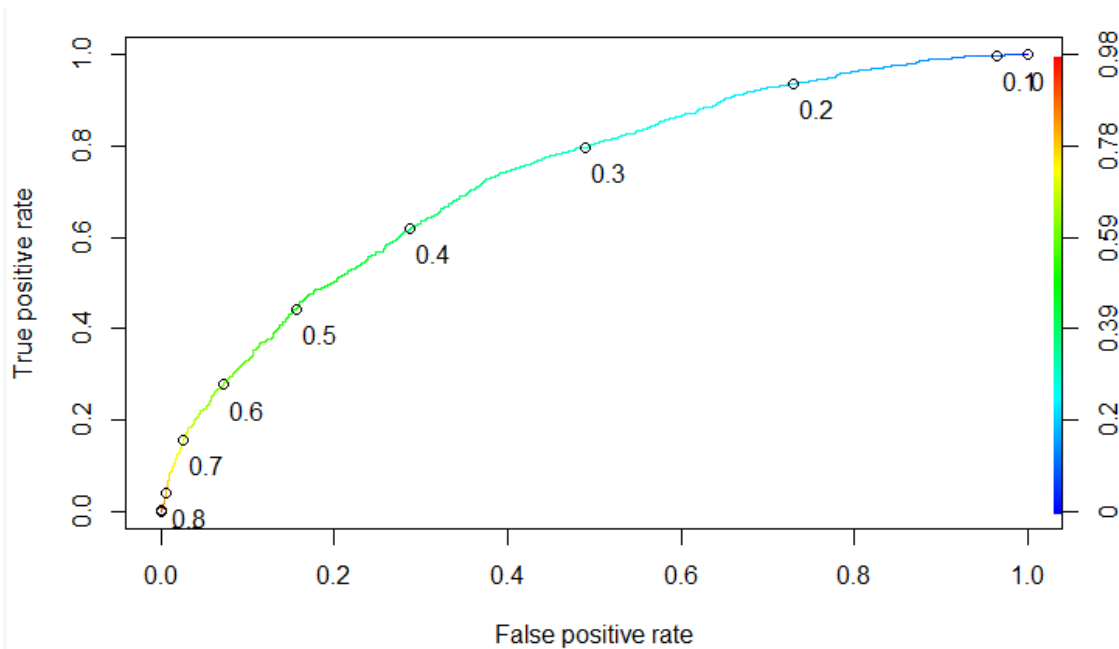
ROC Curve & AUC

As a last step, we are going to plot the ROC curve and calculate the AUC (area under the curve) which are typical performance measurements for a binary classifier.

```
####Roc curve

p <- predict(modell,test, type = "response")
pr = prediction(p, test$vol_leave)
prf = performance(pr, measure = "tpr", x.measure = "fpr")
plot(prf)

Auc Curve
auc = performance(pr, measure = "auc")
auc = auc@y.values[[1]]
auc
```



```
auc = performance(pr, measure = "auc")
auc = auc@y.values[[1]]
auc
## [1] 0.7326298
```

Analysis:

1. The ROC is a curve generated by plotting the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings.
2. The AUC is the area under the ROC curve.
3. As a rule of thumb, a model with good predictive ability should have an AUC closer to 1 (1 is ideal) than to 0.5.
4. Based on the value of AUC for our dataset, we can say that it has good predictive ability.

The ideal model has an AUC of 1 and a non-predictive model embracing the askew line would have an AUC of .5 (since it slices our space down the middle). Our logistic regression model falls into the fair (read: “Not that bad”) category.

THE DECISION TREE MODEL



RUTGERS
THE STATE UNIVERSITY
OF NEW JERSEY

Decision Trees (DT) make everything easier by telling us which predictors are useful. More technically, decision trees are robust to issues like predictors with skewed distributions because the underlying math does not depend on estimating parameter values and having nice, normally distributed predictors. we will create our model using the training data and then use the resulting model to make predictions with our test data.

```
####Decision tree
set.seed(42)
fit<-rpart(vol_leave~role+age+sex+area+perf,data=train,method = "class")
fit

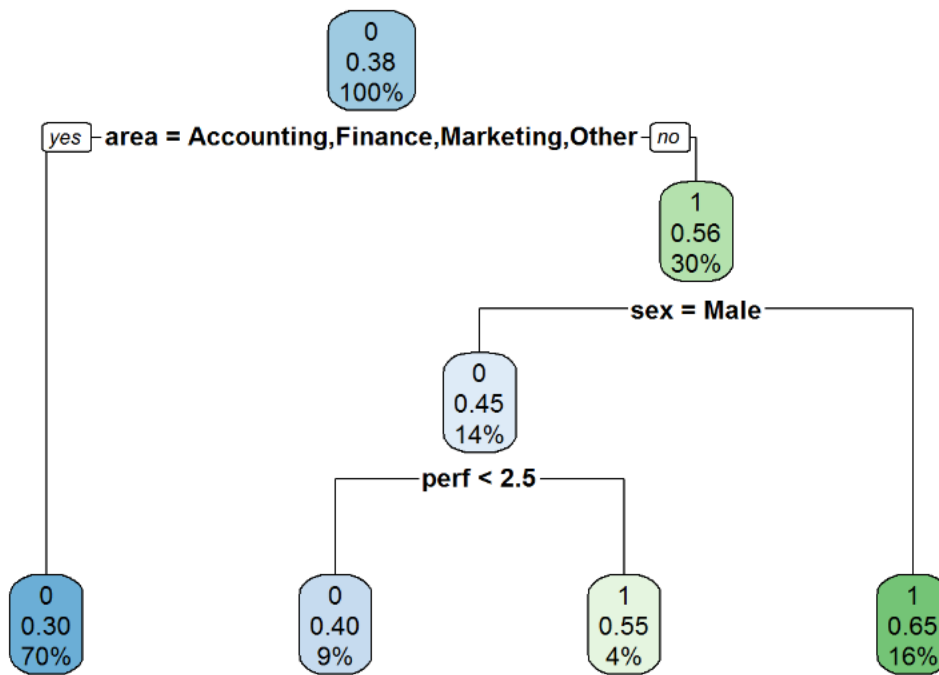
## n= 7407
##
## node), split, n, loss, yval, (yprob)
##      * denotes terminal node
##
##  1) root 7407 2824 0 (0.6187390 0.3812610)
##    2) area=Accounting,Finance,Marketing,Other 5185 1581 0 (0.6950820 0.3049180)
##      *
##    3) area=Sales 2222  979 1 (0.4405941 0.5594059)
##      6) sex=Male 1006  450 0 (0.5526839 0.4473161)
##        12) perf< 2.5 697  279 0 (0.5997131 0.4002869) *
##        13) perf>=2.5 309  138 1 (0.4466019 0.5533981) *
##        7) sex=Female 1216  423 1 (0.3478618 0.6521382) *
```

Plot The Tree

```
par(mar = c(5,4,1,2))
rpart.plot(fit, sub = NULL, main = "Basic Decision Tree")
```



Basic Decision Tree



- INTERPRETING THE MODEL OUTPUT

Analysis:

1. The first node is alluded to as the root. The '0' alludes to the dominate case. Here, 62% of those in our training data have 0 (Stay) for the response variable and 38% have a 1 (Leave).
2. Below that, we see our first decision node. In the event that our workers are in the Accounting, Finance, Marketing, or Other regions, then we say 'yes' and take the left branch. On the off chance that the answer is 'no' (i.e. they are in Sales), then we take the right branch.
3. After the left branch, we see that it ends into a solitary node. Think of this node like a bucket for all of those who are not in Sales. For all of these people, the most common response is '0' (Stay), with 70% employee who will stay in the company and only 30% in this bucket will leave the company. The '70%' reported in the bottom of the node tells us that this single bucket accounts for 70% of the total sample we are modeling.
4. On following the right branch, we see that the most well-known reaction is '1' for the employee who will leave the company.



Moreover, the node is likewise letting us know 42% of employees in this bucket will stay while 58% will leave.

5. Proceeding with the right branch is further, if the worker is male, we say 'yes' and go to the left side. On the off chance that the worker is female, we go right.
6. For females, we wind up in a terminating node that has a dominant response of 1 (33% - Stay and 67% - Leave). This ending node represents 16% of the aggregate populace.
7. For male, we further go down to performance variable. If the performance is less than 2.5 we go left else we go right.
8. For performance less than 2.5, we wind up in a terminating node that has a dominant response of 0 (59% - Stay and 41% - Leave). This ending node represents 16% of the aggregate populace.
9. For performance greater than 2.5, we wind up in a terminating node that has a dominant response of 1 (33% - Stay and 67% - Leave). This ending node represents 4% of the aggregate populace.

Assessing the predictive ability of the model

```
###Confusion Matrix
t_pred = predict(fit, test, type = 'class')
confMat = table(actual = test$vol_leave, prediction = t_pred)
confMat

##      prediction
## actual      0      1
##      0 2036   256
##      1  895   517

##Accuracy
accuracy = sum(diag(confMat)) / sum(confMat)

accuracy

## [1] 0.6892549
```

The accuracy is 0.689

- ❖ Logistic regression is better than decision tree in predicting the output response variable.

CONCLUSION

To play a more important and vital part in the organization, the HR function needs to move past beyond mere reporting to precise expectation. Rather than simply creating receptive reports, it needs to grasp advanced analytics and predictive techniques that bolster key organizational objectives.

Utilization of predictive analysis in HR involves using important information to take care of particular business issues. Predictive analyses helps associations contain HR-related expenses while building up a high performing workforce. The bits of knowledge inferred can enhance business execution and additionally employee engagement and satisfaction.

Predictive Analysis may be an unknown territory for HR, therefore to fully realize its benefits, HR personnel need to collaborate with other business units and customer-facing functions to understand how they leverage data and analytics to create value. By doing so, HR departments can facilitate superlative employee experiences that lead to sustained long-term benefits for the organization.



