

Busara Data Analysis

mburu

February 16, 2019

Contents

Task 1	2
<i>Summary Statistics Age and Income</i>	2
<i>% gender Gap</i>	2
<i>Single Ladies Nyeri</i>	3
<i>Summary Statistics Single Ladies Nyeri</i>	3
<i>Number of Juniors</i>	3
<i>Difference in mean income between male and female</i>	3
<i>Function to plot categorical variables</i>	4
<i>Function to plot categorical variables test 1</i>	4
<i>Function to plot categorical variables test 2</i>	5
Task 2	6
<i>Read Files</i>	6
<i>Head output data frame</i>	7
<i>Head output kenya data</i>	7
<i>Head output kenya data and saving files</i>	7
Task 3	8
<i>Project Motivatovation</i>	8
<i>Average Amount saved at the Bank</i>	8
<i>Average Amount saved at the Agent</i>	10
<i>Average Amount saved Mobile money</i>	11
<i>Demographic characteristics of those who have only made one deposit</i>	13
Gender	13
Region	13
Age	14
Task 3	14
<i>Project Motivation</i>	14
<i>Data Cleaning</i>	14
<i>Variable Selection</i>	15
<i>Visualize Categorical variables</i>	15
<i>Visualize numeric variables</i>	16
<i>One-Hot Encoding for categorical variables with more than 2 levels</i>	17
<i>Scale variables</i>	17
<i>Split test and train sets</i>	18
<i>Fit Logistic Regression</i>	18
<i>Confusion Matrix Logistic regression</i>	19
<i>Accuracy Logistic regression</i>	19
<i>Area under curve</i>	20
<i>Cross Validation SVM</i>	20
<i>Confusion Matrix SVM</i>	21
<i>Area under curve</i>	21
<i>Accuracy SVM</i>	21
<i>Validation Curves</i>	22

```
library(tidyverse)
library(readxl)
library(data.table)
library(knitr)
```

Task 1

```
xyz <- setDT(read_csv("XYZ.csv"))

xyz_sub <- xyz[, .(Gender, Age, Income)]

xyz_subm <- melt(xyz_sub, id.vars = "Gender")
```

Summary Statistics Age and Income

```
xyz_subm %>% group_by(variable) %>%
  summarise(Average = mean(value), Median = median(value),
            Min = min(value), Max = max(value)) %>%

  kable()
```

variable	Average	Median	Min	Max
Age	33.506	33	18	50
Income	5498.844	5557	1000	9897

% gender Gap

Yes there are 5.2% more men than women

```
gender <- xyz %>% group_by(Gender) %>%
  summarise(freq = n()) %>%
  mutate(Perc = round(freq/sum(freq) * 100, 2))

gender %>% kable()
```

Gender	freq	Perc
Female	237	47.4
Male	263	52.6

```
gender[, -2] %>% spread(Gender, Perc) %>%
  mutate(Percentage_Gender_Gap = Male - Female) %>% kable()
```

Female	Male	Percentage_Gender_Gap
47.4	52.6	5.2

Single Ladies Nyeri

```
single_nyeri <- xyz[Gender == "Female" & Marital_Status == "Single" & County == "Nyeri",]  
cat("The Number of single ladies in Nyeri is ", nrow(single_nyeri))
```

```
## The Number of single ladies in Nyeri is 12
```

Summary Statistics Single Ladies Nyeri

```
single_nyeri %>%  
  summarise(Average_Age = mean(Age), Median_Income = median(Income)) %>%  
  kable()
```

Average_Age	Median_Income
36	5557

Number of Juniors

```
juniors_26 <- xyz[!grepl("Operations|Data", Department) & xyz$Age < 26 & grepl("Junior", Role),]  
cat("The Number of juniors ", nrow(juniors_26))
```

```
## The Number of juniors 28
```

```
juniors_26 %>% group_by(Department) %>%  
  summarise(freq = n()) %>%  
  mutate(Perc = round(freq/sum(freq) * 100, 2)) %>%  
  kable()
```

Department	freq	Perc
Associate	5	17.86
Finance	11	39.29
Operations	8	28.57
Research Analyst	4	14.29

Difference in mean income between male and female

The Operations has the biggest difference in mean income

```
income_gender <- xyz %>% group_by(Gender, Department) %>%  
  summarise(Average = mean(Income))  
  
income_gender_dcast <- dcast(Department ~ Gender, data = income_gender)  
  
income_gender_dcast %>% mutate(Difference = Male - Female) %>%  
  kable()
```

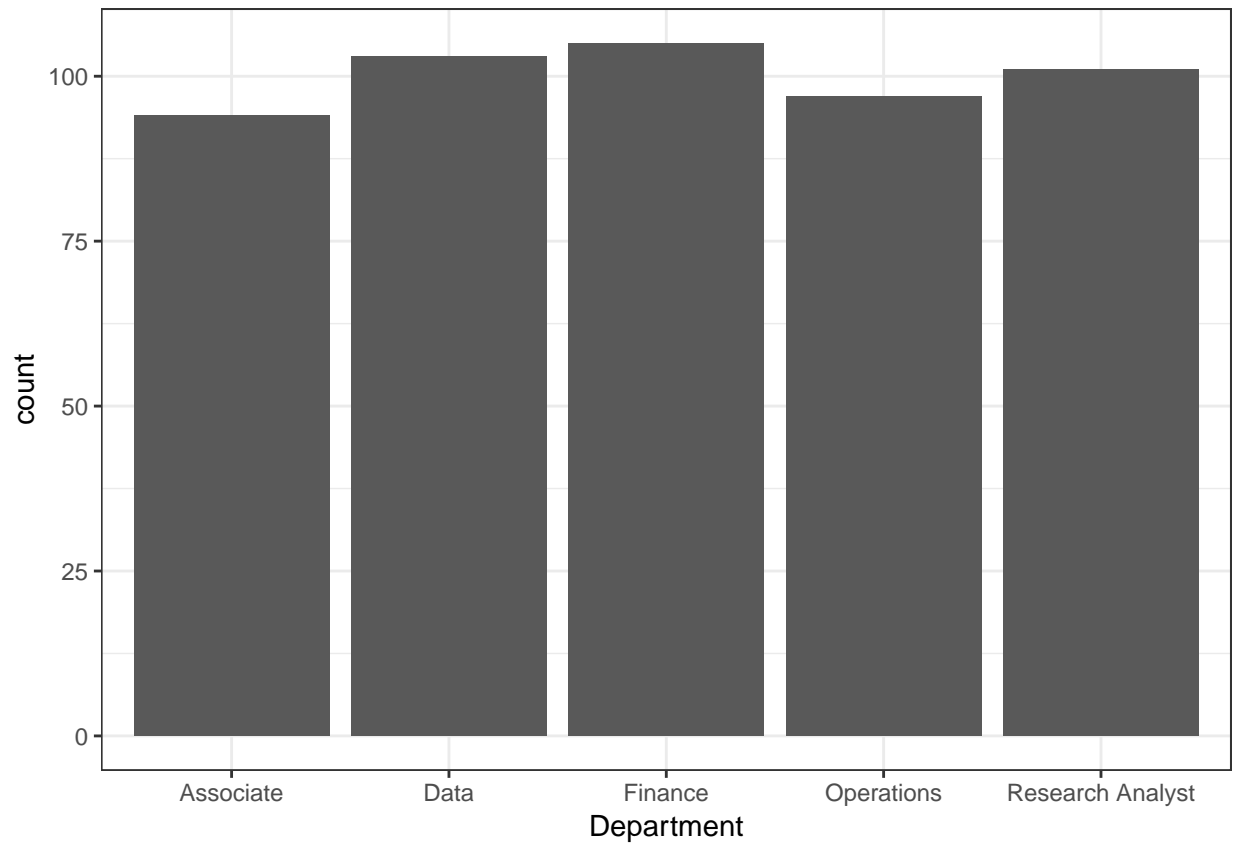
Department	Female	Male	Difference
Associate	5345.047	5071.941	-273.1053
Data	5613.420	5270.396	-343.0238
Finance	5574.714	5936.750	362.0357
Operations	5043.286	6284.854	1241.5685
Research Analyst	5264.522	5533.327	268.8055

Function to plot categorical variables

```
bar_plot <- function(data, ...) {
  #load ggplot2
  #function takes a data frame
  #and other arguments that ggplot
  #function from ggplot2 takes
  # the other arguments are aesthetic mappings
  require(ggplot2)
  ggplot(data) + geom_bar(aes(...))+
    theme_bw()
}
```

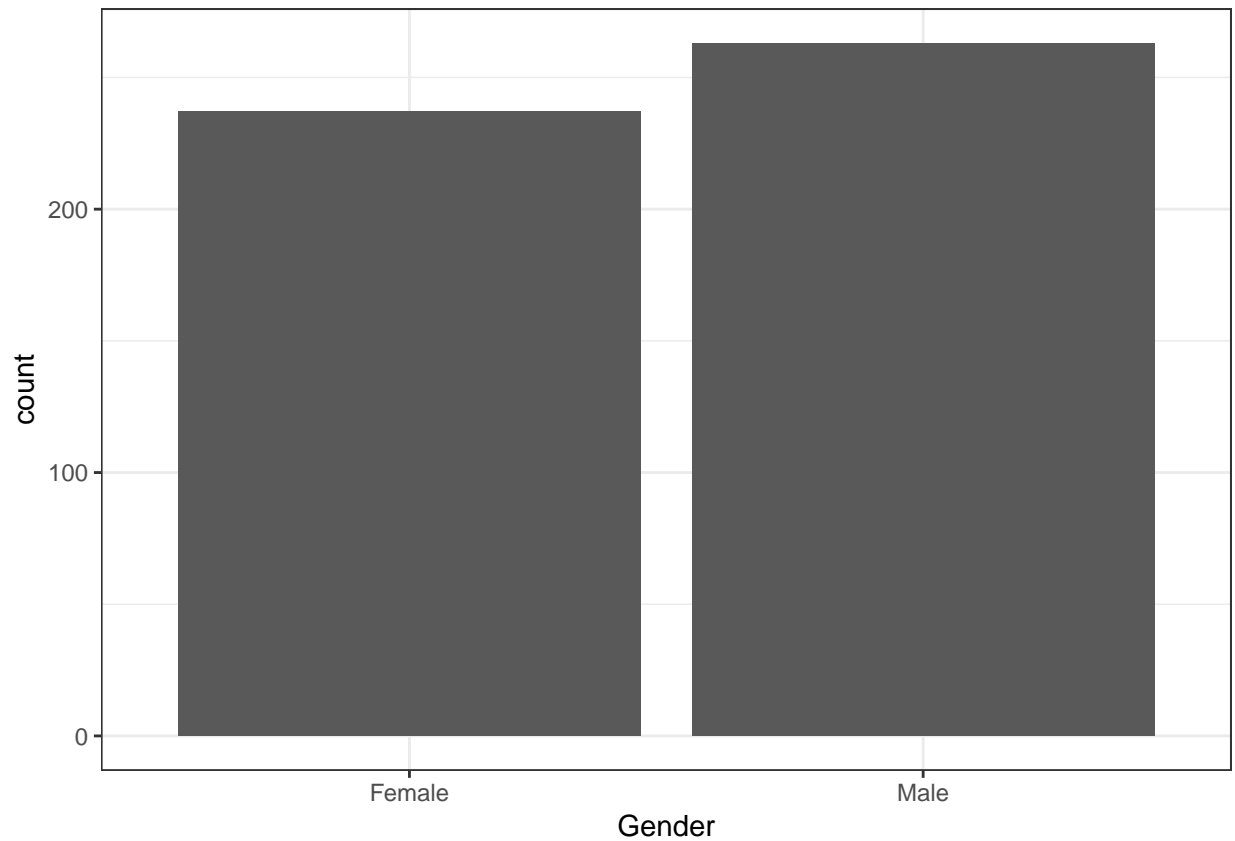
Function to plot categorical variables test 1

```
bar_plot(xyz, Department)
```



Function to plot categorical variables test 2

```
bar_plot(xyz, Gender)
```



Task 2

Read Files

```
in.dir = paste(getwd(), "/Education", sep = "")
setwd(in.dir)
my_files <- dir(pattern = "^Chi|^Sch|^Persi|Secon|^Progr|Pri")

library(readxl)
list_files <- list()

for (i in 1:length(my_files)) {

  x = read_excel(my_files[i])
  id = grep("Country Name", x$`Data Source`)
  nms <- x[id,]
  names(x) <- nms %>% as.character()
  list_files[[i]] <- x[-c(1:id),]
  cat("...")

}
```

```
## .....
df_world <- rbindlist(list_files) %>% setDT()

df_world_melt <- melt(df_world, id.vars = names(df_world)[1:4])

nms2 <- Hmisc::Cs(Country_Name, Country_Code,
                 Indicator_Name, Indicator_Code, Year, Indicator_value)

names(df_world_melt) <- nms2

df_world_melt[, Year := as.numeric(as.character(df_world_melt$Year))]
```

Head output data frame

```
head(df_world_melt) %>% kable()
```

Country_Name	Country_Code	Indicator_Name	Indicator_Code	Year	Indicator_value
Aruba	ABW	Children out of school, primary, female	SE.PRM.UNER.FE	1960	NA
Afghanistan	AFG	Children out of school, primary, female	SE.PRM.UNER.FE	1960	NA
Angola	AGO	Children out of school, primary, female	SE.PRM.UNER.FE	1960	NA
Albania	ALB	Children out of school, primary, female	SE.PRM.UNER.FE	1960	NA
Andorra	AND	Children out of school, primary, female	SE.PRM.UNER.FE	1960	NA
Arab World	ARB	Children out of school, primary, female	SE.PRM.UNER.FE	1960	NA

Head output kenya data

```
kenya_2011 <- df_world_melt[Country_Name == "Kenya" & Year >= 2011]

head(kenya_2011) %>% kable()
```

Country_Name	Country_Code	Indicator_Name	Indicator_Code
Kenya	KEN	Children out of school, primary, female	SE.PRM.UNER.FE
Kenya	KEN	Persistence to last grade of primary, female (% of cohort)	SE.PRM.PRSL.FE.ZS
Kenya	KEN	Persistence to last grade of primary, male (% of cohort)	SE.PRM.PRSL.MA.ZS
Kenya	KEN	Primary completion rate, female (% of relevant age group)	SE.PRM.CMPT.FE.ZS
Kenya	KEN	Primary completion rate, male (% of relevant age group)	SE.PRM.CMPT.MA.ZS
Kenya	KEN	Progression to secondary school, female (%)	SE.SEC.PROG.FE.ZS

Head output kenya data and saving files

```
write.csv(head(kenya_2011, 15), file = "kenya data.csv", row.names = F)

kenya_2011_na <- kenya_2011[!is.na(kenya_2011$Indicator_value),]
write.csv(head(kenya_2011_na, 15), file = "kenya data without na.csv", row.names = F)

head(kenya_2011_na) %>% kable()
```

Country_Name	Country_Code	Indicator_Name	Indicator_Code
Kenya	KEN	Children out of school, primary, female	SE.PRM.UNER.FE
Kenya	KEN	School enrollment, primary (gross), gender parity index (GPI)	SE.ENR.PRIM.FM.ZS
Kenya	KEN	School enrollment, primary, female (% gross)	SE.PRM.ENRR.FE
Kenya	KEN	School enrollment, primary, male (% gross)	SE.PRM.ENRR.MA
Kenya	KEN	Primary completion rate, female (% of relevant age group)	SE.PRM.CMPT.FE.ZS
Kenya	KEN	Primary completion rate, male (% of relevant age group)	SE.PRM.CMPT.MA.ZS

Task 3

Project Motivatovation

```

figari_sheet1 <- read_excel("Figari Bank.xlsx" ) %>% setDT()

figari_sheet2 <- read_excel("Figari Bank.xlsx", sheet = 2 ) %>% setDT()

figari_sheet2[, Dates := as.Date(Dates, origin = "1900-01-01")]

figari_sheet2[, year := year(Dates)]

figari_sheet2[, month := months(Dates)]

figari_sheet2[, week_day := weekdays(Dates)]

figari_sheet2[, week_no := week(Dates)]
figari_sheet2[, day_month := format(Dates, "%d")]

figari_sheet2_m <- melt(figari_sheet2[, c(3:9), with = F], id.vars = c("Amount", "Saving Mode"))

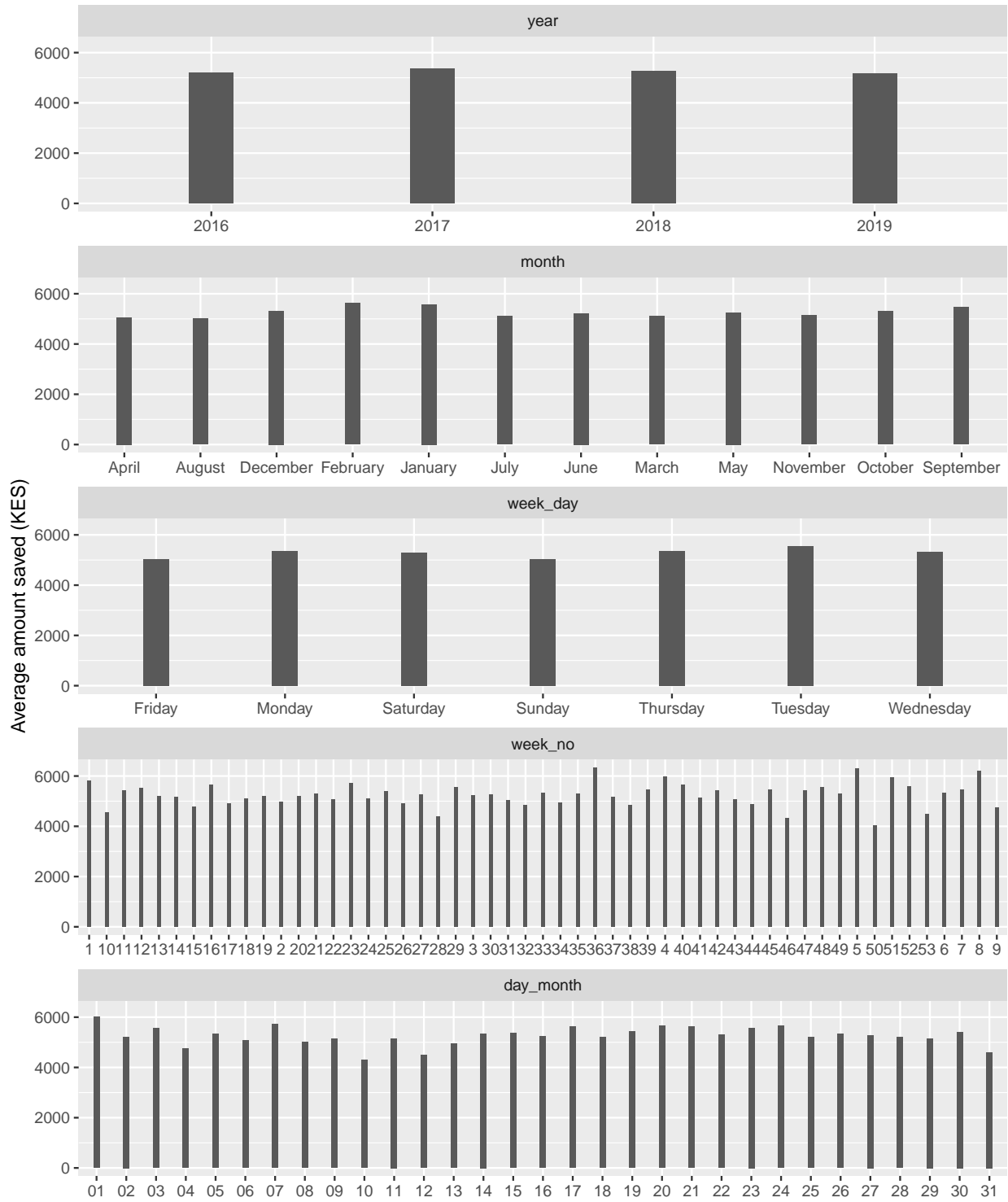
figari_dat <- figari_sheet2_m %>% group_by(`Saving Mode`, variable, value) %>%
  summarise(Average = mean(Amount))
titles <- levels(as.factor(figari_dat$`Saving Mode`))
titles <- paste("Average Savings for", titles)
figari_dat_split <- split(figari_dat, figari_dat$`Saving Mode`)
plots_figari <- list()
for ( i in 1:length(figari_dat_split)) {
  this = figari_dat_split[[i]]
  plots_figari[[i]] <- ggplot(this, aes(value, Average)) +
    geom_bar(stat = "identity", width = 0.2) +
    facet_wrap(~variable, scales = "free_x", ncol = 1)+
    labs(x = "", y = "Average amount saved (KES)", title = titles[i])# +
    #theme(axis.text.x = element_text(angle = 30, vjust = 1, hjust = 1))
}

```

Average Amount saved at the Bank

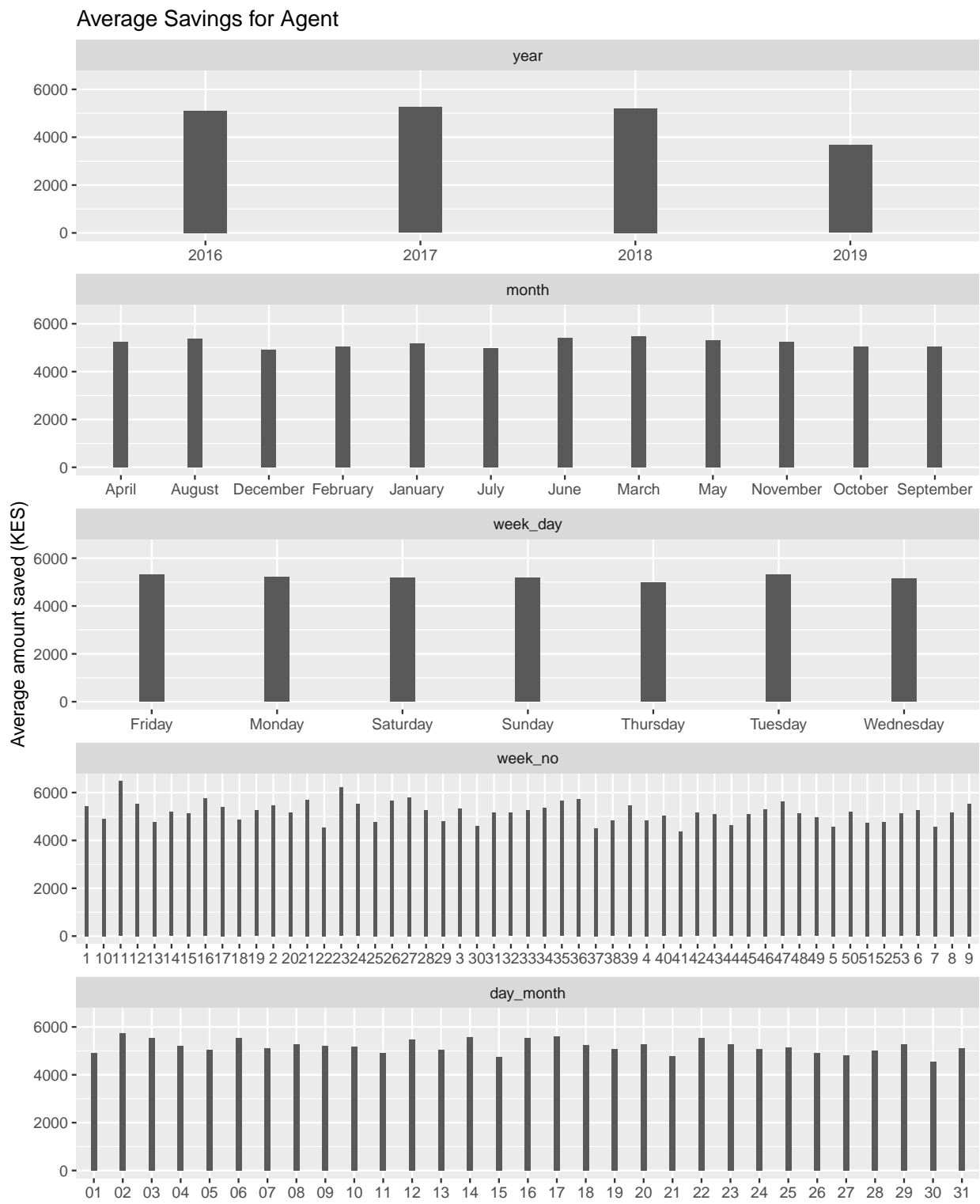
```
plots_figari[[2]]
```


Average Savings for Bank



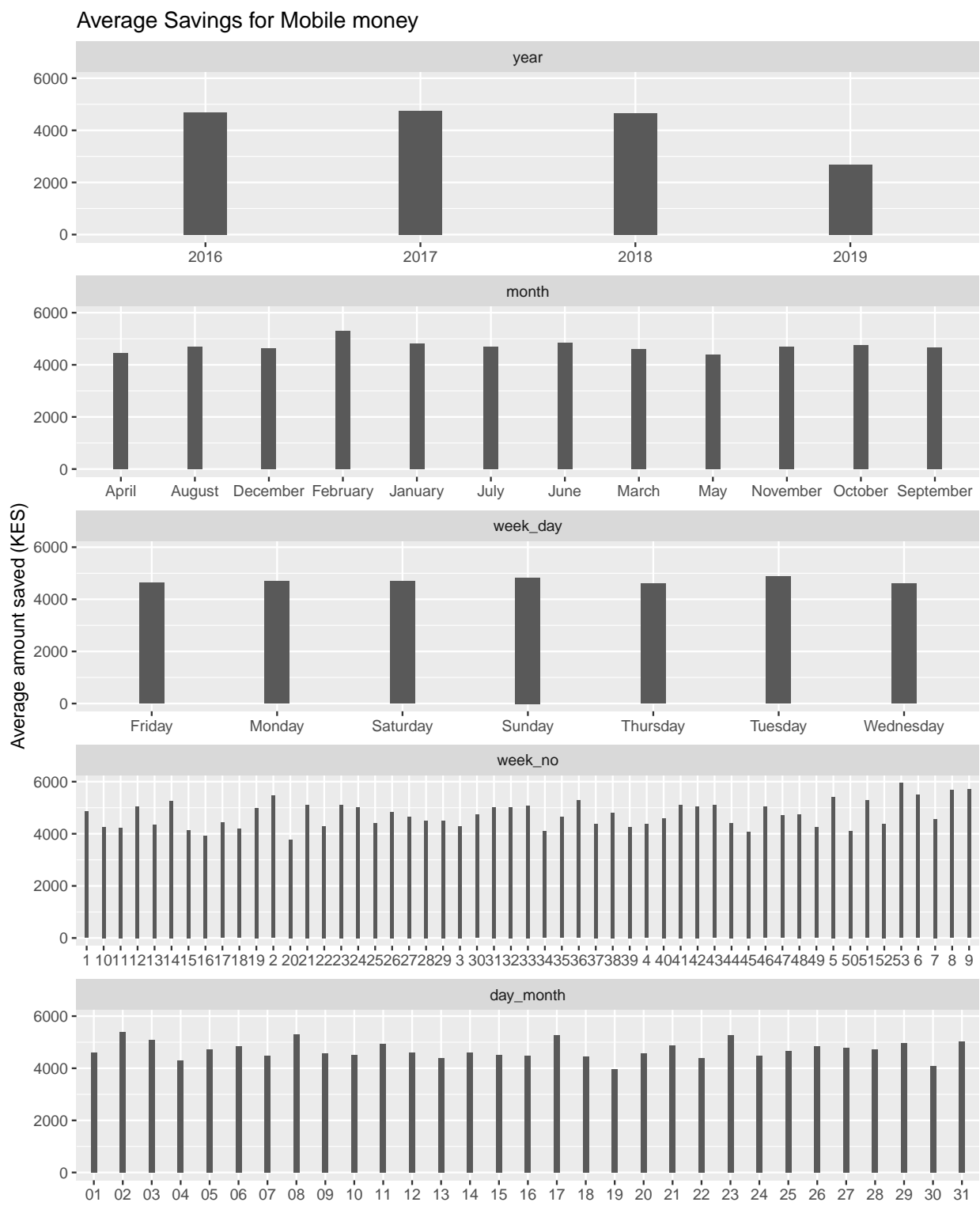
Average Amount saved at the Agent

```
plots_figari[[1]]
```

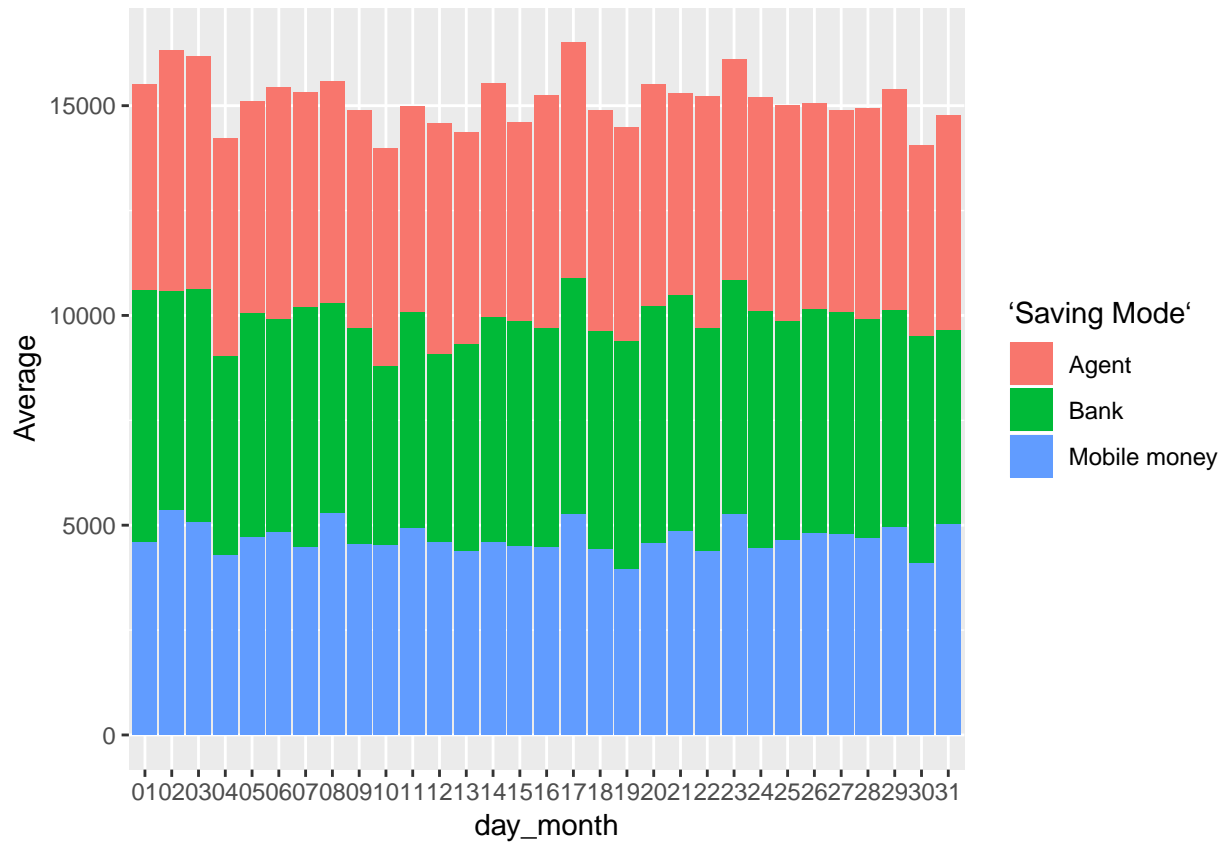


Average Amount saved Mobile money

```
plots_figari[[3]]
```

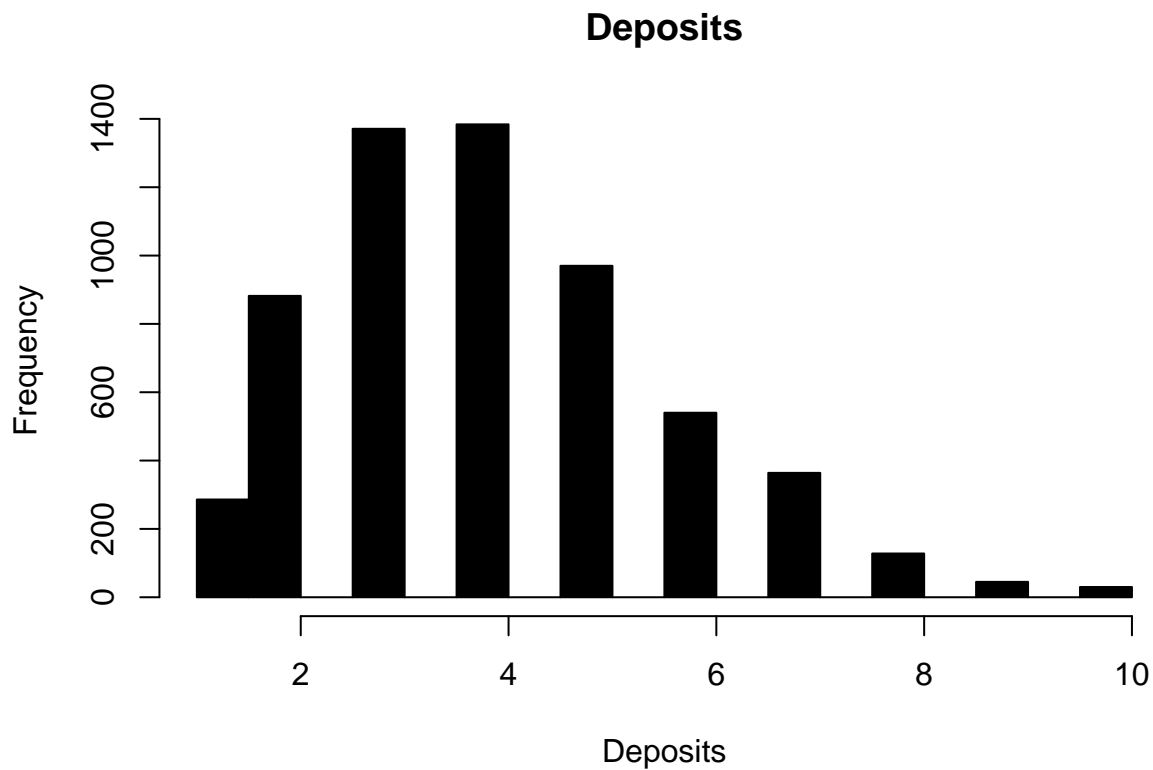


```
end_month <- figari_sheet2 %>% group_by(day_month, `Saving Mode`) %>%
  summarise(Average = mean(Amount))
#The number of times one deposits
ggplot(end_month, aes(day_month, Average, fill = `Saving Mode`)) +
  geom_bar(stat = "identity")
```



```
deposits <- figari_sheet2[, .(`Customer ID`)] %>% group_by(`Customer ID`) %>%
  group_by(freq = n()) %>% setDT()

#approximately poison
hist(deposits$freq, col = "black", main = "Deposits", xlab = "Deposits")
```



```
names(deposits)[1] = names(figari_sheet1)[1]
figari_deposits <- merge(deposits, figari_sheet1, by = "CustomerID")

figari_deposits_one <- figari_deposits[freq == 1]
```

Demographic characteristics of those who have only made one deposit

Gender

```
figari_deposits_one %>% group_by(Gender) %>%
  summarise(freq= n()) %>%
  mutate(Perc =round(freq/sum(freq) *100, 2) ) %>%
  kable()
```

Gender	freq	Perc
Female	141	49.3
Male	145	50.7

Region

```
figari_deposits_one %>% group_by(Region) %>%
  summarise(freq= n()) %>%
```

```
mutate(Perc =round(freq/sum(freq) *100, 2) ) %>%
kable()
```

Region	freq	Perc
Bondo	25	8.74
Gatitu	53	18.53
Kawangware	21	7.34
Kayole	13	4.55
Kibera	31	10.84
Kilimani	41	14.34
Kirinyaga	18	6.29
Rongai	10	3.50
Ruai	45	15.73
Taita	29	10.14

Age

```
figari_deposits_one %>%
  summarise(Mean= round(mean(Age), 2), Median = median(Age))
```

```
##      Mean Median
## 1 54.71      56
```

Task 3

Project Motivation

Data has the potential to transform business and drive the creation of business value. It can be used for a range of tasks such visualization relationships between variables to predicting if an event will occur. The later is one of the heavily reaserched areas in recent times. The reason for this is that data has grown exponentially and so does the computing power. Banks and financial institutions used data analytics for a range of value such as fraud detetction customer segment, recruiting, credit scoring and so on.

In this study I will use Bogoza data set to build a credit model where an applicat will be avaluated on whether they will default or not.

High accuracy for this model will be required because predicting false positives will eventually cause a business to make a loss and false negatives means that the financial instituion looses business.

Data Cleaning

First step is data cleaning. This ensures that columns are consistent. For instance the target variable had values such as Y y yes where all of them represent yes.

```
#some algorithms like xgboost take numeric data
#you can convert binary vars to 1,0
# and form dummie variables using library dummies
#for variables with more than 2 categories
borogoza <- setDT(read_csv( "Bagorogoza Loan.csv"))
```

```

borogoza[, Target := ifelse(grepl("y|Y", Target), 1, 0)]

borogoza[, Gender := ifelse(grepl("^m$|^male$", tolower(Gender)), 0, 1)]

borogoza[, Married := ifelse(grepl("Yes",Married), 1, 0)]

borogoza[, Education := ifelse(grepl("not", tolower(Education)), 0, 1)]

borogoza[, Self_Employed := ifelse(grepl("Yes",Self_Employed), 1, 0)]

borogoza[, Property_Area := ifelse(grepl("rural",tolower(Property_Area)), "Rural", Property_Area)]

borogoza[, Property_Area := ifelse(grepl("semi",tolower(Property_Area)), "Semi-urban", Property_Area)]

borogoza[, Property_Area := ifelse(grepl("^urban$",tolower(Property_Area)), "Urban", Property_Area)]

```

Variable Selection

Visualize Categorical variables

Visualization and summary statistics is an important step before fitting any model as this will give you a glimpse of how the variables are associated with target variable. In this case I will use stacked barplot as from them you can see if the proportions of defaulters and non defaulters is equal in different categories of a variable. From the graphs we can see that the proportion of defaulters and non defaulters is different for the different credit history categories. This is also seen in the property area. From the categorical variables we can therefore conclude that one of the best predictors is credit history.

```

numeric_vars <- Hmisc::Cs(ApplicantIncome,CoapplicantIncome, LoanAmount )

nms_bo <- names(borogoza)[-1]

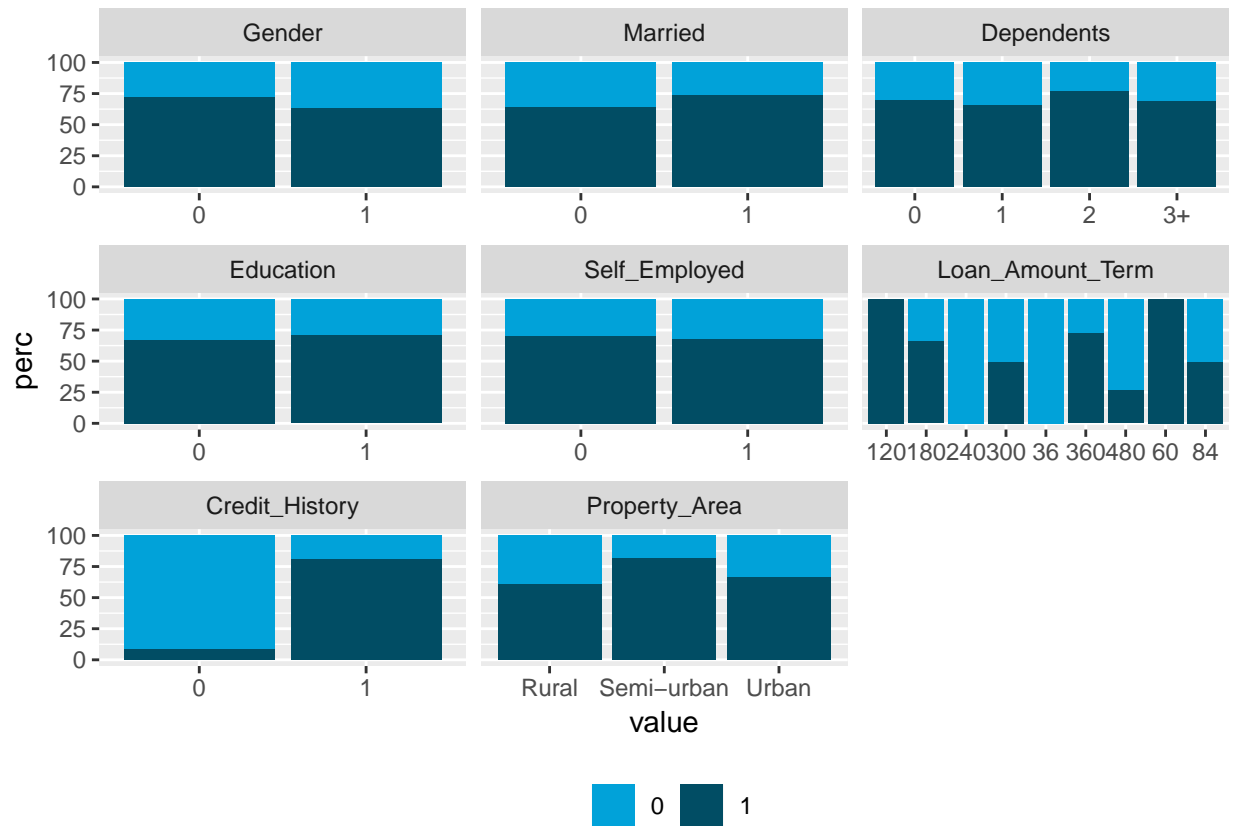
cat_vars <- nms_bo[!nms_bo %in% numeric_vars]

borogoza_catm <- melt(borogoza[, cat_vars, with = F], id.vars = "Target")

borogoza_catm_perc <-borogoza_catm %>% group_by(variable, value, Target) %>%
  summarise(freq= n()) %>% mutate(perc =round(freq/sum(freq) *100, 2) )

library(ggthemes)
ggplot(borogoza_catm_perc, aes(value, perc, fill = factor(Target) )) +
  geom_bar(stat = "identity") +facet_wrap(~variable, scales = "free_x")+
  scale_fill_economist(name = "")+
  theme(legend.position = "bottom")

```

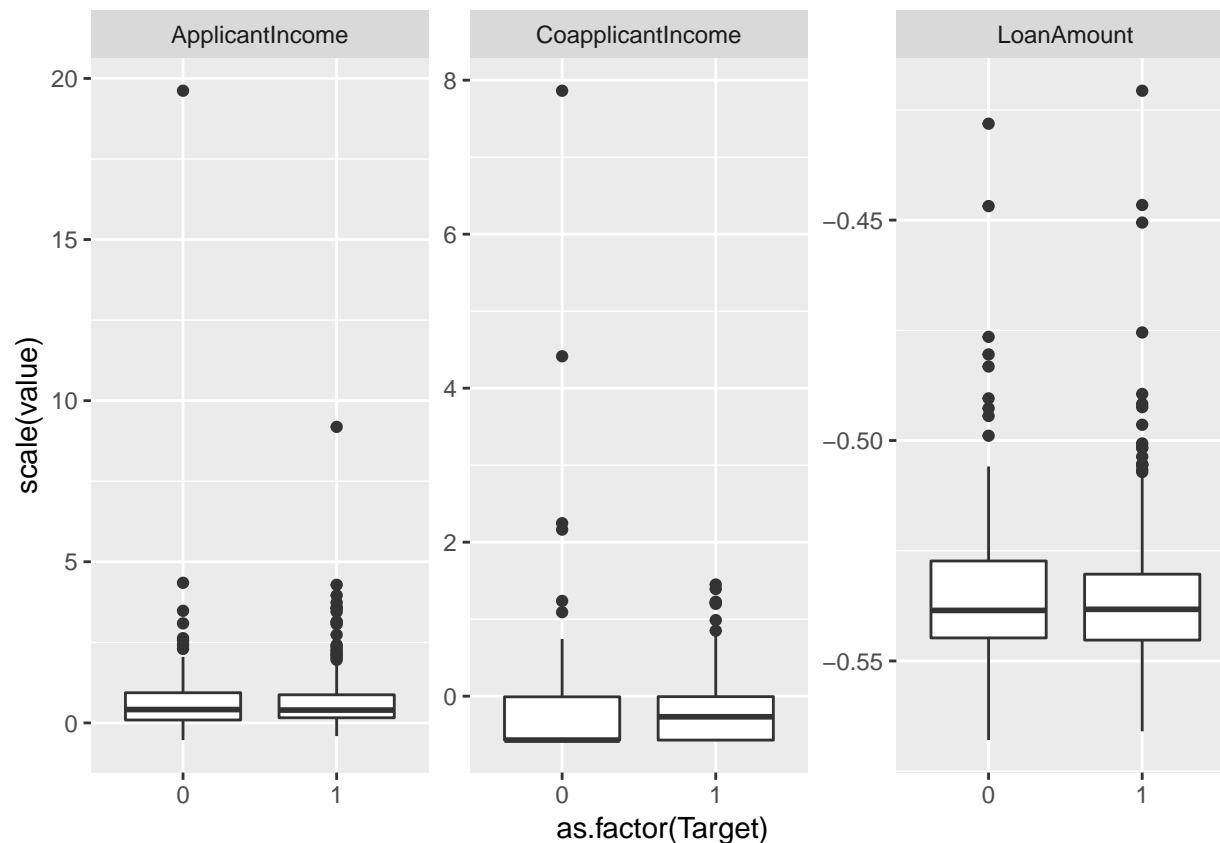


Visualize numeric variables

For the numeric variables boxplot help us visualize which distribution is different from the other. Non overlapping boxplot for defaulters and non defaulters may indicate that the mean/median values in the two groups was significantly different. From this we can see that it's unlikely that education and self employment affect loan repayment and for this we drop this two variables

```
borogoza_numm <- melt(borogoza[, c(numeric_vars, "Target"), with = F], id.vars = "Target")

ggplot(borogoza_numm, aes(as.factor(Target), scale(value))) +
  geom_boxplot() + facet_wrap(~variable, scales = "free_y")
```

One-Hot Encoding for categorical variables with more than 2 levels

In this step variables with more than two categories are converted to dummies variables. The first column in each category is dropped as it's linearly dependent with the second column.

```
chars <- unlist(lapply(borogoza[, -1, with = F], is.character))

chars <- nms_bo[chars]

library(dummies)
borogoza_dummy <- dummy.data.frame(borogoza, names = c(chars, "Loan_Amount_Term")) %>%
  setDT()

borogoza_dummy[, Loan_ID := NULL]
borogoza_dummy[, Loan_Amount_Term36 := NULL]
borogoza_dummy[, `Property_AreaSemi-urban` := NULL]
borogoza_dummy[, `Dependents1` := NULL]
```

Scale variables

It's important to scale your variables since it leads to faster convergence and since some algorithm use distances to find decision boundary this means that variables with big values will have a big influence.

```
xvars <- names(borogoza_dummy)[!names(borogoza_dummy) %in% "Target"]
borogoza_dummy[, (xvars) := lapply(.SD, function(x) scale(x)), .SDcols = xvars ]
```

Split test and train sets

This is important as it helps evaluate your model on data it has never seen. The model will be trained on one set(training set) and tested using test set.

```
set.seed(200) # for reproducibility
train_sample <- sample(1:nrow(borogoza_dummy), round(0.7*nrow(borogoza_dummy)))

train <- borogoza_dummy[train_sample,]
test <- borogoza_dummy[-train_sample,]
```

Fit Logistic Regression

Logistic regression was fit to predict the probability of someone defaulting. The advantages of logistic regression is interpretable, ie you can see the association between a predictor and response value, it also gives a probability. This is very important when you want to have your own cut off point eg you want to label someone as a defaulter if you the predicted probability is more than 0.7. This increases precision but lowers recall. Using stepwise selection the model was used to select the variables that best predict loan default.

```
fit_glm <- glm(Target ~ Married + CoapplicantIncome + Loan_Amount_Term60 +
  Loan_Amount_Term120 + Loan_Amount_Term180 + Loan_Amount_Term300 +
  Loan_Amount_Term360 + Credit_History + Property_AreaRural +
  Property_AreaUrban ,data = train, family = binomial)

borogoza_dummy <- borogoza_dummy[, .(Target,Married , CoapplicantIncome , Loan_Amount_Term60 ,
  Loan_Amount_Term120 , Loan_Amount_Term180 , Loan_Amount_Term300 ,
  Loan_Amount_Term360 , Credit_History , Property_AreaRural ,
  Property_AreaUrban)]

train <- borogoza_dummy[train_sample,]
test <- borogoza_dummy[-train_sample,]
summary(fit_glm) %>% xtable::xtable() %>% kable()
```

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	1.0943288	10.9772973	0.0996902	0.9205903
Married	0.3517280	0.1702977	2.0653720	0.0388878
CoapplicantIncome	-0.2832434	0.1712880	-1.6536094	0.0982069
Loan_Amount_Term60	1.2882921	104.8970941	0.0122815	0.9902010
Loan_Amount_Term120	1.5046767	89.5711230	0.0167987	0.9865972
Loan_Amount_Term180	1.1707354	0.3708054	3.1572773	0.0015925
Loan_Amount_Term300	0.3136355	0.2170364	1.4450824	0.1484347
Loan_Amount_Term360	1.3749642	0.4128089	3.3307522	0.0008661
Credit_History	1.3791744	0.2294197	6.0115776	0.0000000
Property_AreaRural	-0.4577834	0.2067118	-2.2145972	0.0267877
Property_AreaUrban	-0.5703610	0.2111979	-2.7006002	0.0069214

```
#MASS::stepAIC(fit_glm)
```

The estimate column shows the log odds. Positive values means that the variable makes it more likely for a person to repay their loan negative values means that the person is less likely to repay.

Confusion Matrix Logistic regression

The confusion matrix evaluate correctly classified cases. A perfect fit will have all values in the main diagonal while the entries of lower/upper triangulars should be zeros. In this case we have 14 cases of false positives and 7 cases of false negatives the accuracy of the model is 0.82 with and f1 score of 0.87. F1 score is a very important evaluation metric where there is unbalanced classes.

```
library(caret)
pred_glm <- predict(fit_glm,newdata = test)

pred_glm <- ifelse(pred_glm>0.7, 1 , 0)

table(test$Target, pred_glm) %>% kable()
```

	0	1
0	19	14
1	9	73

Accuracy Logistic regression

```
library(broom)
library(pROC)
table(test$Target, pred_glm) %>%
  confusionMatrix(positive = "1") %>%
  tidy() %>% kable()
```

term	class	estimate	conf.low	conf.high	p.value
accuracy	NA	0.8000000	0.7151708	0.86879	0.1642093
kappa	NA	0.4880975	NA	NA	0.4042485
sensitivity	1	0.8390805	NA	NA	NA
specificity	1	0.6785714	NA	NA	NA
pos_pred_value	1	0.8902439	NA	NA	NA
neg_pred_value	1	0.5757576	NA	NA	NA
precision	1	0.8902439	NA	NA	NA
recall	1	0.8390805	NA	NA	NA
f1	1	0.8639053	NA	NA	NA
prevalence	1	0.7565217	NA	NA	NA
detection_rate	1	0.6347826	NA	NA	NA
detection_prevalence	1	0.7130435	NA	NA	NA
balanced_accuracy	1	0.7588259	NA	NA	NA

Area under curve

This is important as it will help you know if the suffers from high false negatives or false positives. A value greater than 0.8 is normally desired in this case we achieve 0.74.

```
roc(as.numeric(test$Target), pred_glm, print.auc=T, print.auc.y=0.5, levels =0:1 )
```

```
##
```

```
## Call:
```

```
## roc.default(response = as.numeric(test$Target), predictor = pred_glm,      levels = 0:1, print.auc = T)
```

```
##
```

```
## Data: pred_glm in 33 controls (as.numeric(test$Target) 0) < 82 cases (as.numeric(test$Target) 1).
```

```
## Area under the curve: 0.733
```

Cross Validation SVM

Next we fit Support vector machine model. We start by finding the best parameters using cross validation. We use 10 fold this where train set is randomly split into 10 sets. In each cases one of the 1 set is used as a valiadtion/test set while the other 9 are used to train the model.

```
library(e1071)
tune.out = tune(svm, as.factor(Target)~., data = train, kernel ="radial",
               type ="C-classification",
               ranges =list (cost=c(0.01, 0.1, 1 ,5 , 10),
                           gamma = c(0.01, 0.1, 1 ,5 )))
```

```
summary(tune.out)
```

```
##
```

```
## Parameter tuning of 'svm':
```

```
##
```

```
## - sampling method: 10-fold cross validation
```

```
##
```

```
## - best parameters:
```

```
##   cost gamma
```

```
##    5  0.01
```

```
##
```

```
## - best performance: 0.1707977
```

```
##
```

```
## - Detailed performance results:
```

```
##   cost gamma   error dispersion
```

```
## 1  0.01  0.01 0.2971510 0.11255905
```

```
## 2  0.10  0.01 0.2971510 0.11255905
```

```
## 3  1.00  0.01 0.1931624 0.08664175
```

```
## 4  5.00  0.01 0.1707977 0.08557146
```

```
## 5 10.00  0.01 0.1707977 0.08557146
```

```
## 6  0.01  0.10 0.2971510 0.11255905
```

```
## 7  0.10  0.10 0.2045584 0.09624777
```

```
## 8  1.00  0.10 0.1709402 0.08574399
```

```
## 9  5.00  0.10 0.1709402 0.08574399
```

```
## 10 10.00  0.10 0.1709402 0.08574399
```

```
## 11 0.01  1.00 0.2971510 0.11255905
```

```
## 12 0.10  1.00 0.2971510 0.11255905
```

```
## 13 1.00  1.00 0.1971510 0.10215705
```

```
## 14  5.00  1.00  0.2008547  0.10672042
## 15 10.00  1.00  0.2045584  0.11367327
## 16  0.01  5.00  0.2971510  0.11255905
## 17  0.10  5.00  0.2971510  0.11255905
## 18  1.00  5.00  0.2451567  0.10166867
## 19  5.00  5.00  0.2340456  0.10593807
## 20 10.00  5.00  0.2415954  0.10204222
```

Confusion Matrix SVM

```
fit_svm <- svm(as.factor(Target)~., data = train, cost =5 , gamma = .01,
              kernel = "radial", type ="C-classification")

pred_svm <-predict(fit_svm, newdata = test)
table(test$Target, pred_svm) %>% kable()
```

	0	1
0	19	14
1	7	75

Area under curve

```
roc(test$Target, as.numeric(pred_svm), print.auc=T, print.auc.y=0.5, levels =0:1 )
```

```
##
```

```
## Call:
```

```
## roc.default(response = test$Target, predictor = as.numeric(pred_svm),      levels = 0:1, print.auc = 'T')
```

```
##
```

```
## Data: as.numeric(pred_svm) in 33 controls (test$Target 0) < 82 cases (test$Target 1).
```

```
## Area under the curve: 0.7452
```

Accuracy SVM

```
table(test$Target, pred_svm) %>%
  confusionMatrix(positive = "1") %>%
  tidy() %>% kable()
```

term	class	estimate	conf.low	conf.high	p.value
accuracy	NA	0.8173913	0.7345109	0.883263	0.1577134
kappa	NA	0.5235747	NA	NA	0.1904303
sensitivity	1	0.8426966	NA	NA	NA
specificity	1	0.7307692	NA	NA	NA
pos_pred_value	1	0.9146341	NA	NA	NA
neg_pred_value	1	0.5757576	NA	NA	NA
precision	1	0.9146341	NA	NA	NA
recall	1	0.8426966	NA	NA	NA
f1	1	0.8771930	NA	NA	NA

term	class	estimate	conf.low	conf.high	p.value
prevalence	1	0.7739130	NA	NA	NA
detection_rate	1	0.6521739	NA	NA	NA
detection_prevalence	1	0.7130435	NA	NA	NA
balanced_accuracy	1	0.7867329	NA	NA	NA

Validation Curves

The two models almost give equal results based on accuracy, f1 score and area under the curve. In this section we will evaluate the models using learning curves to see if they suffer from high variance or bias. In this case the model suffers from high bias. It's evident that adding more data won't solve accuracy problems. In this case additional features would help.

```
sets <- seq(from = 50, to = nrow(train), by = 50)
sets[length(sets)] <- nrow(train)
train.err <- c()
test.err <- c()
for (i in 1:length(sets)) {

  traini = train[1:sets[i],]
  fit_svm <- svm(as.factor(Target)~., data = traini, cost = 5, gamma = .01,
    kernel = "radial", type = "C-classification")

  pred_train = predict(fit_svm, newdata = traini)
  train.err[i] = 1 - mean(pred_train == traini$Target)
  pred_test <- predict(fit_svm, newdata = test)
  test.err[i] = 1 - mean(test$Target == pred_test)

  cat(i, " ")

}

## 1 2 3 4 5

matplot(sets, cbind(test.err, train.err), pch = 19, col = c("red", "blue"),
  type = "b", ylab = "Error", xlab = "Train sample size", main = "SVM Training and Validation error",
  legend("topright", legend = c("Test", "Train"), pch = 19, col = c("red", "blue")))
```



Deployment

Other model like Xgboost which uses boosting and bagging could first be used to see if the model performs better on this data. The problem could after this be intergrated with a loan evaluation software where it can help loan officers decide if the will award a loan.