

# HW6

Yufei

3/21/2020

## R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com> (<http://rmarkdown.rstudio.com>).

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
#load the packages
library(readr)
library(readxl)
library(forecast)
```

```
## Registered S3 method overwritten by 'quantmod':
##   method              from
##   as.zoo.data.frame zoo
```

```
library(tidyverse)
```

```
## — Attaching packages ————— tid
yverse 1.2.1 —
```

```
## ✓ ggplot2 3.2.1      ✓ purrr    0.3.3
## ✓ tibble  2.1.3      ✓ dplyr    0.8.3
## ✓ tidyr   1.0.0      ✓ stringr  1.4.0
## ✓ ggplot2 3.2.1      ✓ forcats  0.4.0
```

```
## — Conflicts — tidyverse
_conflicts() —
## ✖ dplyr::filter() masks stats::filter()
## ✖ dplyr::lag() masks stats::lag()
```

```
library(caret)
```

```
## Loading required package: lattice
```

```
##
## Attaching package: 'caret'
```

```
## The following object is masked from 'package:purrr':
##
## lift
```

```
library(rpart)
library(caret)
library(e1071)
library(data.table)
```

```
##
## Attaching package: 'data.table'
```

```
## The following objects are masked from 'package:dplyr':
##
## between, first, last
```

```
## The following object is masked from 'package:purrr':
##
## transpose
```

```
library(randomForest)
```

```
## randomForest 4.6-14
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##  
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:dplyr':  
##  
##      combine
```

```
## The following object is masked from 'package:ggplot2':  
##  
##      margin
```

```
library(leaps)  
library(MASS)
```

```
##  
## Attaching package: 'MASS'
```

```
## The following object is masked from 'package:dplyr':  
##  
##      select
```

```
library(readr)  
library(corrplot)
```

```
## corrplot 0.84 loaded
```

```
library(gridExtra)
```

```
##  
## Attaching package: 'gridExtra'
```

```
## The following object is masked from 'package:randomForest':  
##  
##      combine
```

```
## The following object is masked from 'package:dplyr':  
##  
##      combine
```

```
library(formattable)
```

```
##  
## Attaching package: 'formattable'
```

```
## The following object is masked from 'package:MASS':  
##  
##      area
```

```
library(readr)  
tele_Churn <- read_csv("~/Desktop/WA_Fn-UseC_-Telco-Customer-Churn.csv")
```

```
## Parsed with column specification:  
## cols(  
##   .default = col_character(),  
##   SeniorCitizen = col_double(),  
##   tenure = col_double(),  
##   MonthlyCharges = col_double(),  
##   TotalCharges = col_double()  
## )
```

```
## See spec(...) for full column specifications.
```

```
#check the missing values
```

```
apply(tele_Churn, function(x) sum(is.na(x)))
```

```
##          customerID          gender  SeniorCitizen          Partn
er
##          0          0          0
0
##          Dependents          tenure  PhoneService  MultipleLin
es
##          0          0          0
0
##  InternetService  OnlineSecurity  OnlineBackup  DeviceProtecti
on
##          0          0          0
0
##          TechSupport  StreamingTV  StreamingMovies          Contra
ct
##          0          0          0
0
##  PaperlessBilling  PaymentMethod  MonthlyCharges  TotalCharg
es
##          0          0          0
11
##          Churn
##          0
```

```
#drop the missing values
```

```
telechurn01 <-na.omit(tele_Churn)
```

```
#change the seniorCitizen column to factor
```

```
telechurn02 <- telechurn01%>%
```

```
  mutate(SeniorCitizen = as.factor(SeniorCitizen))
```

```
telechurn02$customerID <- NULL
```

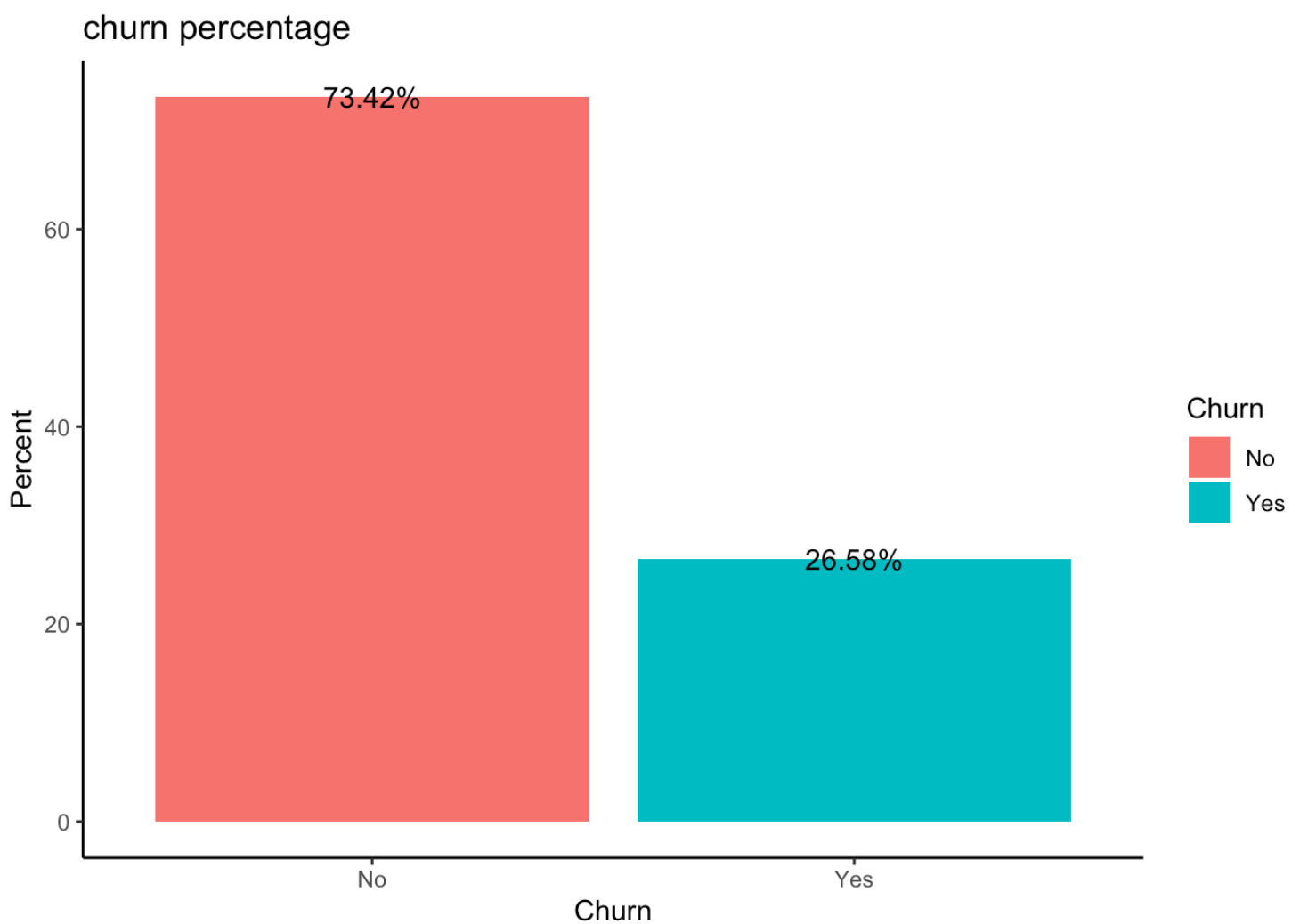
```
head(telechurn02)
```

```
## # A tibble: 6 x 20
##   gender SeniorCitizen Partner Dependents tenure PhoneService MultipleLines
##   <chr>   <fct>           <chr>   <chr>           <dbl> <chr>           <chr>
## 1 Female 0                Yes      No                1 No              No
##   phone ser...
## 2 Male  0                No       No               34 Yes            No
## 3 Male  0                No       No                2 Yes            No
## 4 Male  0                No       No               45 No              No
##   phone ser...
## 5 Female 0                No       No                2 Yes            No
## 6 Female 0                No       No                8 Yes            Yes
## # ... with 13 more variables: InternetService <chr>, OnlineSecurity
##   <chr>,
##   OnlineBackup <chr>, DeviceProtection <chr>, TechSupport <chr>
##   ,
##   StreamingTV <chr>, StreamingMovies <chr>, Contract <chr>,
##   PaperlessBilling <chr>, PaymentMethod <chr>, MonthlyCharges <
##   dbl>,
##   TotalCharges <dbl>, Churn <chr>
```

```
str(telechurn02)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':    7032 obs. of  20 variables:
## $ gender      : chr  "Female" "Male" "Male" "Male" ...
## $ SeniorCitizen : Factor w/ 2 levels "0","1": 1 1 1 1 1 1 1 1 1 1 1 ...
## $ Partner      : chr  "Yes" "No" "No" "No" ...
## $ Dependents   : chr  "No" "No" "No" "No" ...
## $ tenure       : num  1 34 2 45 2 8 22 10 28 62 ...
## $ PhoneService : chr  "No" "Yes" "Yes" "No" ...
## $ MultipleLines : chr  "No phone service" "No" "No" "No phone service" ...
## $ InternetService : chr  "DSL" "DSL" "DSL" "DSL" ...
## $ OnlineSecurity : chr  "No" "Yes" "Yes" "Yes" ...
## $ OnlineBackup  : chr  "Yes" "No" "Yes" "No" ...
## $ DeviceProtection: chr  "No" "Yes" "No" "Yes" ...
## $ TechSupport   : chr  "No" "No" "No" "Yes" ...
## $ StreamingTV   : chr  "No" "No" "No" "No" ...
## $ StreamingMovies : chr  "No" "No" "No" "No" ...
## $ Contract      : chr  "Month-to-month" "One year" "Month-to-month" "One year" ...
## $ PaperlessBilling: chr  "Yes" "No" "Yes" "No" ...
## $ PaymentMethod : chr  "Electronic check" "Mailed check" "Mailed check" "Bank transfer (automatic)" ...
## $ MonthlyCharges : num  29.9 57 53.9 42.3 70.7 ...
## $ TotalCharges   : num  29.9 1889.5 108.2 1840.8 151.7 ...
## $ Churn          : chr  "No" "No" "Yes" "No" ...
```

```
#Explortary
#Tele churn percentage
telechurn02 %>%
  group_by(Churn) %>%
  summarise(Number = n()) %>%
  mutate(Percent = prop.table(Number)*100) %>%
ggplot(aes(x=Churn, y=Percent)) +
  geom_col(aes(fill = Churn)) +
  labs(title = "churn percentage") +
  theme(plot.title = element_text(hjust = 0.5)) +
  geom_text(aes(label = sprintf("%.2f%%", Percent))) +
  scale_colour_brewer(palette = "Set1") + theme_classic()
```



*#Customer Behaviors based on Churn and Not Churn*

*#1 Churn is based on the tenure*

```
telechurn02%>%
```

```
  group_by(tenure, Churn) %>%
```

```
  summarise(Number = n()) %>%
```

```
  ggplot(aes(x=tenure, y=Number)) +
```

```
  geom_line(aes(col = Churn)) +
```

```
  labs(x = "tenure (month)",
```

```
       y = "number of customer",
```

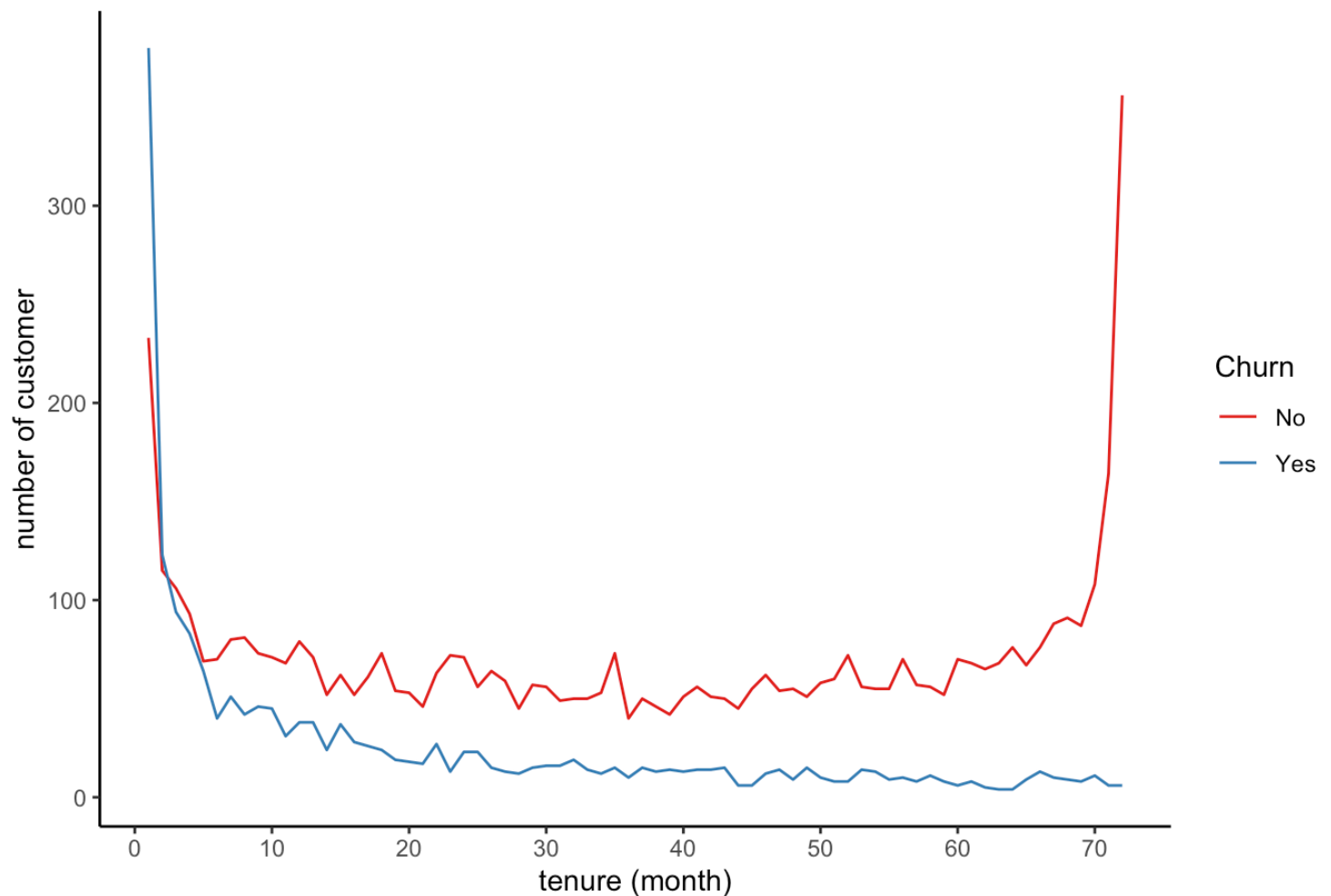
```
       title = "churn based on tenure") +
```

```
  scale_x_continuous(breaks = seq(0, 100, 10)) +
```

```
  scale_colour_brewer(palette = "Set1") + theme_classic()
```



churn based on tenure



*#Correlation between numeric variables*

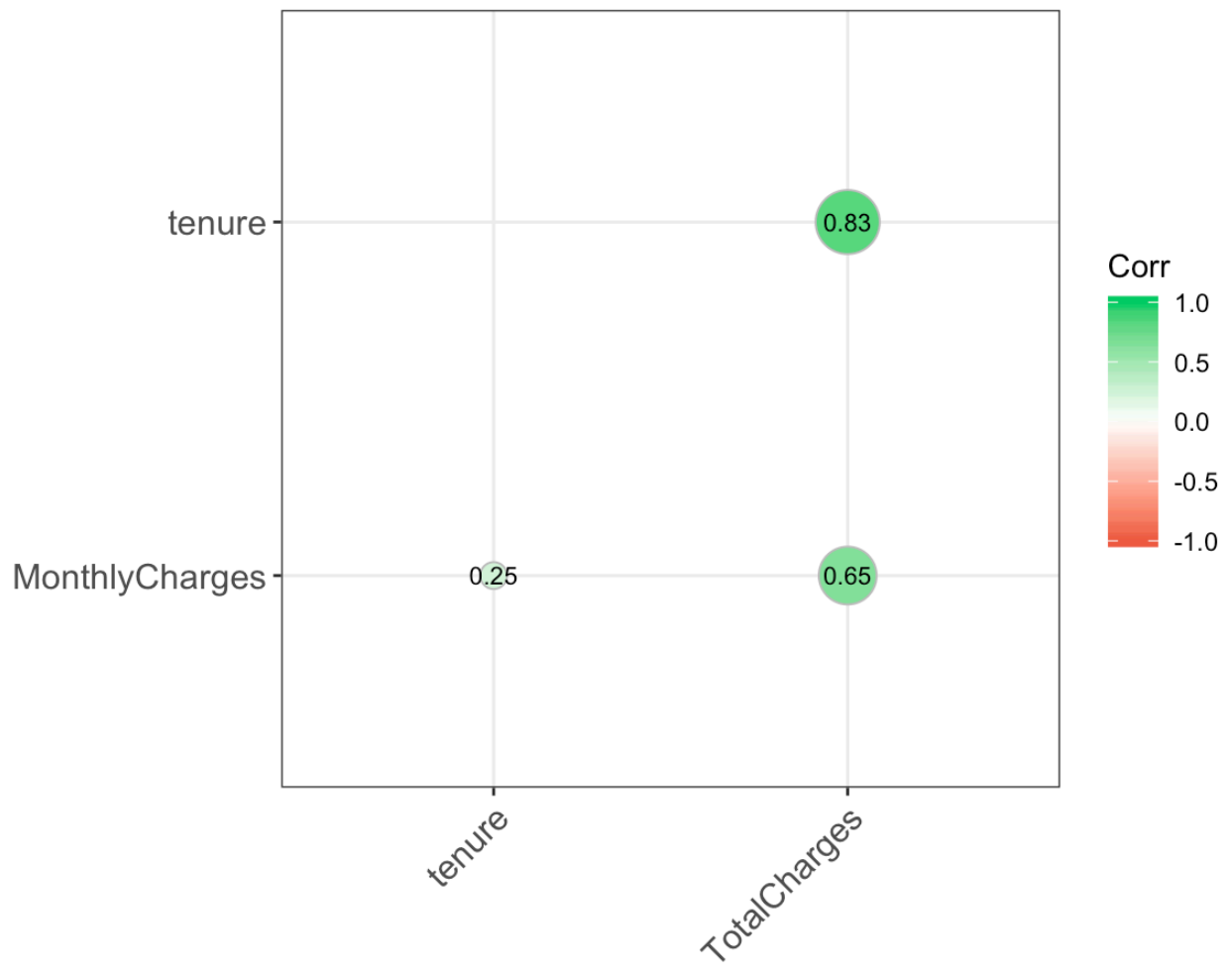
```
library(ggcorrplot)
```

```
numeric_var <- sapply(telechurn02, is.numeric)
```

```
matrix <- cor(telechurn02[,numeric_var])
```

```
ggcorrplot(matrix, hc.order = TRUE,  
            type = "lower",  
            lab = TRUE,  
            lab_size = 3,  
            method="circle",  
            colors = c("tomato2", "white", "springgreen3"),  
            title="Correlogram of numeric variables",  
            ggtheme=theme_bw)
```

Correlogram of numeric variables



*#Density plots for numerical variables*

*#1.Tenure*

```
a1 <-  
  ggplot(telechurn02,aes(x=tenure)) +  
  geom_density(aes(col = Churn))+  
  labs(x = "tenure",  
       y = "density",  
       title = "churn based on tenure") +  
  scale_colour_brewer(palette = "Set1") + theme_classic()+  
  xlim(0,100)
```

*#2.Monthly charges*

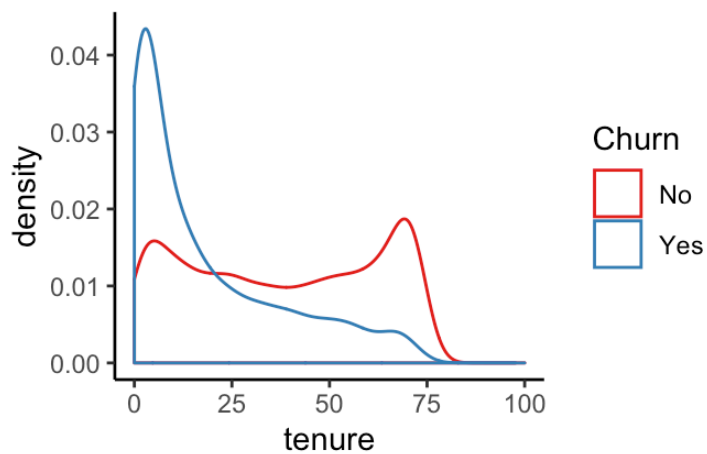
```
a2 <-  
  ggplot(telechurn02,aes(x=MonthlyCharges)) +  
  geom_density(aes(col = Churn)) +  
  labs(x = "monthly charges",  
       y = "density",  
       title = "churn based on monthly charges") +  
  scale_colour_brewer(palette = "Set1") + theme_classic()+  
  xlim(0,150)
```

*#3 Total charges*

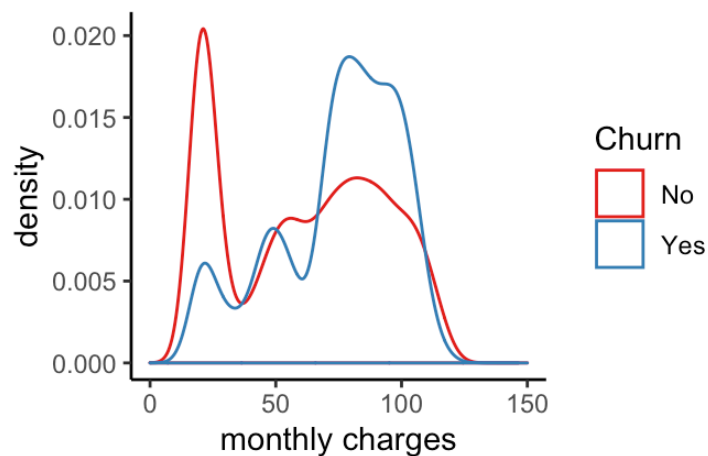
```
a3 <-  
  ggplot(telechurn02,aes(x=TotalCharges)) +  
  geom_density(aes(col = Churn)) +  
  labs(x = "total charges",  
       y = "density",  
       title = "churn based on total charges") +  
  scale_colour_brewer(palette = "Set1") + theme_classic()+  
  xlim(0,10000)
```

```
grid.arrange(a1,a2,a3,nrow=2)
```

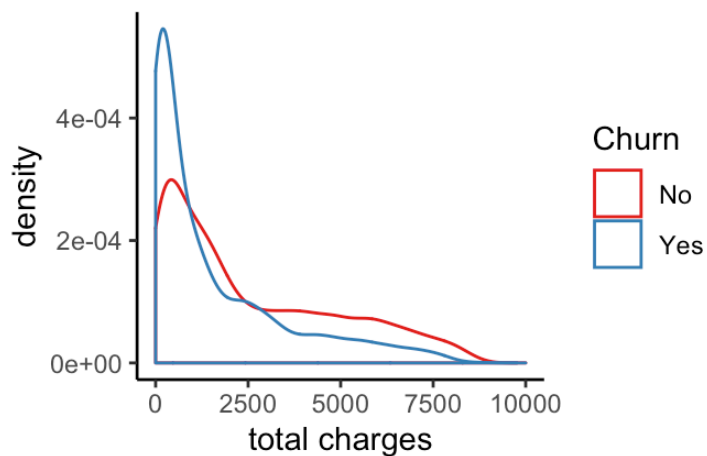
churn based on tenure



churn based on monthly charges



churn based on total charges



*#Different customer chracteristics*

```
p1 <- ggplot(data=telechurn02, aes(x=Churn, group = gender)) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  labs(title = "Gender")+  
  facet_grid(~gender) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

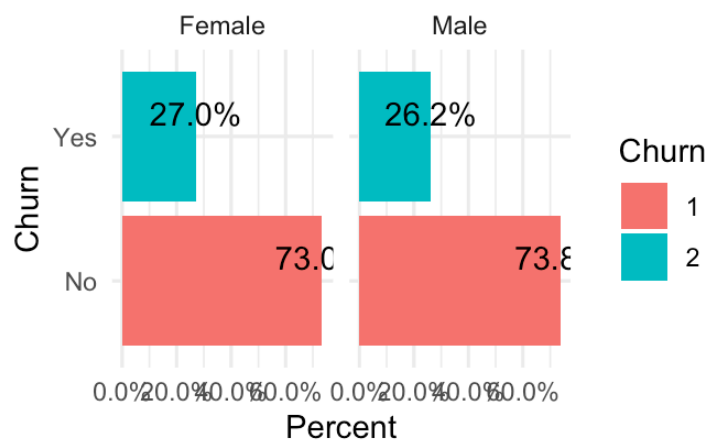
*#1 is the senior citizen; 0 is not the senior citizen*

```
p2 <- ggplot(data=telechurn02, aes(x=Churn, group =SeniorCitizen)) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  labs(title = "Senior Citizen")+  
  facet_grid(~SeniorCitizen) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

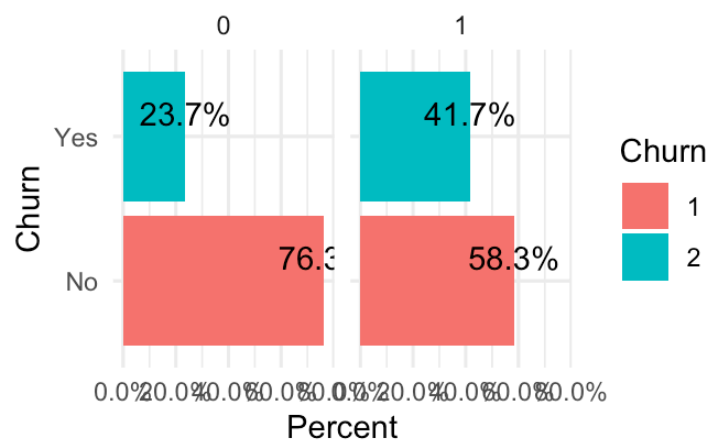
*#Yes means that the customer has partner;*

```
p3<- ggplot(data=telechurn02, aes(x=Churn, group =Partner)) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  labs(title = "Partner")+  
  facet_grid(~Partner) +  
  scale_y_continuous(labels = scales::percent)+ coord_flip() + theme_minimal()  
grid.arrange(p1,p2,p3, nrow=2)
```

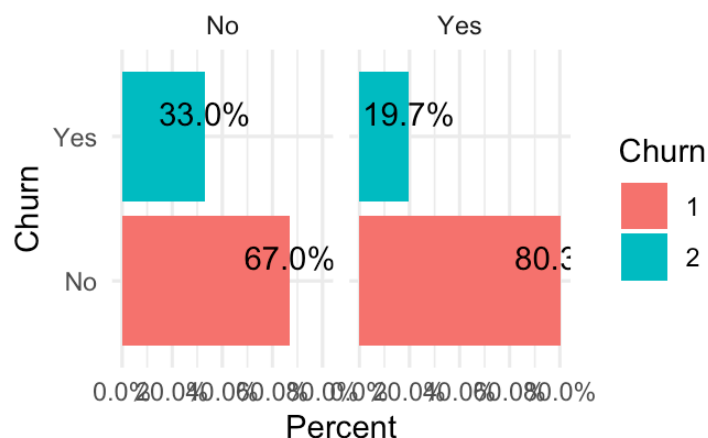
## Gender



## Senior Citizen



## Partner



### *#Paperless Billing*

```
p1 <- ggplot(data=telechurn02, aes(x=Churn, group = PaperlessBilling)) +
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +
  geom_text(aes( label = scales::percent(..prop..),
    y= ..prop.. ), stat= "count", vjust = -.5) +
  labs(y = "Percent", fill="Churn") +
  facet_grid(~PaperlessBilling) +
  scale_y_continuous(labels = scales::percent)+
  coord_flip() + theme_minimal()
```

### *#Payment Methods*

```
p2 <- ggplot(data=telechurn02, aes(x=Churn, group = PaymentMethod)) +
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +
  geom_text(aes( label = scales::percent(..prop..),
```

```

        y= ..prop.. ), stat= "count", vjust = -.5) +
labs(y = "Percent", fill="Churn") +
facet_grid(~PaymentMethod) +
scale_y_continuous(labels = scales::percent)+
coord_flip() + theme_minimal()

```

### *#Contract Methods*

```

p3 <- ggplot(data=telechurn02, aes(x=Churn, group = Contract)) +
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +
  geom_text(aes( label = scales::percent(..prop..),
                y= ..prop.. ), stat= "count", vjust = -.5) +
labs(y = "Percent", fill="Churn") +
facet_grid(~Contract) +
scale_y_continuous(labels = scales::percent)+
coord_flip() + theme_minimal()

```

### *#internet service*

```

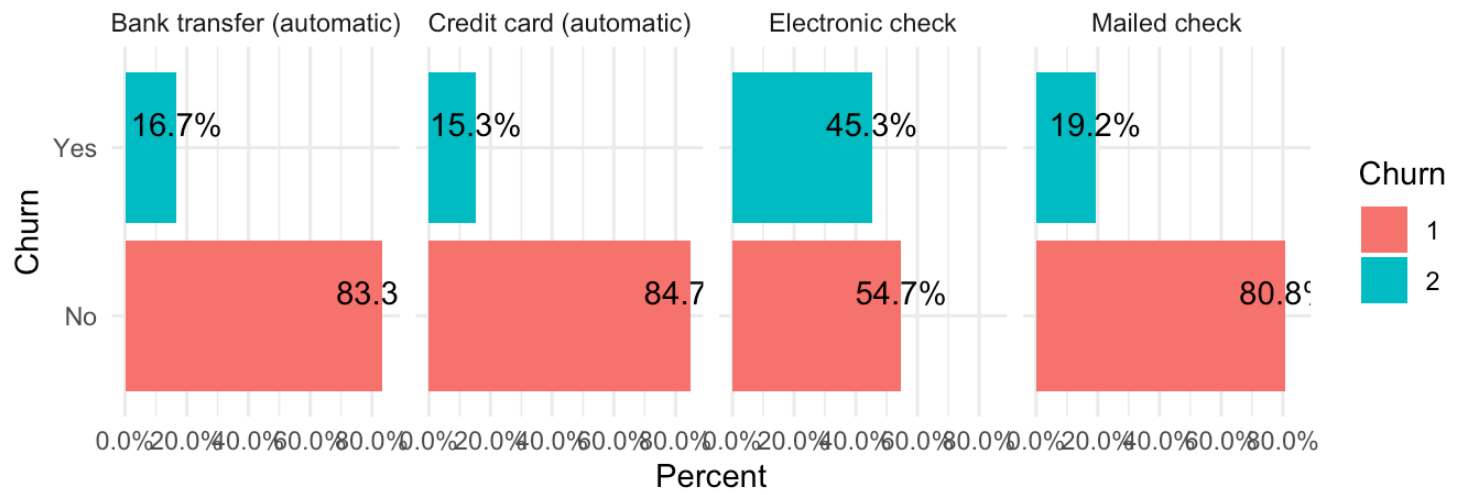
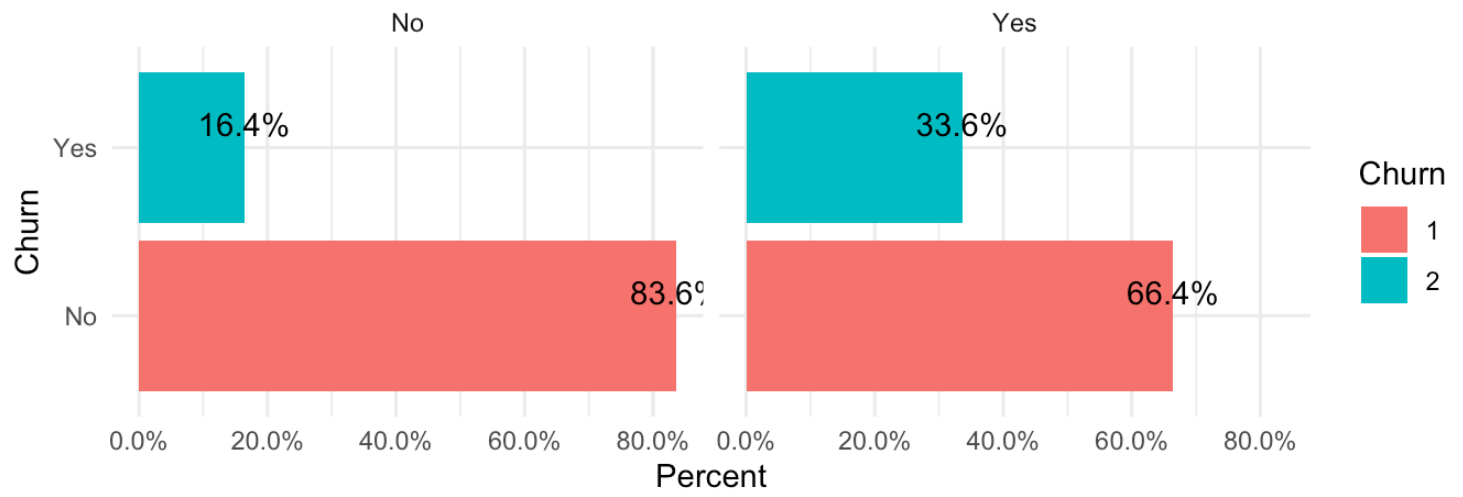
p4 <- ggplot(data=telechurn02, aes(x=Churn, group = InternetService)
) +
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +
  geom_text(aes( label = scales::percent(..prop..),
                y= ..prop.. ), stat= "count", vjust = -.5) +
labs(y = "Percent", fill="Churn") +
facet_grid(~InternetService) +
scale_y_continuous(labels = scales::percent)+
coord_flip() + theme_minimal()

```

```

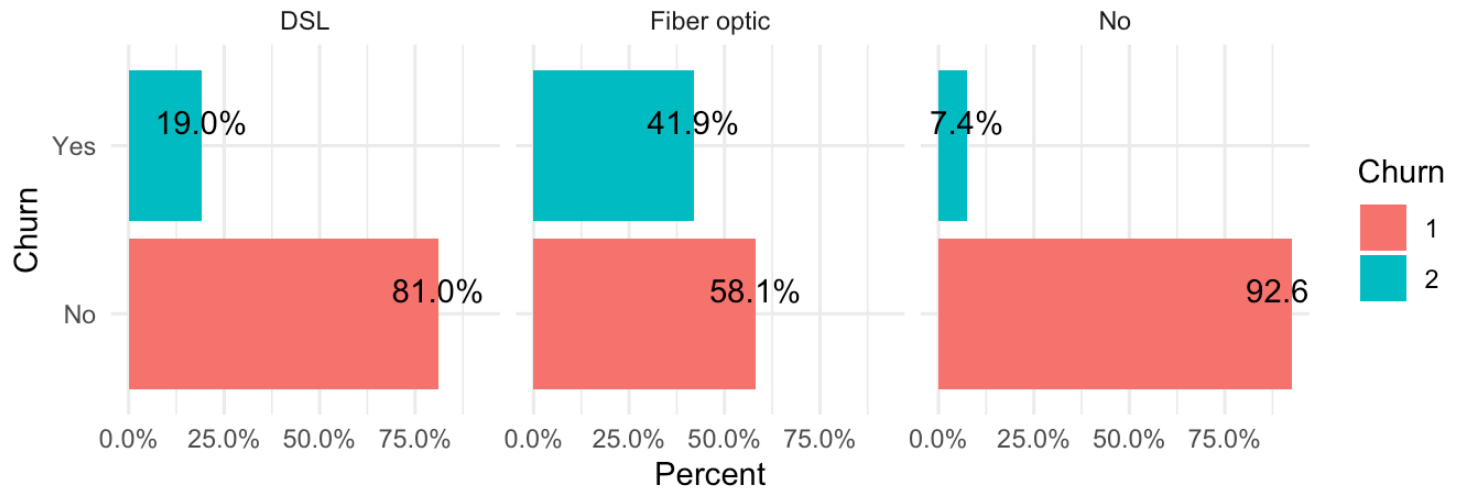
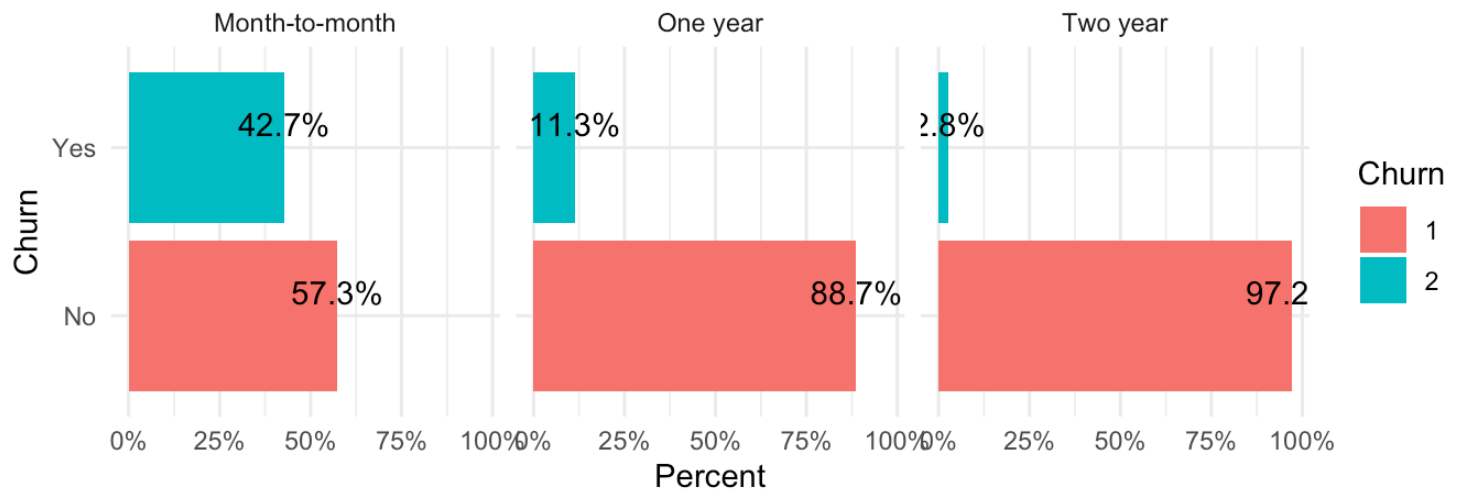
grid.arrange(p1,p2)

```



```
grid.arrange(p3,p4)
```





```
#churn based on the service
```

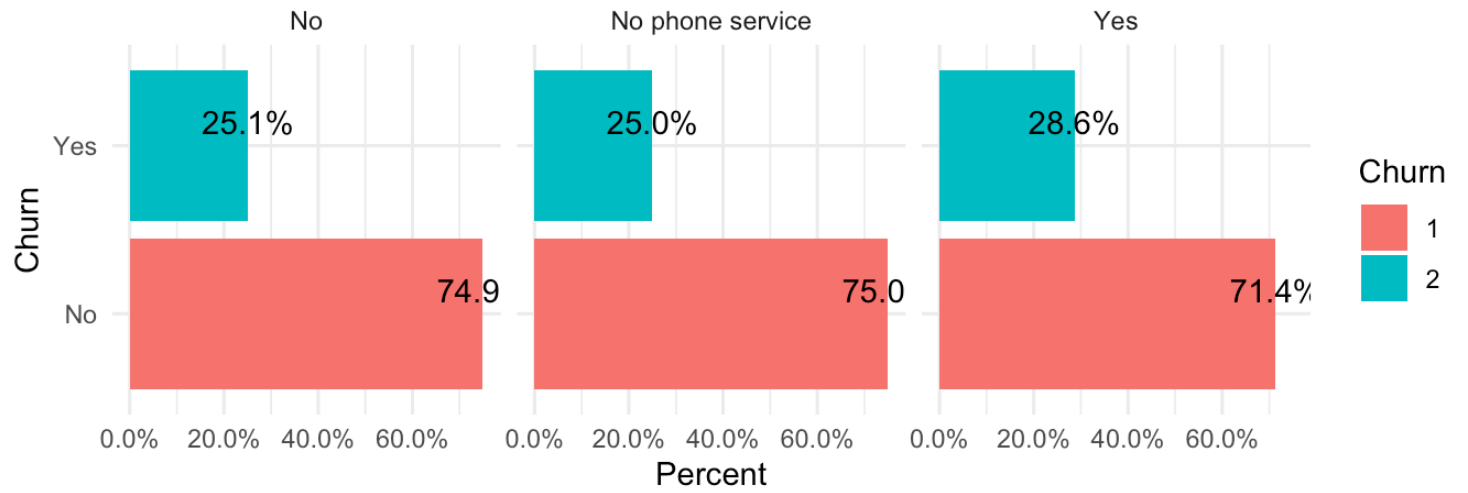
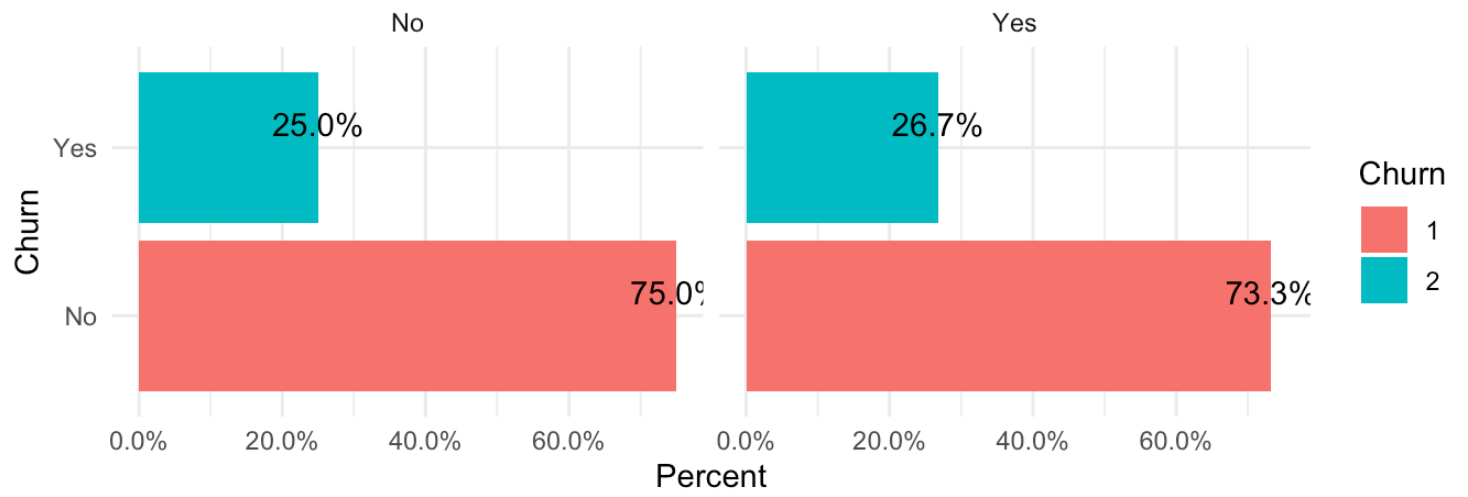
```
#phone service
```

```
p1 <- ggplot(data=telechurn02, aes(x=Churn, group = PhoneService)) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  facet_grid(~PhoneService) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

```
#Multiple lines
```

```
p2 <- ggplot(data=telechurn02, aes(x=Churn, group = MultipleLines)) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  facet_grid(~MultipleLines) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

```
grid.arrange(p1,p2)
```



*#online security*

```
p4 <- ggplot(data=telechurn02, aes(x=Churn, group = OnlineSecurity))  
+  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  facet_grid(~OnlineSecurity) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

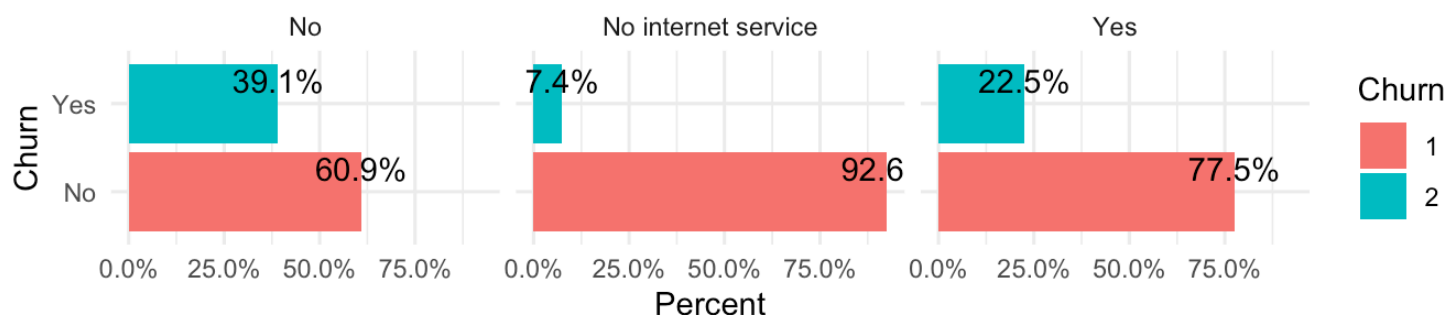
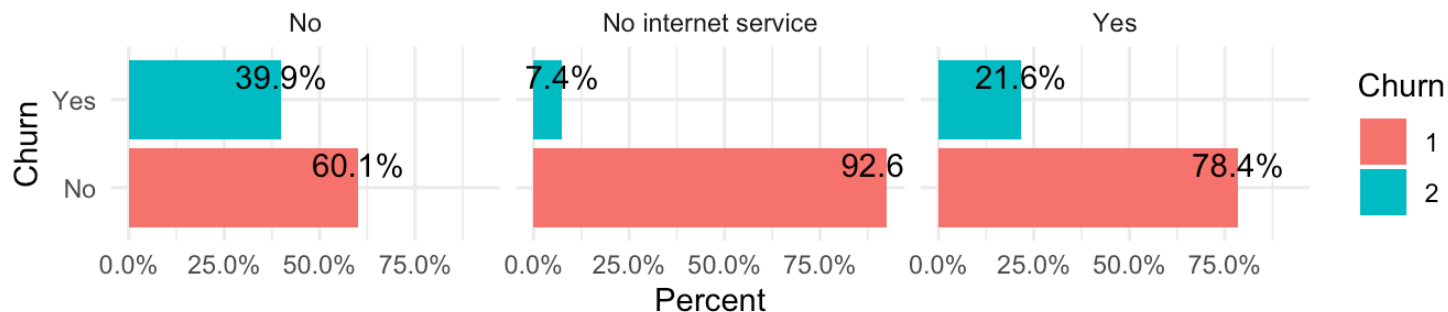
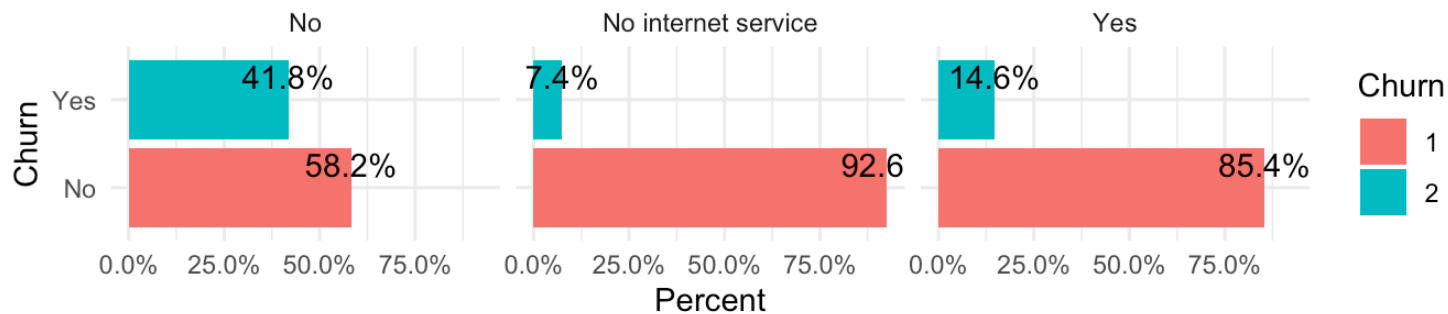
*#online backup*

```
p5 <- ggplot(data=telechurn02, aes(x=Churn, group = OnlineBackup)) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  facet_grid(~OnlineBackup) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

*#device protection*

```
p6 <- ggplot(data=telechurn02, aes(x=Churn, group = DeviceProtection  
) ) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  facet_grid(~DeviceProtection) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

```
grid.arrange(p4,p5,p6)
```



```
#TechSupport
```

```
p7 <- ggplot(data=telechurn02, aes(x=Churn, group = TechSupport)) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  facet_grid(~TechSupport) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

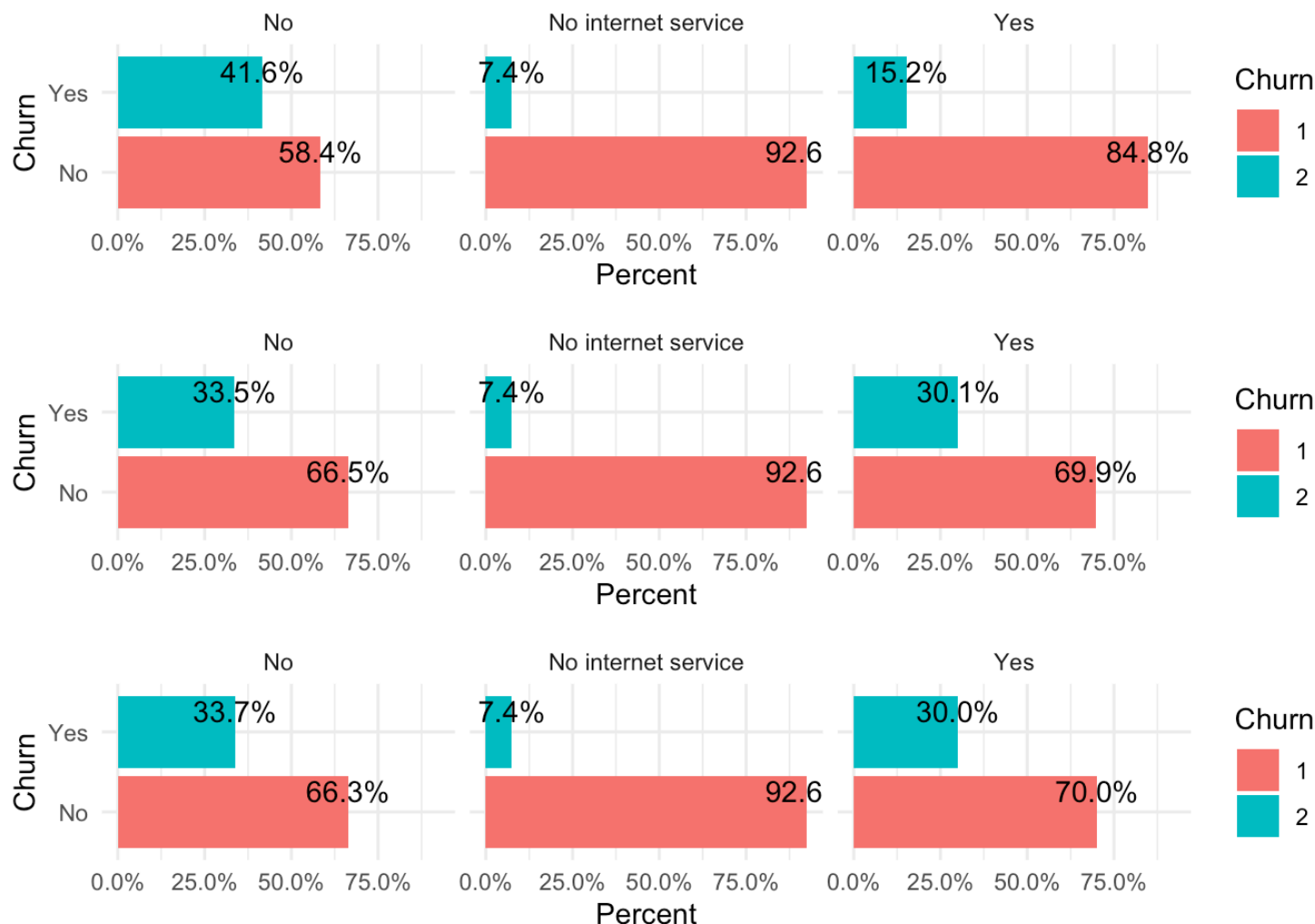
```
#StreamingTV
```

```
p8 <- ggplot(data=telechurn02, aes(x=Churn, group = StreamingTV)) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  facet_grid(~StreamingTV) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

```
#StreamingMovies
```

```
p9 <- ggplot(data=telechurn02, aes(x=Churn, group = StreamingMovies)  
 ) +  
  geom_bar(aes(y = ..prop.., fill = factor(..x..)), stat="count") +  
  geom_text(aes( label = scales::percent(..prop..),  
                y= ..prop.. ), stat= "count", vjust = -.5) +  
  labs(y = "Percent", fill="Churn") +  
  facet_grid(~StreamingMovies) +  
  scale_y_continuous(labels = scales::percent)+  
  coord_flip() + theme_minimal()
```

```
grid.arrange(p7,p8,p9)
```



```
#churn prediction
table(telechurn02$Churn)
```

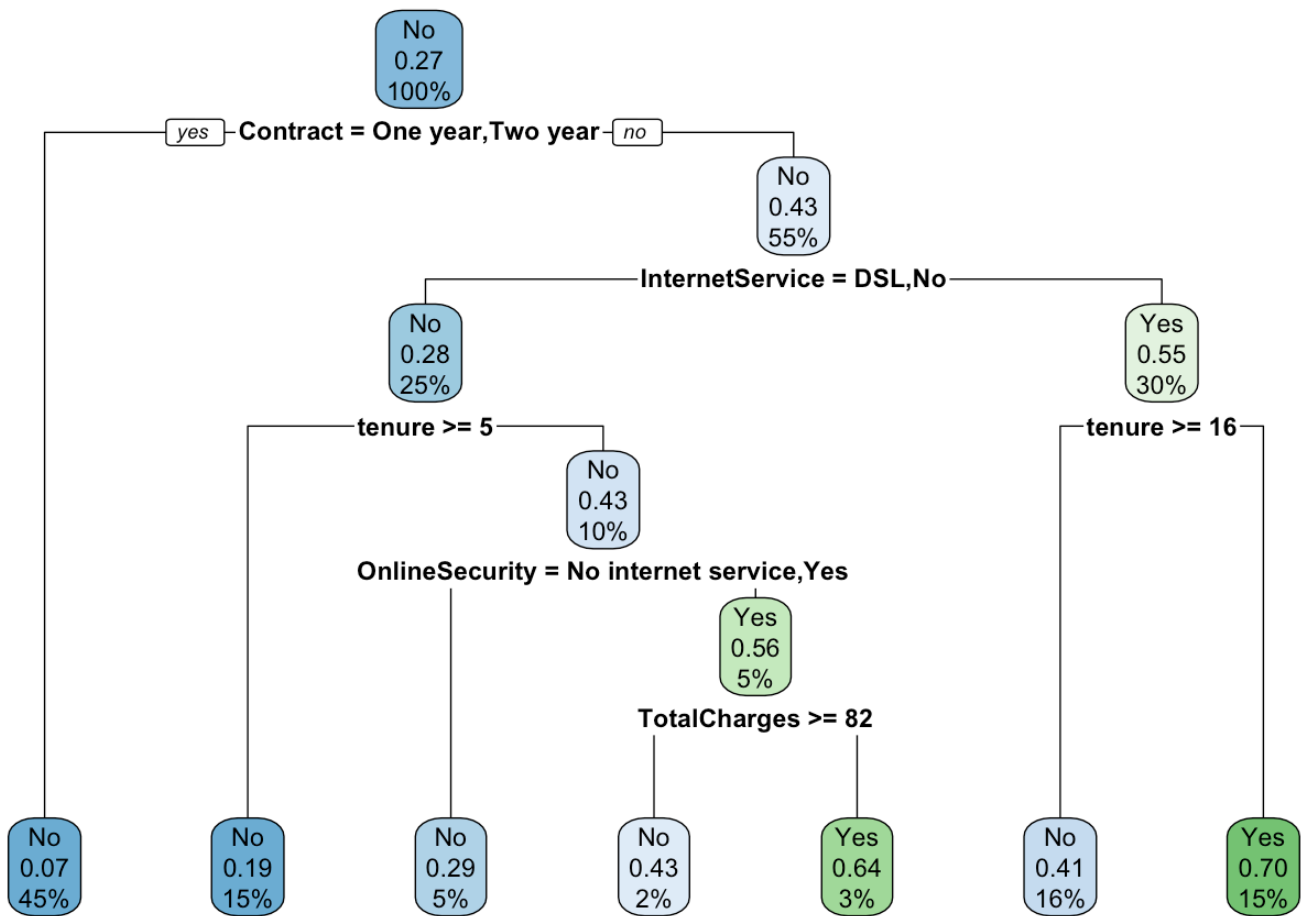
```
##
##      No   Yes
## 5163 1869
```

```
#chr variables to factor variables
index01 <- sapply(telechurn02,is.character)
telechurn02[index01] <- lapply(telechurn02[index01],as.factor)
str(telechurn02)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':    7032 obs. of  20 variables:
## $ gender      : Factor w/ 2 levels "Female","Male": 1 2 2 2
1 1 2 1 1 2 ...
## $ SeniorCitizen : Factor w/ 2 levels "0","1": 1 1 1 1 1 1 1 1
1 1 ...
## $ Partner      : Factor w/ 2 levels "No","Yes": 2 1 1 1 1 1 1
1 2 1 ...
## $ Dependents   : Factor w/ 2 levels "No","Yes": 1 1 1 1 1 1 2
1 1 2 ...
## $ tenure       : num  1 34 2 45 2 8 22 10 28 62 ...
## $ PhoneService : Factor w/ 2 levels "No","Yes": 1 2 2 1 2 2 2
1 2 2 ...
## $ MultipleLines : Factor w/ 3 levels "No","No phone service",.
.: 2 1 1 2 1 3 3 2 3 1 ...
## $ InternetService : Factor w/ 3 levels "DSL","Fiber optic",...: 1
1 1 1 2 2 2 1 2 1 ...
## $ OnlineSecurity : Factor w/ 3 levels "No","No internet service
",...: 1 3 3 3 1 1 1 3 1 3 ...
## $ OnlineBackup   : Factor w/ 3 levels "No","No internet service
",...: 3 1 3 1 1 1 3 1 1 3 ...
## $ DeviceProtection: Factor w/ 3 levels "No","No internet service
",...: 1 3 1 3 1 3 1 1 3 1 ...
## $ TechSupport    : Factor w/ 3 levels "No","No internet service
",...: 1 1 1 3 1 1 1 1 3 1 ...
## $ StreamingTV     : Factor w/ 3 levels "No","No internet service
",...: 1 1 1 1 1 3 3 1 3 1 ...
## $ StreamingMovies : Factor w/ 3 levels "No","No internet service
",...: 1 1 1 1 1 3 1 1 3 1 ...
## $ Contract       : Factor w/ 3 levels "Month-to-month",...: 1 2
1 2 1 1 1 1 1 2 ...
## $ PaperlessBilling: Factor w/ 2 levels "No","Yes": 2 1 2 1 2 2 2
1 2 1 ...
## $ PaymentMethod   : Factor w/ 4 levels "Bank transfer (automatic
)",...: 3 4 4 1 3 3 2 4 3 1 ...
## $ MonthlyCharges : num  29.9 57 53.9 42.3 70.7 ...
## $ TotalCharges    : num  29.9 1889.5 108.2 1840.8 151.7 ...
## $ Churn           : Factor w/ 2 levels "No","Yes": 1 1 2 1 2 2 1
1 2 1 ...
```







```

library(caret)
dt_test <- predict(object = dtmodel,
                    newdata = teletest,
                    type = "class")

#confusion matrix for prediction
dt1 <- table(Predicted = dt_test, Actual = teletest$Churn)
dt1

```

```

##           Actual
## Predicted  No  Yes
##         No  944 210
##         Yes   88 163

```

```
#Evaluation
accuracy = sum(944+163)/ length(teletest$Churn)
precision = dt1[1,1]/sum(dt1[,1])
recall = dt1[1,1]/sum(dt1[1,])
f = 2 * (precision * recall) / (precision + recall)

cat(paste("Accuracy:\t", format(accuracy, digits=2), "\n",sep=" "))
```

```
## Accuracy:      0.79
```

```
cat(paste("Precision:\t", format(precision, digits=2), "\n",sep=" ")
)
```

```
## Precision:     0.91
```

```
cat(paste("Recall:\t\t", format(recall, digits=2), "\n",sep=" "))
```

```
## Recall:        0.82
```

```
cat(paste("F-measure:\t", format(f, digits=2), "\n",sep=" "))
```

```
## F-measure:     0.86
```

```
#Visualize ROC curve for Decision Tree Model
library(ROCR)
```

```
## Loading required package: gplots
```

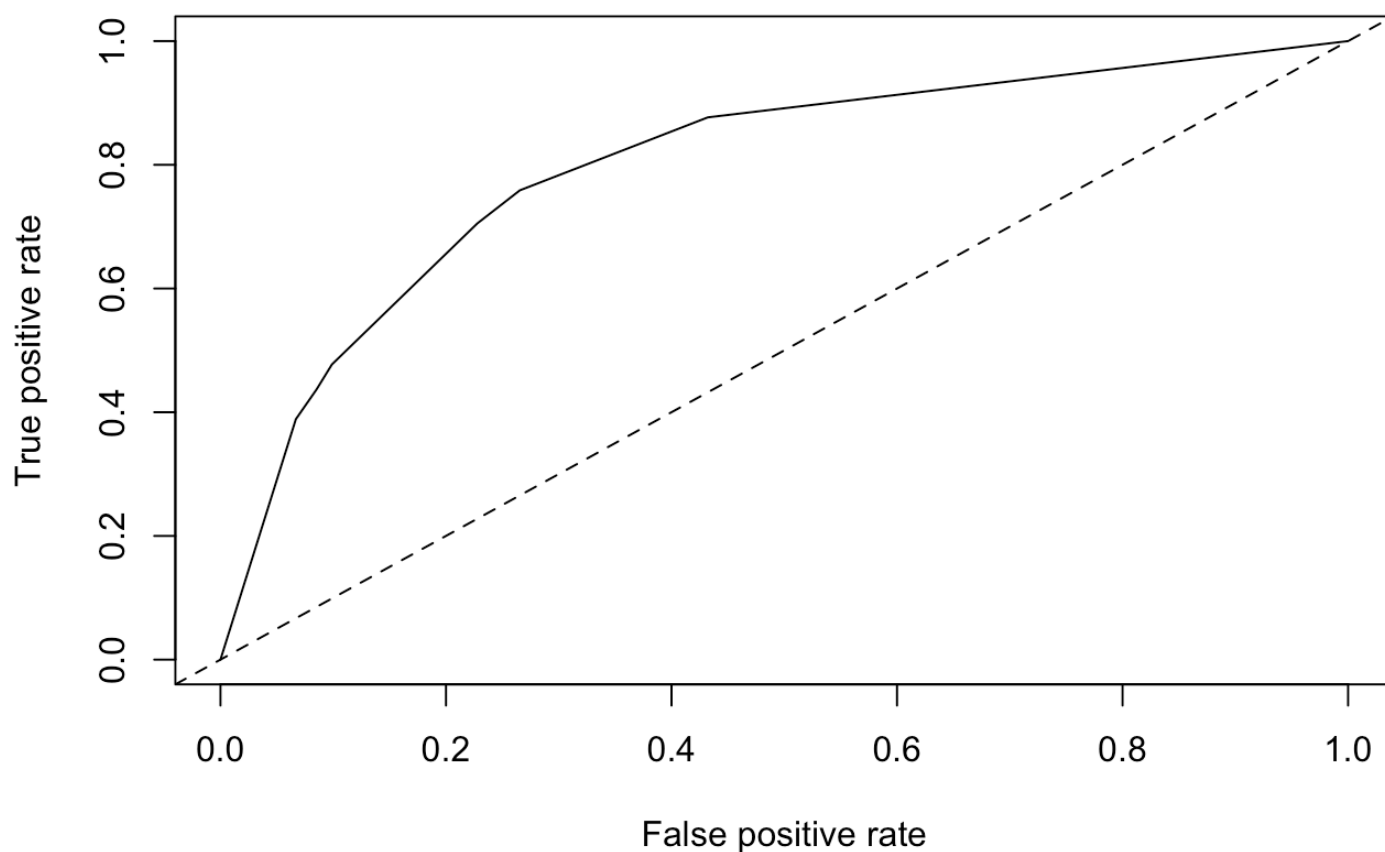
```
##
## Attaching package: 'gplots'
```

```
## The following object is masked from 'package:stats':
```

```
##
```

```
## lowess
```

```
Pred.cart = predict(dtmmodel, newdata =teletest, type = "prob")[,2]  
Pred2 = prediction(Pred.cart, teletest$Churn)  
plot(performance(Pred2, "tpr", "fpr"))  
abline(0, 1, lty = 2)
```



```
#Get the AUC
```

```
auc <- performance(Pred2, measure = "auc")
```

```
auc@y.values[[1]]
```

```
## [1] 0.7997407
```

#AUC=0.799