1 Female No No phone service Yes ## 2 Male 0 No No 34 No ## 3 Male No 2 Yes No No No phone service 0 No No 45 ## 4 Male ## 5 Female ## 6 Female 0 No No Yes Yes InternetService OnlineSecurity OnlineBackup DeviceProtection TechSupport ## 1 DSL No Yes No DSL ## 2 Yes No ## 3 Yes DSL Yes No No ## 4 DSL Yes Yes Yes ## 5 Fiber optic No No No No Fiber optic No No No StreamingTV StreamingMovies Contract PaperlessBilling ## ## 1 No Month-to-month ## 2 No No One year No ## 3 No Month-to-month Yes ## 4 No No No One year ## 5 No No Month-to-month Yes ## 6 Yes Month-to-month Yes ## PaymentMethod MonthlyCharges TotalCharges Churn ## 1 Electronic check 29.85 29.85 No 1889.50 ## 2 Mailed check 56.95 No ## 3 Mailed check 53.85 108.15 Yes 42.30 ## 4 Bank transfer (automatic) 1840.75 No ## 5 Electronic check 70.70 151.65 Yes ## 6 Electronic check 99.65 820.50 Yes #missing value imputation df\$TotalCharges <- as.numeric(df\$TotalCharges)</pre> miss = which(is.na(df\$TotalCharges) == TRUE) df\$TotalCharges[miss] <- median(df\$TotalCharges, na.rm = TRUE)</pre> str(df) ## 'data.frame': 7043 obs. of 20 variables: : chr "Female" "Male" "Male" "... ## \$ SeniorCitizen : int 0000000000... ## \$ Partner : chr "Yes" "No" "No" "No" ... ## \$ Dependents : chr "No" "No" "No" "No" ... \$ tenure : int 1 34 2 45 2 8 22 10 28 62 ... : chr "No" "Yes" "Yes" "No" ... \$ PhoneService ## \$ MultipleLines : chr "No phone service" "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" ... : chr "No" "No" "No" "No" ... ## \$ StreamingTV ## \$ StreamingMovies : chr "No" "No" "No" "No" ... ## \$ Contract : chr "Month-to-month" "One year" "Month-to-month" "One year" ... ## \$ PaperlessBilling: chr "Yes" "No" "Yes" "No" ... : chr "Electronic check" "Mailed check" "Mailed check" "Bank transfer (automatic)" ... ## \$ PaymentMethod ## \$ MonthlyCharges : num 29.9 57 53.9 42.3 70.7 ... ## \$ TotalCharges : num 29.9 1889.5 108.2 1840.8 151.7 ... : chr "No" "No" "Yes" "No" ... ## \$ Churn #No Service to No for(i in (which(colnames(df) == 'OnlineSecurity') : which(colnames(df) == 'StreamingMovies'))){ df[i] <- as.factor(ifelse(df[i] != 'Yes', 'No', 'Yes'))</pre> df\$InternetService <- as.factor(ifelse(df\$InternetService != 'No', 'Yes', 'No'))</pre> df\$MultipleLines <- as.factor(ifelse(df\$MultipleLines != 'Yes', 'No', 'Yes'))</pre> df\$SeniorCitizen <- as.factor(df\$SeniorCitizen)</pre> for(i in 1:ncol(df)){ if(class(df[,i]) == 'character'){ df[,i] <- as.factor(df[,i])</pre> } str(df) ## 'data.frame': 7043 obs. of 20 variables: : Factor w/ 2 levels "Female", "Male": 1 2 2 2 1 1 2 1 1 2 ... ## \$ gender ## \$ SeniorCitizen : Factor w/ 2 levels "0", "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 : int 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/ 2 levels "No", "Yes": 1 1 1 1 1 2 2 1 2 1 ... ## \$ InternetService : Factor w/ 2 levels "No", "Yes": 2 2 2 2 2 2 2 2 2 ... ## \$ OnlineSecurity : Factor w/ 2 levels "No", "Yes": 1 2 2 2 1 1 1 2 1 2 ... ## \$ OnlineBackup : Factor w/ 2 levels "No", "Yes": 2 1 2 1 1 1 2 1 1 2 ... ## \$ DeviceProtection: Factor w/ 2 levels "No", "Yes": 1 2 1 2 1 2 1 1 2 1 ... ## \$ TechSupport : Factor w/ 2 levels "No", "Yes": 1 1 1 2 1 1 1 1 2 1 ... ## \$ StreamingTV : Factor w/ 2 levels "No", "Yes": 1 1 1 1 1 2 2 1 2 1 ... ## \$ StreamingMovies : Factor w/ 2 levels "No", "Yes": 1 1 1 1 1 2 1 1 2 1 ... : 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 ... #Correlation between numeric variables cr < -cor(df[, c(5, 18, 19)])corrplot(cr, method="circle") MonthlyCharges TotalCharges tenure 0.6 0.4 0.2 MonthlyCharges -0.2 -0.4 -0.6 TotalCharges -0.8 #EDA $p1 \leftarrow ggplot(df, aes(x = Churn, fill = Churn)) + facet_grid(\sim gender) + geom_bar() + ggtitle("Churn - Gender") + them$ p2 <- ggplot(df, aes(x = Churn, fill = Churn)) +facet_grid(~SeniorCitizen)+ geom_bar() + ggtitle("Churn - SeniorC itizen") + theme_bw() $p3 < -ggplot(df, aes(x = Churn, fill = Churn)) + facet_grid(~Dependents) + geom_bar() + ggtitle("Churn - Dependent)$ s") + theme_bw() p4 <- ggplot(df, aes(x = Churn, fill = Churn)) +facet_grid(~Partner)+ geom_bar() + ggtitle("Churn - Partner") + t p5 <- ggplot(df, aes(x = Churn, fill = Churn)) +facet_grid(~PhoneService)+ geom_bar() + ggtitle("Churn - PhoneSer vice")+ theme_bw() $p6 < -ggplot(df, aes(x = Churn, fill = Churn)) + facet_grid(~InternetService) + geom_bar() + ggtitle("Churn - InternetService)) + geom_bar() + ggtitle("Churn - InternetService))$ netService") + theme_bw() p7 <- ggplot(df, aes(x = Churn, fill = Churn)) +facet_grid(~PaperlessBilling)+ geom_bar() + ggtitle("Churn - Pape rlessBilling") + theme_bw() p8 <- ggplot(df, aes(x = Churn, fill = Churn)) +facet_grid(~PaymentMethod)+ geom_bar() + ggtitle("Churn - Payment Method") + theme_bw() ggpubr::ggarrange(p1, p2, p3, p4, p5, p6, p7, p8, nrow = 3, ncol = 3)Churn - Gender Churn - SeniorCitizen Churn - Dependents Male Female Yes Churn Churn Churn 3000 -2000 -1000 -4000 2000 -1000 -3000 -2000 -1000 -No Yes NoYes NoYes NoYes NoYes NoYes NoYes Churn Churn Churn Churn - PhoneService Churn - InternetServic Churn - Partner No Yes Yes Yes Churn Churn Churn 4000 3000 -2000 -1000 -4000 -3000 -2000 -1000 -2000 t No Yes Yes NoYes NoYes NoYes NoYes NoYes NoYes Churn Churn Churn Churn - PaperlessBilli Churn - Payment Method Yes Churn Churn conut 1000 500 2000 to 1000 t No Yes NoYes NoYes N/ces/N/ces/N/ces Churn Churn #dummification attach(df) to_dummy <- data.frame(Contract,PaymentMethod)</pre> dmy <- dummyVars(" ~ .", data = to_dummy)</pre> df2 <- data.frame(predict(dmy, newdata = to_dummy))</pre> df2 <- df2[, !(colnames(df2) %in% c("Contract.Month.to.month", "PaymentMethod.Bank.transfer..automatic."))] df <- df[,!(colnames(df) %in% c("Contract", "PaymentMethod", "TotalCharges"))]</pre> df <- cbind(df, df2)</pre> gender SeniorCitizen Partner Dependents tenure PhoneService MultipleLines 0 ## 1 Female 1 No ## 2 Male Yes No ## 3 0 No Yes No Male Male No 0 Yes ## 5 Female No No No ## 6 Female 0 No Yes InternetService OnlineSecurity OnlineBackup DeviceProtection TechSupport ## 1 ## 2 Yes Yes No Yes No ## 3 Yes Yes Yes No No ## 4 Yes Yes No Yes Yes Yes No No ## 6 Yes No No Yes No ## StreamingTV StreamingMovies PaperlessBilling MonthlyCharges Churn ## 1 29.85 No No Yes ## 2 No 56.95 No ## 3 No No Yes 53.85 Yes ## 4 No No No 42.30 No ## 5 No No Yes 70.70 Yes Yes Yes 99.65 Contract.One.year Contract.Two.year PaymentMethod.Credit.card..automatic. ## ## 1 ## 2 0 0 ## 4 1 ## 5 ## 6 PaymentMethod.Electronic.check PaymentMethod.Mailed.check ## 1 ## 2 ## 3 0 1 ## 4 ## 5 0 attach(df) ## The following objects are masked from df (pos = 3): Churn, Dependents, DeviceProtection, gender, InternetService, MonthlyCharges, MultipleLines, OnlineBackup, OnlineSecurity, ## PaperlessBilling, Partner, PhoneService, SeniorCitizen, StreamingMovies, StreamingTV, TechSupport, tenure dim(df) ## [1] 7043 22 #Feature Selection regfit.full=regsubsets(Churn~., data=df, nvmax=21) reg.summary=summary(regfit.full) names(reg.summary) "adjr2" "cp" ## [1] "which" "rsq" "rss" "bic" "outmat" "obj" which.min(reg.summary\$bic) ## [1] 12 plot(reg.summary\$bic,xlab="No. of Variables",ylab=expression(paste("BIC")),type="1") points(12, reg.summary\$bic[12], col="red", cex=2, pch=20) -1400 BIC -2200 5 10 15 20 No. of Variables names(coef(regfit.full, 12))[-1] ## [1] "SeniorCitizen1" "tenure" [3] "PhoneServiceYes" "OnlineSecurityYes" "DeviceProtectionYes" ## [5] "OnlineBackupYes" "PaperlessBillingYes" ## [7] "TechSupportYes" ## [9] "MonthlyCharges" "Contract.One.year" ## [11] "Contract.Two.year" "PaymentMethod.Electronic.check" #Final Dataset data=df[,-c(1,3,4,7,8,13,14,20,22)] dim(data) ## [1] 7043 13 #train-test split: index1=sample(1:nrow(data),floor(0.7*nrow(data))) train=data[index1,] remaining=data[-index1,] index2=sample(1:nrow(remaining), floor(2/3*nrow(remaining))) crossval=remaining[index2,] test=remaining[-index2,] actual_churn=crossval\$Churn #logistic regression logistic.fit=glm(Churn~., data=train, family="binomial") logistic.predict=rep("No", nrow(crossval)) predicted_prob=predict(logistic.fit,newdata=crossval,type="response") logistic.predict[predicted_prob>0.5]="Yes" table(logistic.predict,actual_churn) ## actual_churn ## logistic.predict No Yes No 938 169 Yes 100 201 mean(logistic.predict==actual_churn) ## [1] 0.8089489 #lda fit lda.fit=lda(Churn~.,data=train) lda.predict=predict(lda.fit,crossval)\$class table(lda.predict,actual_churn) actual_churn ## lda.predict No Yes No 930 169 Yes 108 201 mean(lda.predict==actual_churn) ## [1] 0.803267 **#Classification Tree** tree.fit=tree(Churn~., train, method="class") summary(tree.fit) ## Classification tree: ## tree(formula = Churn ~ ., data = train, method = "class") ## Variables actually used in tree construction: ## [1] "Contract.Two.year" "Contract.One.year" ## [3] "PaymentMethod.Electronic.check" "tenure" ## [5] "MonthlyCharges" ## Number of terminal nodes: 8 ## Residual mean deviance: 0.8856 = 4359 / 4922 ## Misclassification error rate: 0.2152 = 1061 / 4930 plot(tree.fit) text(tree.fit,pretty=0,cex=1) Contract.Two.year < 0.5 Contract.One.year < 0.5 No PaymentMethod.Electronic.check < 0.5 MonthlyCharges < 81.95 tenure < 7.5 MonthlyCharges < 29.825 tenure < 10.5 No No No Yes No No Yes tree.predict=predict(tree.fit,crossval,type="class") table(predicted_churn=tree.predict,actual_churn) actual_churn ## predicted_churn No Yes No 920 226 Yes 118 144 mean(tree.predict==actual_churn) ## [1] 0.7556818 #Random Forest rf.fit=randomForest(Churn~., data=train, ntree=200, mtry=4) rf.predict=predict(rf.fit,crossval) table(predicted_churn=rf.predict,actual_churn) actual_churn ## predicted_churn No Yes No 917 178 Yes 121 192 mean(rf.predict==actual_churn) ## [1] 0.787642 misclassification_rate_logistic=(mean(logistic.predict!=actual_churn))*100 misclassification_rate_lda=(mean(lda.predict!=actual_churn))*100 misclassification_rate_tree=(mean(tree.predict!=actual_churn))*100 misclassification_rate_forest=(mean(rf.predict!=actual_churn))*100 paste("Misclassification Error Rate for Logistic Regression is", misclassification_rate_logistic, "%") ## [1] "Misclassification Error Rate for Logistic Regression is 19.1051136363636 %" paste("Misclassification Error Rate for Linear Discriminant Analysis is is", misclassification_rate_lda, "%") ## [1] "Misclassification Error Rate for Linear Discriminant Analysis is is 19.6732954545455 %" paste("Misclassification Error Rate for Decision Tree is", misclassification_rate_tree, "%") ## [1] "Misclassification Error Rate for Decision Tree is 24.4318181818182 %" paste("Misclassification Error Rate for Random Forest is", misclassification_rate_forest, "%") ## [1] "Misclassification Error Rate for Random Forest is 21.2357954545455 %" #Choice is Logistic Regression #Fit on test dataset predicted_prob=predict(logistic.fit, newdata=test, type="response") logistic.predict.test=rep("No",nrow(test)) logistic.predict.test[predicted_prob>0.5]="Yes" actual.churn.test=test\$Churn table(logistic.predict.test,actual.churn.test) actual.churn.test ## logistic.predict.test No Yes No 474 72 Yes 56 103 misclassification.final=mean(logistic.predict.test!=actual.churn.test)*100 paste("Misclassification Error Rate for final model is", misclassification.final, "%") ## [1] "Misclassification Error Rate for final model is 18.1560283687943 %" #Assessing final model accuracy via ROC curve ROC=roc(actual.churn.test,predicted_prob) ## Setting levels: control = No, case = Yes ## Setting direction: controls < cases</pre> plot(ROC, col="blue") 1.0 0.8 Sensitivity 1.0 0.5 0.0 Specificity auc(ROC) ## Area under the curve: 0.8427

paste("Title: Will the Telcom Customer Churn?- A Classification Analysis")

[1] "Title: Will the Telcom Customer Churn?- A Classification Analysis"

paste("Authors: Rajdeep Saha & Soumik Karmakar")

[1] "Authors: Rajdeep Saha & Soumik Karmakar"

Loading required package: lattice

Loading required package: carData

Attaching package: 'randomForest'

Type 'citation("pROC")' for a citation.

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

The following object is masked from 'package:ggplot2':

The following objects are masked from 'package:stats':

df <- read.csv('C:/Users/user/OneDrive/Desktop/Self Project/WA_Fn-UseC_-Telco-Customer-Churn.csv')</pre>

No

No

No

No No

Yes

Yes

Yes

No

No

No Month-to-month

No Month-to-month

No Month-to-month

Yes Month-to-month

One year

One year

29.85

108.15 Yes

820.50 Yes

1889.50

1840.75

: chr "7590-VHVEG" "5575-GNVDE" "3668-QPYBK" "7795-CF0CW" ...

: chr "Month-to-month" "One year" "Month-to-month" "One year" ...

\$ PaymentMethod : chr "Electronic check" "Mailed check" "Mailed check" "Bank transfer (automatic)" ...

"SeniorCitizen"

"PhoneService"

"OnlineBackup"

"StreamingMovies"

"MonthlyCharges"

"Partner"

"Contract"

"MultipleLines"

"TotalCharges"

"DeviceProtection"

MultipleLines

151.65

1 34

2

45

8

No

Yes

No

No

Contract PaperlessBilling

No

No

Yes

2

Yes

Yes

No

Yes

Yes

No

Yes

No

Yes

Yes

Yes

No

Yes

No

Yes

customerID gender SeniorCitizen Partner Dependents tenure PhoneService

0 No

No

No

No

29.85

56.95

53.85

42.30

70.70

99.65

: chr "Female" "Male" "Male" "...

: chr "Yes" "No" "No" "No" ... : chr "No" "No" "No" "No" ...

: chr "No" "Yes" "Yes" "No" ... ## \$ MultipleLines : chr "No phone service" "No" "No phone service" ...

: chr "Yes" "No" "Yes" "No" ...

: chr "No" "No" "No" "Yes" ...

: chr "No" "No" "Yes" "No" ...

: num 29.9 1889.5 108.2 1840.8 151.7 ...

: chr "No" "No" "No" "No" ...

: int 1 34 2 45 2 8 22 10 28 62 ...

PaymentMethod MonthlyCharges TotalCharges Churn

MultipleLines InternetService OnlineSecurity OnlineBackup DeviceProtection

0

0

DSL

DSL

DSL

DSL

Fiber optic

Fiber optic

7043 obs. of 21 variables:

\$ SeniorCitizen : int 00000000000...

\$ InternetService : chr "DSL" "DSL" "DSL" "DSL" ... ## \$ OnlineSecurity : chr "No" "Yes" "Yes" "Yes" ...

\$ DeviceProtection: chr "No" "Yes" "No" "Yes" ...

\$ StreamingMovies : chr "No" "No" "No" "No" ...

\$ PaperlessBilling: chr "Yes" "No" "Yes" "No" ...

\$ MonthlyCharges : num 29.9 57 53.9 42.3 70.7 ...

"gender"

"tenure"

"StreamingTV"

gender SeniorCitizen Partner Dependents tenure PhoneService

[9] "InternetService" "OnlineSecurity"

[17] "PaperlessBilling" "PaymentMethod"

df <- df[-which(colnames(df) == 'customerID')]</pre>

rm(list=ls()) set.seed(1) library(ggplot2) library(leaps) library(caret)

library(car)

library(corrplot)

library(tree) library(MASS)

##

##

#data access

1 7590-VHVEG Female

2 5575-GNVDE Male

4 7795-CFOCW Male

5 9237-HQITU Female

6 9305-CDSKC Female

1 No phone service

4 No phone service

Male

No

No

Yes

No

No

No

Yes

No

4 Bank transfer (automatic)

TechSupport StreamingTV StreamingMovies

No

No

No

No

Yes

Electronic check

Electronic check

Electronic check

Mailed check

Mailed check

3 3668-QPYBK

head(df)

2

3

5

1

2

3

4

5

6

1

2

3

5

6

dim(df)

str(df)

[1] 7043 21

'data.frame':

\$ customerID

\$ Dependents ## \$ tenure

\$ PhoneService

\$ OnlineBackup

\$ TotalCharges

\$ Churn

n <- nrow(df)</pre>

colnames(df)

#id column remove

[1] "customerID"

[13] "TechSupport"

[21] "Churn"

head(df)

[5] "Dependents"

\$ TechSupport ## \$ StreamingTV

\$ gender

\$ Partner

corrplot 0.89 loaded

library(randomForest)

randomForest 4.6-14

margin

Attaching package: 'pROC'

cov, smooth, var

library(pROC)