

Modelos y resultados

William Alexander Aguirre A

2024-04-13

1. Modelos de clasificación

En esta sección se muestran los resultados de los modelos más relevantes encontrados para clasificar como pobres o no pobres a los hogares. La variable objetivo a predecir es *Pobre* y como predictores se utilizaron las variables que se encuentran en Predictores.

Capturamos la data que se encuentra en el repositorio con las diferentes variables que hemos traído a partir de la información de personas y las transformaciones de interés. Para ver mayor detalle del procesamiento de la data se puede consultar el script Creación de variables de interés.R.

Predicciones con Logit.

Inicialmente se entrena un modelo con todos los predictores que no poseen valores perdidos en los set de entrenamiento y prueba. Inicialmente se hizo un entrenamiento con la data sin Bogotá para utilizar la variable Dominio, sin embargo, al no ser significativa se retiró la variable y trabajamos con el set de entrenamiento completo. También se retiraron las variables que tienen multicolinealidad, estos predictores se pueden consultar en Predictores Logit.

```
logit<-glm(as.formula(  
  paste("Pobre~",  
        paste(predictores_logit, collapse = " + "))),data = train_hogares)  
logit$coefficients
```

##	(Intercept)	informalidad_jefe1
##	1.05474964	0.10543097
##	edad_jefe_joven0	Ina_jefe
##	-0.03725487	-0.70229898
##	Des_jefe	Oc_jefe
##	-0.55403206	-0.78493194
##	ocupacion_jefe	educacion_jefe
##	-0.02485087	-0.03643114
##	sexo_jefe0	Clase2
##	0.03674382	0.02705082
##	P5000	P5010
##	-0.03408300	-0.04901314
##	P5090	Nper
##	0.02273378	0.01101098
##	Npersug	desempleo_jefe
##	0.06093213	0.14045400
##	tipo_casa0	Personas_habitacion_round1
##	0.02298655	-0.07324089

Se evalua el modelo dentro de muestra

```
train_hogares$prob_logit<-predict(logit,newdata = train_hogares)
train_hogares$logit_predict<-ifelse(train_hogares$prob_logit>0.5,1,0)

matrix_logit<-table(train_hogares$logit_predict,
                    train_hogares$Pobre)

confusionMatrix(matrix_logit,
                 positive = "1",
                 mode="prec_recall")
```

```
## Confusion Matrix and Statistics
##
##
##      0      1
## 0 128512 26870
## 1   3424  6154
##
##              Accuracy : 0.8164
##              95% CI : (0.8145, 0.8182)
##      No Information Rate : 0.7998
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.2186
##
##  Mcnemar's Test P-Value : < 2.2e-16
##
##              Precision : 0.64251
##              Recall : 0.18635
##              F1 : 0.28891
##              Prevalence : 0.20019
##              Detection Rate : 0.03731
##      Detection Prevalence : 0.05806
##              Balanced Accuracy : 0.58020
##
##      'Positive' Class : 1
##
```

Se realiza la predicción fuera de muestra y se genera el submission para Kaggle utilizando la sintaxis que se encuentra en Submission Logit

```
sub3<-test_hogares %>% select(id,logit_predict)
sub3<-sub3 %>% rename(pobre=logit_predict)
write_csv(x = sub3,"C:/Users/HP-Laptop/Documents/GitHub/Curso-Big-Data/Taller 2/2.Entregables/Submission
```

También se genera un modelo con las interacciones entre las variables predictores para interacciones. Evaluamos el modelo dentro de muestra.

```
interacciones<-c("informalidad_jefe","edad_jefe_joven","sexo_jefe","desempleo_jefe","tipo_casa")
```

```
logit2<-glm(as.formula(
  paste("Pobre~",
    paste(interacciones, collapse = " * "))),data = train_hogares)

train_hogares$prob_logit2<-predict(logit2,newdata = train_hogares)
train_hogares$logit_predict2<-ifelse(train_hogares$prob_logit2>0.5,1,0)
summary(train_hogares$prob_logit2)
```

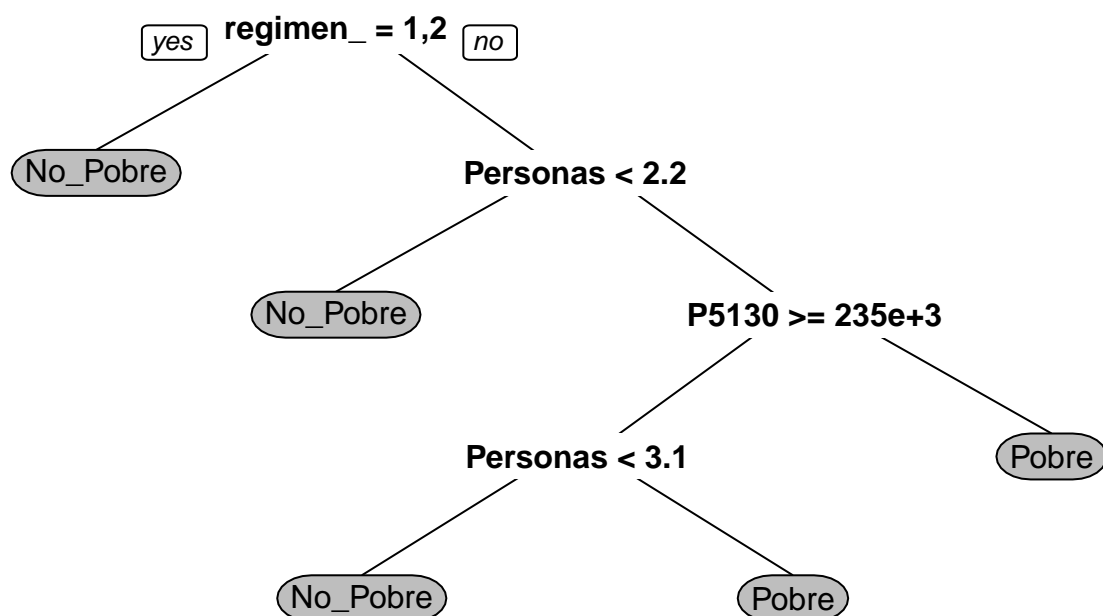
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.02374 0.13061 0.17777 0.20019 0.27360 0.43697
```

La probabilidad maxima que predice la interacción entre estas variables es inferior a 0.5 por lo que no se tomó en cuenta para hacer predicciones fuera de muestra.

Predicciones con Árbol

Para entrenar el árbol se utilizó la siguiente sintaxis.

```
arbol<-rpart(formula = as.formula(paste("pobre_texto~",
                                         paste(predictores,
                                               collapse = " + "))),
             data=train_hogares,
             method = "class",
             parms = list(split="Gini"))
prp(arbol,box.col = "gray")
```



Se evalúan los resultados dentro de muestra

```
train_hogares$predic_arbol<-predict(arbol,newdata =train_hogares,
                                   type = "class")
```

```
confusionMatrix(data = train_hogares$predic_arbol,
                 reference = train_hogares$pobre_texto,
                 positive="Pobre",
                 mode = "prec_recall")
```

Confusion Matrix and Statistics

##

Reference

Prediction No_Pobre Pobre

No_Pobre 128147 25715

Pobre 3789 7309

##

Accuracy : 0.8211

95% CI : (0.8193, 0.823)

No Information Rate : 0.7998

P-Value [Acc > NIR] : < 2.2e-16

##

Kappa : 0.2564

##

McNemar's Test P-Value : < 2.2e-16

##

```
## Precision : 0.65859
## Recall : 0.22132
## F1 : 0.33131
## Prevalence : 0.20019
## Detection Rate : 0.04431
## Detection Prevalence : 0.06728
## Balanced Accuracy : 0.59630
##
## 'Positive' Class : Pobre
##
```

Se realiza la predicción fuera de muestra y se genera el submission para Kaggle. La sintaxis utilizada para generar el submission se encuentra en Submission Árbol.

```
test_hogares$predic_arbol<-predict(arbol,newdata =test_hogares,type = "class")
table(test_hogares$predic_arbol)
```

```
##
## No_Pobre      Pobre
##      61469      4699
```

```
sub9<-test_hogares %>% select(id,predic_arbol)
sub9<-sub9 %>% rename(pobre=predic_arbol)
sub9$pobre<-ifelse(sub9$pobre=="Pobre",1,0)
table(test_hogares$predic_arbol)
table(sub9$pobre)
write_csv(x = sub9,"C:/Users/HP-Laptop/Documents/GitHub/Curso-Big-Data/Taller 2/2.Entregables/Submission")
```

Predicciones con Random Forest

El algoritmo de Random Forest no nos permite utilizar variables con valores perdidos por lo que se excluyen de los predictores algunas variables que no aplican para toda la muestra. Los predictores utilizados para entrenar el modelo se pueden consultar en Predictores Random Forest. Se realiza un entrenamiento inicial de un modelo de Random Forest utilizando la siguiente sintaxis.

```
RF<- ranger(formula = as.formula(paste("pobre_texto~",
                                     paste(predictores, collapse = " + "))),
            data = train_hogares,
            num.trees= 500, ## Numero de arboles a estimar
            mtry= 4,      # N. var aleatoriamente seleccionadas en cada partición.
            min.node.size = 1, ## Numero minimo de observaciones en un nodo
            importance="impurity")
```

RF

```
## Ranger result
##
## Call:
## ranger(formula = as.formula(paste("pobre_texto~", paste(predictores, collapse = " + "))), data =
##
## Type:                Classification
## Number of trees:     500
```

```
## Sample size: 164960
## Number of independent variables: 21
## Mtry: 4
## Target node size: 1
## Variable importance mode: impurity
## Splitrule: gini
## OOB prediction error: 16.91 %
```

Para definir la cantidad optima de variables seleccionada por cada partición se prueba aumentar la cantidad de arboles para identificar que hay una reducción importante del OOB predictor error.

```
RF1000<- ranger(formula = as.formula(paste("pobre_texto~",
                                           paste(predictores, collapse = " + "))),
               data = train_hogares,
               num.trees= 1000, ## Aumentamos de 500 a 1000
               mtry= 4,
               min.node.size = 1,
               importance="impurity")
```

RF1000

```
## Ranger result
```

```
##
```

```
## Call:
```

```
## ranger(formula = as.formula(paste("pobre_texto~", paste(predictores, collapse = " + "))), data
```

```
##
```

```
## Type: Classification
```

```
## Number of trees: 1000
```

```
## Sample size: 164960
```

```
## Number of independent variables: 21
```

```
## Mtry: 4
```

```
## Target node size: 1
```

```
## Variable importance mode: impurity
```

```
## Splitrule: gini
```

```
## OOB prediction error: 16.95 %
```

Se observa que no hay reducción importante del OOB predictor error. Para continuar de afinar hiperparametros del modelo aumentamos la cantidad de minima de observaciones por nodo.

```
RF_NODE100<- ranger(formula = as.formula(paste("pobre_texto~",
                                                paste(predictores, collapse = " + "))),
                   data = train_hogares,
                   num.trees= 500,
                   mtry= 4,
                   min.node.size = 100, #Aumentamos de 1 a 100
                   importance="impurity")
```

RF_NODE100

```
## Ranger result
```

```
##
```

```
## Call:
```

```
## ranger(formula = as.formula(paste("pobre_texto~", paste(predictores, collapse = " + "))), data
```

```
##
```

```
## Type: Classification
```

```
## Number of trees:          500
## Sample size:              164960
## Number of independent variables: 21
## Mtry:                     4
## Target node size:         100
## Variable importance mode:  impurity
## Splitrule:                gini
## OOB prediction error:     16.92 %
```

Vemos que tampoco existe un cambio importante en el OOB predictor error. Finalmente probamos aumentar la cantidad de variables por arbol.

```
RF_5VAR<- ranger(formula = as.formula(paste("pobre_texto~",
                                             paste(predictores, collapse = " + "))),
                  data = train_hogares,
                  num.trees= 500,
                  mtry= 5,
                  min.node.size = 1,
                  importance="impurity")
```

```
## Growing trees.. Progress: 17%. Estimated remaining time: 2 minutes, 27 seconds.
## Growing trees.. Progress: 35%. Estimated remaining time: 1 minute, 57 seconds.
## Growing trees.. Progress: 51%. Estimated remaining time: 1 minute, 28 seconds.
## Growing trees.. Progress: 68%. Estimated remaining time: 57 seconds.
## Growing trees.. Progress: 85%. Estimated remaining time: 28 seconds.
## Growing trees.. Progress: 99%. Estimated remaining time: 1 seconds.
```

```
RF_6VAR<- ranger(formula = as.formula(paste("pobre_texto~",
                                             paste(predictores, collapse = " + "))),
                  data = train_hogares,
                  num.trees= 500,
                  mtry= 6,
                  min.node.size = 1,
                  importance="impurity")
```

```
## Growing trees.. Progress: 14%. Estimated remaining time: 3 minutes, 4 seconds.
## Growing trees.. Progress: 29%. Estimated remaining time: 2 minutes, 34 seconds.
## Growing trees.. Progress: 43%. Estimated remaining time: 2 minutes, 4 seconds.
## Growing trees.. Progress: 57%. Estimated remaining time: 1 minute, 32 seconds.
## Growing trees.. Progress: 72%. Estimated remaining time: 1 minute, 0 seconds.
## Growing trees.. Progress: 86%. Estimated remaining time: 30 seconds.
## Growing trees.. Progress: 99%. Estimated remaining time: 3 seconds.
```

```
RF_7VAR<- ranger(formula = as.formula(paste("pobre_texto~",
                                             paste(predictores, collapse = " + "))),
                  data = train_hogares,
                  num.trees= 500,
                  mtry= 7,
                  min.node.size = 1,
                  importance="impurity")
```

```
## Growing trees.. Progress: 13%. Estimated remaining time: 3 minutes, 35 seconds.
```

```
## Growing trees.. Progress: 26%. Estimated remaining time: 3 minutes, 0 seconds.
## Growing trees.. Progress: 39%. Estimated remaining time: 2 minutes, 27 seconds.
## Growing trees.. Progress: 52%. Estimated remaining time: 1 minute, 54 seconds.
## Growing trees.. Progress: 64%. Estimated remaining time: 1 minute, 28 seconds.
## Growing trees.. Progress: 74%. Estimated remaining time: 1 minute, 7 seconds.
## Growing trees.. Progress: 86%. Estimated remaining time: 35 seconds.
## Growing trees.. Progress: 98%. Estimated remaining time: 4 seconds.
```

```
RF_5VAR$prediction.error
```

```
## [1] 0.1704292
```

```
RF_6VAR$prediction.error
```

```
## [1] 0.1723448
```

```
RF_7VAR$prediction.error
```

```
## [1] 0.174703
```

```
## [1] 0.1703322
```

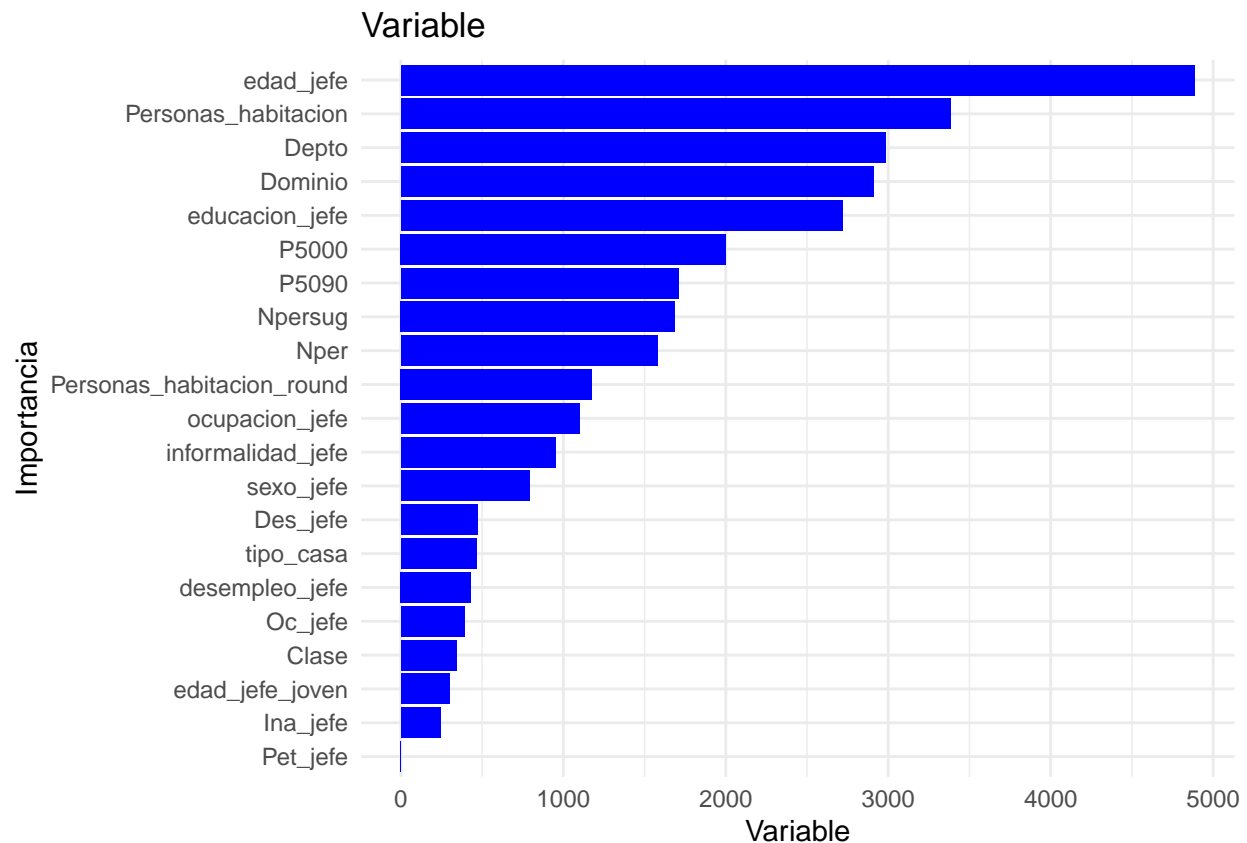
```
## [1] 0.1721266
```

```
## [1] 0.1749333
```

Observamos que después de cinco variables por árbol el OBB predictor error deja de reducirse. Por lo que se toma la decisión de trabajar con el primer árbol.

Podemos observar también la importancia que tienen las variables en el árbol para identificar que algunas cobran mayor relevancia en comparación de un árbol sencillo.

```
imp<-importance(RF)
imp<-data.frame(variables=names(imp),importancia=imp)
ggplot(imp, aes(x = reorder(variables, importancia) , y =importancia )) +
  geom_bar(stat = "identity", fill = "blue") +
  labs(title = "Variable ", x = "Importancia", y="Variable") +
  theme_minimal() +
  coord_flip()
```

Analizamos las predicciones realizadas dentro de muestra

```
confusionMatrix(data = RF$predictions,
  reference = train_hogares$pobre_texto,
  positive = "Pobre",
  mode = "prec_recall")
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction No_Pobre  Pobre
## No_Pobre   126035  21999
## Pobre       5901   11025
##
##           Accuracy : 0.8309
##           95% CI : (0.8291, 0.8327)
## No Information Rate : 0.7998
## P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.3538
##
## McNemar's Test P-Value : < 2.2e-16
##
##           Precision : 0.65136
##           Recall : 0.33385
##           F1 : 0.44144
```

```
##           Prevalence : 0.20019
##       Detection Rate : 0.06683
## Detection Prevalence : 0.10261
##       Balanced Accuracy : 0.64456
##
##       'Positive' Class : Pobre
##
```

Predecimos fuera de muestra y generamos el submission para Kaggle. Podemos consultar la sintaxis con la que se dió formato a las predicción en formato de submission con Random Forest.

2. Modelos de regresión del ingreso

En esta sección se muestran los resultados de los modelos más relevantes encontrados para predecir el ingreso con el objetivo de apartir de estas predicciones poder clasificar como pobres o no pobres a los hogares, tomando el umbral de pobreza como límite de clasificación entre los dos grupos de hogares.

Predicciones con regresión lineal

Para realizar la predicción del ingreso de los hogares utilizamos como variable dependiente el ingreso total de la unidad de gasto con imputación de arriendo a propietarios y usufructuario (*Ingtotugarr*). Las variables predictoras se pueden consultar en predictores para regresión lineal.

Para correr el modelo se utilizó la siguiente sintaxis:

```
modelo_lm<-lm(as.formula( #Entrenamos el modelo
  paste("Ingtotug~",paste(predictores, collapse = " + "))),
  data =train_hogares_sin_bogota)
modelo_lm$coefficients
```

```
##      (Intercept)      desempleo_jefe      educacion_jefe
##      -682920.86      -104313.55      467822.23
##      sexo_jefe0      Clase2      DominioBARRANQUILLA
##      -379481.53      -520708.26      379685.66
## DominioBUCARAMANGA      DominioCALI      DominioCARTAGENA
##      582796.78      449737.26      144594.93
## DominioCUCUTA      DominioFLORENCIA      DominioIBAGUE
##      14418.17      -223839.86      98501.99
## DominioMANIZALES      DominioMEDELLIN      DominioMONTERIA
##      506834.87      672627.01      -27551.28
## DominioNEIVA      DominioPASTO      DominioPEREIRA
##      158617.95      59890.86      240486.36
## DominioPOPAYAN      DominioQUIBDO      DominioRESTO URBANO
##      -124724.49      -346242.25      -183826.42
## DominioRIOHACHA      DominioRURAL      DominioSANTA MARTA
##      26308.66      NA      184594.80
## DominioSINCELEJO      DominioTUNJA      DominioVALLEDUPAR
##      165588.93      168525.54      -180587.08
## DominioVILLAVICENCIO      P5090      Ina_jefe
##      335741.72      -287730.31      1629924.17
##      Des_jefe      Oc_jefe
##      636896.58      1516954.22
```

Realizamos las predicciones en el conjunto de entrenamiento para poder clasificar

```
train_hogares_sin_bogota$modelo_lm_ingreso<-predict(object = modelo_lm,
                                                    newdata = train_hogares_sin_bogota)

## Warning in predict.lm(object = modelo_lm, newdata = train_hogares_sin_bogota):
## prediction from rank-deficient fit; attr(*, "non-estim") has doubtful cases

#Clasificamos los hogares en pobres o no según el ingreso predicho
train_hogares_sin_bogota$modelo_lm_predict<-ifelse(
  train_hogares_sin_bogota$modelo_lm_ingreso<train_hogares_sin_bogota$Lp,
  1,0)
matrix_lm<-table(train_hogares_sin_bogota$modelo_lm_predict,train_hogares_sin_bogota$Pobre)
confusionMatrix(matrix_lm,
                  positive="1",
                  mode = "prec_recall")

## Confusion Matrix and Statistics
##
##
##           0          1
## 0 119888  29453
## 1   2512   2540
##
##              Accuracy : 0.793
##              95% CI : (0.7909, 0.795)
##    No Information Rate : 0.7928
##    P-Value [Acc > NIR] : 0.4317
##
##              Kappa : 0.0854
##
##  Mcnemar's Test P-Value : <2e-16
##
##              Precision : 0.50277
##              Recall    : 0.07939
##              F1        : 0.13713
##              Prevalence : 0.20722
##              Detection Rate : 0.01645
##              Detection Prevalence : 0.03272
##              Balanced Accuracy : 0.52943
##
##              'Positive' Class : 1
##
```

Luego realizamos la predicción en el conjunto de prueba y exportamos con el formato adecuado para realizar submitir en Keaggle. La sintaxis utilizada para esto se encuentra en formato para submission con regresión lineal.

```
sub1<-test_hogares %>% select(id,modelo_lm_predict)
sub1<-sub1 %>% rename(pobre=modelo_lm_predict)
write_csv(x = sub1,"C:/Users/HP-Laptop/Documents/GitHub/Curso-Big-Data/Taller 2/2.Entregables/Submission
```

Posteriormente se realizan diferentes combinaciones de variables predictivas para encontrar aquellas que tengan una mejor capacidad de predicción.

Predicciones con Regularización.

Para realizar estos modelos utilizamos solo las variables que no poseen NA en train ni en test, estos se pueden consultar en predictores regularización. Para aplicar regularización tenemos que trabajar con nuestros datos en formato matriz, la sintaxis utilizada también se puede observar en transformación en matriz para regularización.

Ridge, Lasso y Elastic Net

Para entrenar los modelos fue necesario afinar los parámetros lambda utilizados para la predicción. En este caso, el parámetro lambda mínimo fue encontrado mediante cross validation utilizando la siguiente sintaxis:

```
cv_ridge <- cv.glmnet(x = X,
                     y = Y, alpha = i)
# Para Ridge i=1, Lasso i=0 Lasso, Elastic Net i=0.75.

coef(cv_ridge, s = "lambda.min")
test_hogares$predict_ingreso_ridge <- predict(cv_ridge,
                                              newx = X_test,
                                              s = "lambda.min")
```

Elastic Net con Cartet

Para identificar los hiperparametros optimos de Elastic Net con Caret se utilizó la siguiente sintaxis.

```
tc_10 <- trainControl(method = "cv", number = 10)

en_caret <- train(x=X,y=Y,method = "glmnet",trControl = tc_10,
                 tuneLength=100)
```

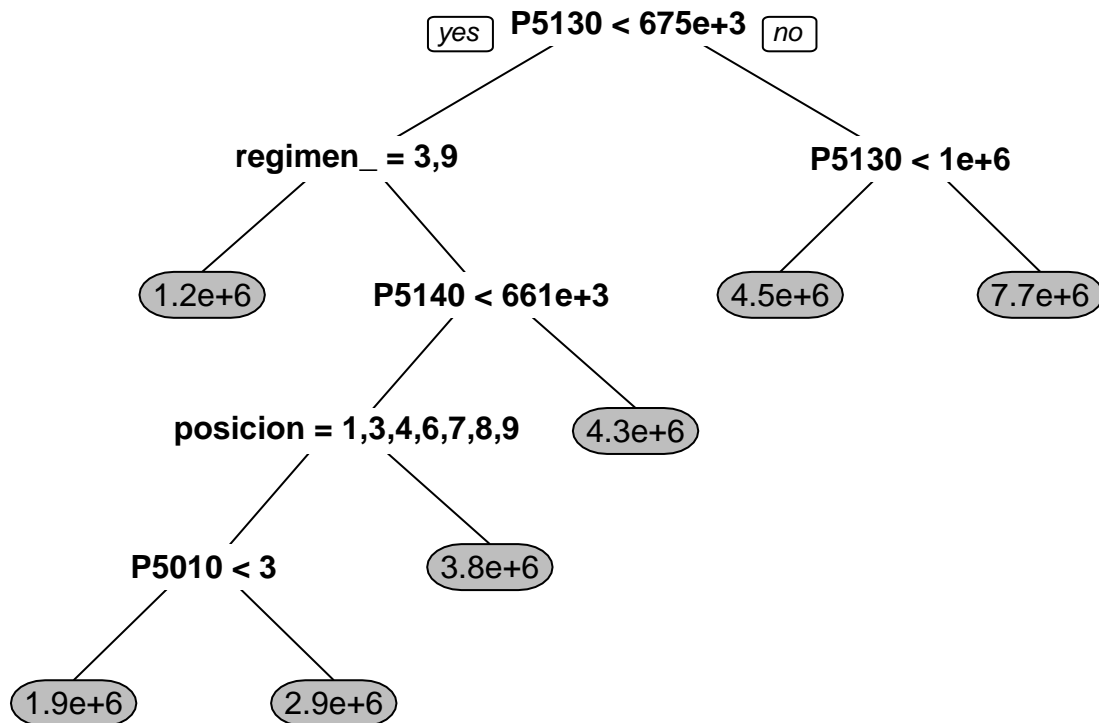
```
## Warning in nominalTrainWorkflow(x = x, y = y, wts = weights, info = trainInfo,
## : There were missing values in resampled performance measures.
```

Predicciones con Árbol

Para predecir el ingreso contamos con mas flexibilidad al momento de utilizar variables con NA, estas variables se pueden consultar en predictores para árbol de regresión . Para realizar la predicción con árboles de decisión se realizó la siguiente sintaxis

```
arbol <- rpart(formula = as.formula(paste("Ingtotugarr~",
                                         paste(predictores,
                                                collapse = " + "))),
              data=train_hogares_sin_bogota,
              parms = list(split="Gini"))

prp(arbol, box.palette = "gray")
```



Luego de entrenar el modelo realizamos la predicción en train y evaluamos su resultado dentro de muestra.

```

train_hogares_sin_bogota$predic_arbol_ingreso<-predict(arbol,
                                                       newdata =train_hogares_sin_bogota)

train_hogares_sin_bogota$predict_arbol<-ifelse(
  train_hogares_sin_bogota$predic_arbol_ingreso<
  train_hogares_sin_bogota$Lp*train_hogares_sin_bogota$Npersug,
  1,0)

matriz_arbol<-table(train_hogares_sin_bogota$predict_arbol,train_hogares_sin_bogota$Pobre)
confusionMatrix(matriz_arbol,
                 positive="1",
                 mode = "prec_recall")

```

```

## Confusion Matrix and Statistics
##
##
##          0          1
## 0 115419  24044
## 1   6981   7949
##
##               Accuracy : 0.7991
##               95% CI : (0.797, 0.801)
##    No Information Rate : 0.7928
##    P-Value [Acc > NIR] : 5.392e-10

```

```
##
##           Kappa : 0.2384
##
## Mcnemar's Test P-Value : < 2.2e-16
##
##           Precision : 0.53242
##           Recall : 0.24846
##           F1 : 0.33881
##           Prevalence : 0.20722
##           Detection Rate : 0.05149
##           Detection Prevalence : 0.09670
##           Balanced Accuracy : 0.59571
##
##           'Positive' Class : 1
##
```

Luego realizamos la predicción fuera de muestra:

```
test_hogares$predic_arbol_ingreso<-predict(arbol,
                                           newdata =test_hogares)

test_hogares$predict_arbol3<-ifelse(
  test_hogares$predic_arbol_ingreso<
    test_hogares$Lp*test_hogares$Npersug,
  1,0)
table(test_hogares$predict_arbol)
```

```
##
##      0      1
## 59622 6546
```