

```
1 #####LIBRARY AND FUNCTIONS#####
2 #For cleaning and reading data
3 library(tidyverse)
4 library(caret)
5 library(themis)
6 library(EZtune)
7 library(MLmetrics)
8
9 #
10 # library(tidymodel)
11 #for plotting
12 theme_set(theme_bw())
13
14
15 #change wd, and import data
16 # setwd("~/Documents/Box Sync/Statistics Master/Fall 2022/STATS504/hw5")
17 setwd("/Users/brad/Downloads/hw5")
18 df <- read.csv("data/nuMoM2bsubset.csv")
19 # load("brad_models.RData") # Saved my global enviornment just in case...
20
21 ##### Groups 4, 9, 14 outcome: dv.hypertension1
22
23 null_outcomes <-c("dv.diabetes1",
24                   "dv.v1epdstotal",
25                   "dv.gestweeks",
26                   "dv.preeclampsia")
27
28 df <- df[,!(names(df) %in% null_outcomes)]
29
30 ##### DATA CLEANING #####
31 test_df = df %>% summarise(across(everything(), list(min,max)))
32 test_df = t(test_df)
33
34 #Convert booleans to "TRUE"/"FALSE"
35 df$emosupport <- df$emosupport == 1
36 df$financialsupport <- df$financialsupport==1
37 df$prenatalsupport <- df$prenatalsupport == 1
38 df$financialsupport <- df$financialsupport ==1
39 df$deliverysupport <- df$deliverysupport ==1
40 df$exercise <- df$exercise == 1
41 df$dv.hypertension1 <- df$dv.hypertension1==1
42 df$kidney1 <- df$kidney1==1
43 df$lupus1 <- df$lupus1 == 1
44 df$collagen1 <- df$collagen1==1
45 df$crohns1 <- df$crohns1 == 1
46 df$pcos1 <- df$pcos1 == 1
47
48
49 #Three level factors....
50 df$familypreeclampsia <- as.factor(df$familypreeclampsia)
51 df$bornearly <- as.factor(df$bornearly)
52
53
54
55 #Higher level factors....
56 df$worryfambaby <- as.factor(df$worryfambaby)
```

```

57 df$worryhealthcare <- as.factor(df$worryhealthcare)
58 df$worrysymptoms <- as.factor(df$worrysymptoms) #what are the levels for this? There should be a codebook.
59 df$discrimination <- as.factor(df$discrimination)
60 df$race <- as.factor(df$race)
61
62 # skimr::skim(df)
63 #####SUMMARY TABLE#####
64 skimmed_df = skimr::skim(df)
65 skimmed_df$n_missing = NULL
66 skimmed_df$complete_rate= NULL
67 skimmed_df$numeric.hist = NULL
68 skimmed_df$factor.ordered = NULL
69
70 #For factors
71 skimmed_df$factor.top_counts[8:18] <- skimmed_df$logical.count[8:18]
72 skimmed_df$factor.top_counts[8:18] <- skimmed_df$logical.count[8:18]
73 skimmed_df$factor.n_unique[8:18] <- 2
74 skimmed_df$numeric.mean[8:18] <- skimmed_df$logical.mean[8:18]
75 skimmed_df$numeric.mean <- round(skimmed_df$numeric.mean, digits= 4)
76 skimmed_df$logical.mean <- NULL
77 skimmed_df$logical.count = NULL
78 skimmed_df$factor.top_counts[19:26] <- paste0(
  ("",skimmed_df$numeric.p25[19:26],",",skimmed_df$numeric.p75[19:26],")")
79 skimmed_df$numeric.sd <- NULL
80 skimmed_df$numeric.p0 <- NULL
81 skimmed_df$numeric.p25 <- NULL
82 skimmed_df$numeric.p50 <- NULL
83 skimmed_df$numeric.p75 <- NULL
84 skimmed_df$numeric.p100 <- NULL
85 # write.csv(skimmed_df, "data/baseline.csv")
86
87
88 rm(test_df, null_outcomes) #drop unused items
89
90 #####IMPUTE VALUES#####
91 colSums(df == 0)
92 #age
93 #psstotla
94 #ssqmean
95 # prepreglbs
96
97 df$age[df$age==0] <- mean(df$age[df$age!=0])
98 df$prepreglbs[df$prepreglbs==0] <- mean(df$prepreglbs[df$prepreglbs!=0])
99
100 #####TEST TRAIN#####
101 set.seed(1123)
102 size = floor(0.3*dim(df))
103
104 id = sample(c(1:7934),replace=F, size = floor(0.3*dim(df)))[1])
105 train<-df[-id,]
106 test<-df[id,]
107
108 # write.csv(test,"test_df.csv")
109 # write.csv(train,"train_df.csv")
110
111 #####ENCODE TRAIN#####
112 dmy <- dummyVars(" ~ .", data = train)
113 train_hot <- data.frame(predict(dmy, newdata = train))
114
115 train_hot_X = train_hot[,!(names(train_hot) %in% c("dv.hypertension1FALSE","dv.hypertension1TRUE"))]

```

```

116
117 #Outcome variable needs to have a valid name. Use make.names() or use
118 train_hot_y = factor(train_hot$dv.hypertension1TRUE,
119                       levels = c(1,0),
120                       labels = c("yes","no"))
121
122
123
124 #Combined data frame .....
125 train_hot_X_y<- train_hot_X
126 train_hot_X_y$dv.hypertension <- train_hot_y
127
128
129 #####UP SAMPLE TRAINING DAT#####
130 #minority class has 50% of observations as majority
131 train_up50<- smote(train_hot_X_y, var ="dv.hypertension", over=0.5, k=10)
132 table(train_up50$dv.hypertension)
133 #majority class has 100% of observations as majority
134 train_up100<- smote(train_hot_X_y, var ="dv.hypertension", over=1, k=10)
135 table(train_up100$dv.hypertension)
136
137 #Export
138 # write.csv(train_up50, "data/train_one_hot_50.csv")
139 # write.csv(train_up100, "data/train_one_hot_100.csv")
140
141
142
143
144 #####columns with no variance in training data #####
145 X_no_var = nearZeroVar(train_hot_X_y)
146 X_no_var_names = names(train_hot_X_y)[X_no_var]
147 X_no_var_names = X_no_var_names[-c(1,37)] #hypertension has low variance, as does race == native
148
149 #Drop columns with no variance....
150 train_hot_X_y <- train_hot_X_y[!(names(train_hot_X) %in% X_no_var_names)]
151 train_up50 <- train_up50[!(names(train_up50) %in% X_no_var_names)]
152 train_up100 <- train_up100[!(names(train_up100) %in% X_no_var_names)]
153
154 # Check the distribution of these variables....
155 # train_hot_X[,X_no_var_names] %>% skimr::skim()
156
157 #####DOWNSAMPLE#####
158 set.seed(1123)
159 train_hot_down = downSample(x=train_hot_X_y,
160                             y=train_hot_X_y$dv.hypertension)
161 train_hot_down$Class <- NULL
162
163 #####ENCODE TEST#####
164 dmy <- dummyVars(" ~ .", data = test)
165 test_hot <- data.frame(predict(dmy, newdata = test))
166
167 test_hot_X = test_hot[!(names(test_hot) %in% c("dv.hypertension1FALSE", "dv.hypertension1TRUE"))]
168
169 #Outcome variable needs to have a valid name. Use make.names() or use
170 test_hot_y = factor(test_hot$dv.hypertension1TRUE,
171                     levels = c(1,0),
172                     labels = c("yes","no"))
173
174 #combine into one data frame....
175 test_hot_X_y = test_hot_X

```

```

176 test_hot_X_y$dv.hypertension = test_hot_y
177
178 #Drop columns with no variance
179 test_hot_X_y <- test_hot_X_y[,!(names(test_hot_X) %in% X_no_var_names)]
180
181
182 #####CLEAN WORKSPACE#####
183 rm(list = c("#test_hot_X_sub",
184             "test_hot_X",
185             "test_hot",
186             "test_hot_y",
187             # "train_hot_X_sub",
188             "train_hot_X",
189             "train_hot",
190             "train_hot_y",
191             "dmy"))
192
193 ##### HYPERPARAMETER TUNING#####
194
195 models <- caret::modellookup() #what models are in the caret package?
196
197 ##### Adaboost.M1#####
198
199 #!!!!!!!!!!!!!!!!!! WARNING THE FOLLOWING CHUNKS TAKE ~40 Minutes to run!!!!!!!!!!!!!!!!!!!!
200
201 fitGrid_ada <- expand.grid(mfinal = c(1,6,9,100),
202                           # mfinal = (1:3)*3,
203                           # maxdepth = c(1:3),
204                           maxdepth = c(1,2,4),
205                           coeflearn = c("Breiman"))
206
207 fitControl_ada <- trainControl(method = "repeatedcv",
208                                repeats = 5,
209                                classProbs = T,
210                                # summaryFunction = twoClassSummary,
211                                summaryFunction = prSummary)
212 #on up sampled
213
214 # using the adaboost.m1 package....
215 set.seed(1123)
216 start_time = Sys.time()
217 ada.mod <- train(x=train_hot_X_y[,-48],
218                 y= train_hot_X_y$dv.hypertension,
219                 method = 'AdaBoost.M1',
220                 trControl = fitControl_ada,
221                 tuneGrid = fitGrid_ada,
222                 metric = "AUC",
223                 verbose = TRUE)
224 total_time <- Sys.time() - start_time
225 total_time
226
227 # Upsampled to be 50% majority class
228 set.seed(1123)
229 start_time <- Sys.time()
230 ada50.mod <- train(x=train_up50[,-48],
231                   y= train_up50$dv.hypertension,
232                   method = 'AdaBoost.M1',
233                   trControl = fitControl_ada,
234                   tuneGrid = fitGrid_ada,
235                   metric = "AUC",

```

```

236         verbose = TRUE)
237 total_time <- Sys.time() - start_time
238 total_time
239
240
241 # upsampled to Matched classes - change metric to ROC.
242 set.seed(1123)
243 total_time <- Sys.time()
244 ada100.mod <- train(x=train_up100[,-48],
245                    y= train_up100$dv.hypertension,
246                    method = 'AdaBoost.M1',
247                    trControl = fitControl_ada,
248                    tuneGrid = fitGrid_ada,
249                    metric = "ROC",
250                    verbose = TRUE)
251 total_time <- Sys.time() - start_time
252 total_time
253
254 # On downsampled - use ROC
255 set.seed(1123)
256 start_time <- Sys.time()
257 adadown.mod <- train(x=train_hot_down[,-48],
258                    y= train_hot_down$dv.hypertension,
259                    method = 'AdaBoost.M1',
260                    trControl = fitControl_ada,
261                    tuneGrid = fitGrid_ada,
262                    metric = "ROC",
263                    verbose = TRUE)
264 total_time <- Sys.time() - start_time
265 total_time
266
267
268 ##### USING GBM PACKAGE #####3
269
270 # set.seed(1123)
271 # start_time <- Sys.time()
272 # ada.mod <- train(x=train_hot_X_y[,-48],
273 #                 y= train_hot_X_y$dv.hypertension,
274 #                 distribution = 'adaboost',
275 #                 method="gbm",
276 #                 trControl = fitControl_ada,
277 #                 tuneGrid = fitGrid_ada,
278 #                 metric = "AUC",
279 #                 verbose = TRUE)
280 # total_time <- Sys.time() - start_time
281 # total_time
282
283 # set.seed(1123)
284 # start_time <- Sys.time()
285 # ada50.mod <- train(x=train_up50[,-48],
286 #                   y= train_up50$dv.hypertension,
287 #                   distribution = 'adaboost',
288 #                   method="gbm",
289 #                   trControl = fitControl_ada,
290 #                   tuneGrid = fitGrid_ada,
291 #                   metric = "AUC",
292 #                   verbose = TRUE)
293 # total_time <- Sys.time() - start_time
294 # total_time
295

```

```

296
297 # set.seed(1123)
298 # total_time <- Sys.time()
299 # ada100.mod <- train(x=train_up100[,-83],
300 #                     y= train_up100$dv.hypertension,
301 #                     distribution = 'adaboost',
302 #                     method="gbm",
303 #                     trControl = fitControl_ada,
304 #                     tuneGrid = fitGrid_ada,
305 #                     metric = "ROC",
306 #                     verbose = TRUE)
307 # total_time <- Sys.time() - start_time
308 # total_time
309
310
311 #Fit downsampled data on finer grid...
312 # fitGrid_ada <- expand.grid(interaction.depth = c(1, 3, 6, 9),
313 #                             n.trees = c(1,10,20,50,100),
314 #                             shrinkage = seq(.0005, .05,.0005),
315 #                             n.minobsinnode = 10)
316 #
317 # fitControl_ada <- trainControl(method = "repeatedcv",
318 #                                 repeats = 5,
319 #                                 classProbs = T,
320 #                                 summaryFunction = twoClassSummary)
321 #on up sampled
322 # set.seed(1123)
323 # start_time <- Sys.time()
324 # adadown.mod <- train(x=train_hot_down[,-83],
325 #                      y= train_hot_down$dv.hypertension,
326 #                      distribution = 'adaboost',
327 #                      method="gbm",
328 #                      trControl = fitControl_ada,
329 #                      tuneGrid = fitGrid_ada,
330 #                      metric = "ROC",
331 #                      verbose = TRUE)
332 # total_time <- Sys.time() - start_time
333 # total_time
334
335
336
337 #####PREDICTION#####
338 # Class Predictions
339 test_predada <- predict(object = ada.mod,newdata = test_hot_X_y[,-48])
340 test_predada50 <- predict(object = ada50.mod,newdata = test_hot_X_y[,-48])
341 test_predada100 <- predict(object = ada100.mod,newdata = test_hot_X_y[,-48])
342 test_predadadown <- predict(object = adadown.mod,newdata = test_hot_X_y[,-48])
343
344 # Probabilities
345 test_predada_p <- predict(object = ada.mod,newdata = test_hot_X_y[,-48],type="prob")
346 test_predada50_p <- predict(object = ada50.mod,newdata = test_hot_X_y[,-48],type = "prob")
347 test_predada100_p <- predict(object = ada100.mod,newdata = test_hot_X_y[,-48],type = "prob")
348 test_predadadown_p <- predict(object = adadown.mod,newdata = test_hot_X_y[,-48],type = "prob")
349
350 #####PRAUC#####
351 ada_prauc = MLmetrics::PRAUC(test_predada_p$yes, test_hot_X_y$dv.hypertension)
352 ada_prauc
353
354 adaup_prauc = MLmetrics::PRAUC(test_predada100_p$yes, test_hot_X_y$dv.hypertension)
355 adaup_prauc

```

```

356 ##### AUROC #####
357 # library(pROC)
358 # ada.roc <- roc(test_hot_X_y$dv.hypertension, test_predada_p$yes)
359 # # plot(ada.roc, print.thres="best", print.thres.best.method="closest.topleft")
360 # ada50.roc <- roc(test_hot_X_y$dv.hypertension, test_predada50_p$yes)
361 # ada100.roc <- roc(test_hot_X_y$dv.hypertension, test_predada100_p$yes)
362 #
363 #
364 # plot(ada50.roc, print.thres="best", print.thres.best.method="closest.topleft")
365 # result.coords <- coords(ada.roc, "best", best.method="closest.topleft", ret=c("ppv", "tpr"))
366 # print(result.coords)#to get threshold and accuracy
367
368
369 #####CONFUSION METRICS#####
370 conf_matada = table("truth"=test_hot_X_y$dv.hypertension,"pred"= test_predada)
371 conf_matada=confusionMatrix(conf_matada, mode = "everything", positive = "yes")
372 conf_matada
373
374
375 conf_matada50 = table("truth"=test_hot_X_y$dv.hypertension,"pred"= test_predada50)
376 conf_matada50=confusionMatrix(conf_matada50, mode = "everything", positive = "yes")
377 conf_matada50
378
379 conf_matada100 = table("truth"=test_hot_X_y$dv.hypertension,"pred"= test_predada100) #Adaboost looks a bit better
380 conf_matada100=confusionMatrix(conf_matada100, mode = "everything", positive ="yes")
381 conf_matada100
382
383 conf_matadadown= table("truth"=test_hot_X_y$dv.hypertension,"pred"= test_predadadown) #Adaboost looks a bit better
384 conf_matadadown=confusionMatrix(conf_matadadown, mode = "everything", positive ="yes")
385 conf_matadadown
386
387 metrics_df = data.frame(ada = conf_matada$byClass,
388                          ada50 = conf_matada50$byClass,
389                          ada.tune100 = conf_matada100$byClass)
390
391 # save.image(file='brad_models.RData')
392 # load("brad_svm_env.RData")
393
394 #####PLOTS #####
395 hyp_race_p = df %>% group_by(race) %>% summarise(p = mean(dv.hypertension1)) %>%
396   ggplot(aes(x=reorder(race, -p), y = p, fill = race))+
397   geom_bar(stat='identity')+
398   scale_y_continuous(labels = scales::percent)+
399   xlab("Race")+
400   ylab("Perc. of Women w/ Hypertension")+
401   guides(fill="none")
402
403
404 ggplot(data = df, aes(x=prepreglbs, fill = race, group = race))+
405   # geom_density(alpha=0.4)+
406   geom_histogram(aes(y=stat(density)))+
407   # scale_y_continuous(labels = scales::percent)+
408   # xlab("Frequency")+
409   # scale_y_continuous(labels = percent )
410   facet_grid(rows = vars(race))
411   # ylab("Perc. w/ Hypertension")+
412   # guides(fill="none")
413
414

```

```
415 ggplot(data=df)+
416   geom_jitter(aes(x=diastolic,
417                   y = systolic,
418                   col = dv.hypertension1),
419               alpha=0.5)+
420   facet_wrap(~dv.hypertension1)+
421   scale_x_continuous(sec.axis = sec_axis(~ . ,
422                                           name = "Has Hypertension",
423                                           breaks = NULL,
424                                           labels = NULL))+
425   # scale_y_continuous(labels = scales::percent)+
426   xlab("Systolic Blood Pressure At First Visit")+
427   ylab("Diastolic Blood Pressure At First Visit")+
428   guides(col="none")
429
430
431
```