

MC884/MO444 - APRENDIZADO DE MÁQUINA

Experimentação com diversas Técnicas

Erik de Godoy Perillo - RA135582

Universidade Estadual de Campinas

24 de outubro de 2016

1 Introdução

O objetivo do trabalho era experimentar com diversas técnicas de *Machine Learning*: *K-nearest neighbours*, *SVM*, *Neural Networks*, *Random Forests*, *GBM*.

1.1 Implementação

A linguagem de implementação escolhida foi o R. Todo o código utilizado no relatório encontra-se na seção 4. Ao longo do documento, linhas do código serão citadas para referência no mesmo. A função `main` (linha 254 da seção 4) executa tudo que é requisitado no enunciado, mostrando os resultados.

2 Metodologia

O pré-processamento pedido é feito na função `pre_proc`, definida na linha 39 da seção 4. As colunas com 30% dos dados faltantes são removidas seguidas das linhas com dados faltantes. Os dados são normalizados com média 0 e desvio padrão 1.

É feita validação *k-fold* externa de 5 *folds* com validações internas de 3 *folds*. Para cada algoritmo, foi feita uma função para achar os melhores parâmetros e uma para a acurácia a cada *fold* externo. Todos os parâmetros requisitados para serem usados nas buscas são definidos a partir da linha 12 da seção 4.

Para todos os algoritmos, é mantida a maior acurácia. No final, todas as maiores acurácias pelo *k-fold* externo são reportadas na linha 380 da função `main`.

- **Knn**: Funções de busca de parâmetros e determinação de acurácia estão nas linhas 81 e 100, respectivamente, da seção 4. O PCA é implementado como pedido, mantendo 80% da variância.
- **SVM**: Funções de busca de parâmetros e determinação de acurácia nas linhas 111 e 126, respectivamente.
- **Neural Network**: Funções de busca de parâmetros e determinação de acurácia nas linhas 148 e 163, respectivamente.
- **Random Forests**: Funções de busca de parâmetros e determinação de acurácia nas linhas 183 e 197, respectivamente.
- **GBM**: Funções de busca de parâmetros e determinação de acurácia nas linhas 214 e 230, respectivamente.

saída relevante do código todo sendo executado pela `main` encontra-se na seção 5.

3 Resultados

Como se observa na saída da seção 5, os algoritmos que obtiveram melhores acurácias, em ordem decrescente, são: **SVM**, **Neural Network**, **GBM**, **Knn**, **Random Forest**. Vale notar, entretanto, que todos os algoritmos obtiveram resultados muito parecidos e então não pode-se tirar conclusões definitivas sobre qual o melhor.

4 Código-fonte

```
1 #packages
2 library(caret)
3 library(e1071)
4 library(class)
5 library(nnet)
6 library(randomForest)
7
8 #file path of data
9 data_filepath <- "~/random/secom/secom.data"
10 labels_filepath <- "~/random/secom/secom_labels.data"
11
12 #k for external k-fold
13 ext_k <- 5
14 #k for internal k-fold
15 inn_k <- 3
16
17 #knn parameters
18 knn_ks <- c(1, 5, 11, 15, 21, 25)
19
20 #svm parameters
21 svm_costs <- c(2^-5, 2^0, 2^5, 2^10)
22 svm_gammas <- c(2^-15, 2^-10, 2^-5, 2^0, 2^5)
23
24 #neural net parameters
25 nn_hidden_layer_sizes <- c(10, 20, 30, 40)
26
27 #random forest parameters
28 rf_nums_features <- c(10, 15, 20, 25)
29 rf_nums_trees <- c(100, 200, 300, 400, 500)
30
31 #gradient boosting machine
32 gbm_nums_trees <- c(30, 70, 100)
33 gbm_learning_rates <- c(0.1, 0.05)
34 gbm_depth <- 5
35
36 #wrapper for sprintf
37 printf <- function(...) cat(sprintf(...))
38
39 pre_proc <- function(x, y, empty_ratio_thr=0.3, scale=TRUE)
40 {
41   #filtering out columns with more holes than specified as limit
42   empty_num_thr <- floor(empty_ratio_thr*nrow(x))
43   col_filter <- apply(x, 2,
44     function(col) {sum(is.na(col)) <= empty_num_thr})
45   new_x <- x[, col_filter]
46
47   #filtering out rows with holes
48   row_filter <- apply(new_x, 1,
49     function(row) {!any(is.na(row))})
50   new_x <- new_x[row_filter, ]
51
52   new_y <- as.matrix(y[row_filter, 1])
53
54   #scaling x
55   if(scale)
56   {
57     new_x <- scale(new_x, center=TRUE, scale=TRUE)
58     #calling pre_proc again to filter out nans since scale produces nans
59     ret <- pre_proc(new_x, new_y, empty_ratio_thr, FALSE)
60     new_x <- ret$x
61     new_y <- ret$y
62   }
63
64   return(list("x"=new_x, "y"=new_y))
65 }
66
```

```

67 #gets minimum number of principal components to keep in order to conserve
68 #min_var of variance
69 pca_min_pcs <- function(pcs, min_var)
70 {
71   #getting k minimum number of components required for minimum variance
72   pcs_var_cumsum <- cumsum(pcs$sdev^2/sum(pcs$sdev^2))
73   min_pcs <- which(pcs_var_cumsum >= min_var)[1]
74   #printing result
75   #printf("\t-Number of components to keep %.2f%% variance: %d\n",
76   # min_var*100, min_pcs)
77
78   return(min_pcs)
79 }
80
81 knn_best_params <- function(x_train, y_train, pca_min_var=0.8, ks=knn_ks,
82   cross=inn_k)
83 {
84   #getting principal components maintaining minimum variance percentage
85   pcs <- prcomp(x_train, scale=FALSE)
86   min_pcs <- pca_min_pcs(pcs, pca_min_var)
87   #getting transformation matrix
88   transf_mat <- t(pcs$rotation[, 1:min_pcs])
89
90   #transforming data into k dimensions while keeping percentage of variance
91   x_train <- as.matrix(x_train) %*% t(transf_mat)
92
93   #performing grid search to find best k
94   control <- tune.control(sampling=c("cross"), cross=cross)
95   tuning <- tune.knn(x_train, y_train, k=ks, tunecontrol=control)
96
97   return(tuning$best.parameters)
98 }
99
100 knn_accuracy <- function(x_train, y_train, x_test, y_test, k)
101 {
102   #getting prediction
103   pred <- knn(x_train, x_test, y_train, k)
104
105   #getting accuracy
106   acc <- sum(pred == y_test)/length(y_test)
107
108   return(acc)
109 }
110
111 svm_best_params <- function(x_train, y_train, costs=svm_costs,
112   gammas=svm_gammas, cross=inn_k)
113 {
114   #converting to numeric values
115   y_train[y_train] = 1
116   y_train[!y_train] = 0
117
118   #getting best parameters
119   control <- tune.control(sampling=c("cross"), cross=cross)
120   tuning <- tune.svm(x_train, y_train, cost=costs, gamma=gammas,
121     tunecontrol=control)
122
123   return(tuning$best.parameters)
124 }
125
126 svm_accuracy <- function(x_train, y_train, x_test, y_test, cost, gamma)
127 {
128   #converting to numeric values
129   y_train[y_train] = 1
130   y_train[!y_train] = 0
131   y_test[y_test] = 1
132   y_test[!y_test] = 0
133
134   #getting model
135   model <- svm(x_train, y_train, scale=FALSE, cost=cost, gamma=gamma)

```

```

136
137 #getting prediction
138 pred <- predict(model, x_test, probability=TRUE)
139 pred[pred >= 0.5] = 1
140 pred[pred < 0.5] = 0
141
142 #getting accuracy
143 acc <- sum(pred == y_test)/length(y_test)
144
145 return(acc)
146 }
147
148 nn_best_params <- function(x_train, y_train,
149 hidden_layer_sizes=nn_hidden_layer_sizes, cross=inn_k)
150 {
151   #converting to numeric values
152   y_train[y_train] = 1
153   y_train[!y_train] = 0
154
155   #getting best parameters
156   control <- tune.control(sampling=c("cross"), cross=cross)
157   tuning <- tune.nnet(x_train, y_train, size=hidden_layer_sizes,
158     tunecontrol=control, MaxNWts=20000)
159
160   return(tuning$best.parameters)
161 }
162
163 nn_accuracy <- function(x_train, y_train, x_test, y_test, size)
164 {
165   #converting to numeric values
166   y_train[y_train] = 1
167   y_train[!y_train] = 0
168   y_test[y_test] = 1
169   y_test[!y_test] = 0
170
171   #getting model
172   model <- nnet(x_train, y_train, size=size, MaxNWts=20000)
173
174   #getting prediction
175   pred <- predict(model, x_test)
176
177   #getting accuracy
178   acc <- sum(pred == y_test)/length(y_test)
179
180   return(acc)
181 }
182
183 rf_best_params <- function(x_train, y_train, mtry=rf_nums_features,
184 ntree=rf_nums_trees, cross=inn_k)
185 {
186   #converting to categorical values
187   y_train <- as.factor(y_train)
188
189   #getting best parameters
190   control <- tune.control(sampling=c("cross"), cross=cross)
191   tuning <- tune.randomForest(x_train, y_train, mtry=mtry, ntree=ntree,
192     tunecontrol=control)
193
194   return(tuning$best.parameters)
195 }
196
197 rf_accuracy <- function(x_train, y_train, x_test, y_test, mtry, ntree)
198 {
199   #converting to categorical values
200   y_train <- as.factor(y_train)
201   y_test <- as.factor(y_test)
202
203   #getting model
204   model <- randomForest(x_train, y=y_train, xtest=x_test, ytest=y_test,

```

```

205     mtry=mtry, ntree=ntree)
206
207     #getting accuracy
208     conf <- ((model$confusion))
209     acc <- (conf[1, 1] + conf[2, 2])/sum(conf)
210
211     return(acc)
212 }
213
214 gbm_best_params <- function(x_train, y_train, depth=gbm_depth,
215 n_trees=gbm_nums_trees, shrinkage=gbm_learning_rates, cross=inn_k)
216 {
217     #converting to numeric values
218     y_train <- as.factor(y_train)
219
220     #getting best parameters
221     fitControl <- trainControl(method="repeatedcv", number=cross, repeats=1)
222     gbmGrid <- expand.grid(interaction.depth=depth, n.trees=n_trees,
223 shrinkage=shrinkage, n.minobsinnode=c(10))
224     tuning <- train(x_train, y_train, method = "gbm", trControl = fitControl,
225 verbose = FALSE, tuneGrid = gbmGrid)
226
227     return(tuning$bestTune)
228 }
229
230 gbm_accuracy <- function(x_train, y_train, x_test, y_test,
231 n_trees, depth, shrinkage, minobs)
232 {
233     #converting to numeric values
234     y_train[y_train] <- 1
235     y_train[!y_train] <- 0
236
237     #getting model
238     model <- gbm.fit(x_train, y_train, n.trees=n_trees, interaction.depth=depth,
239 shrinkage=shrinkage, n.minobsinnode=minobs, distribution="adaboost",
240 verbose=FALSE)
241
242     #getting prediction
243     pred <- predict(model, x_test, n.trees=n_trees, probability=TRUE)
244     pred[pred >= 0.5] = 1
245     pred[pred < 0.5] = 0
246
247     #getting accuracy
248     acc <- sum(pred == y_test)/length(y_test)
249
250     return(acc)
251 }
252
253 #main method for whole challenge
254 main <- function()
255 {
256     #reading data
257     x <- read.csv(data_filepath, header=FALSE, sep=" ")
258     x <- as.matrix(x)
259     y <- read.csv(labels_filepath, header=FALSE, sep=" ")
260     #transforming y into logical matrix
261     y <- y == 1
262     y <- as.matrix(y[, 1])
263
264     #pre-processing data
265     ret <- pre_proc(x, y)
266     x <- ret$x
267     y <- ret$y
268
269     knn_best_acc <- 0
270     svm_best_acc <- 0
271     nn_best_acc <- 0
272     rf_best_acc <- 0
273     gbm_best_acc <- 0

```

```

274 folds <- createFolds(y, k=ext_k)
275 i <- 0
276 for(fold in folds)
277 {
278   printf("on fold n. %d...\n", i)
279
280   #getting train/test folds
281   x_train <- x[-fold, ]
282   y_train <- as.matrix(y[-fold, ])
283   x_test <- x[fold, ]
284   y_test <- as.matrix(y[fold, ])
285
286   #KNN
287   printf("--- on KNN ---\n")
288   printf("\tselecting parameters... ")
289   knn_best <- knn_best_params(x_train, y_train)
290   knn_k <- knn_best$k
291   printf("done. k=%d\n", knn_k)
292   printf("\tgetting accuracy... ")
293   knn_acc <- knn_accuracy(x_train, y_train, x_test, y_test, knn_k)
294   printf("done. accuracy=%.6f", knn_acc)
295   if(knn_acc > knn_best_acc)
296   {
297     printf(" (best so far!)\n")
298     knn_best_acc <- knn_acc
299   }
300   printf("\n")
301
302   #SVM
303   printf("--- on SVM ---\n")
304   printf("\tselecting parameters... ")
305   svm_best <- svm_best_params(x_train, y_train)
306   svm_cost <- svm_best$cost
307   svm_gamma <- svm_best$gamma
308   printf("done. cost=%.6f, gamma=%.6f\n", svm_cost, svm_gamma)
309   printf("\tgetting accuracy... ")
310   svm_acc <- svm_accuracy(x_train, y_train, x_test, y_test,
311     svm_cost, svm_gamma)
312   printf("done. accuracy=%.6f", svm_acc)
313   if(svm_acc > svm_best_acc)
314   {
315     printf(" (best so far!)\n")
316     svm_best_acc <- svm_acc
317   }
318   printf("\n")
319
320   #NEURAL NETWORK
321   printf("--- on NEURAL NETWORK ---\n")
322   printf("\tselecting parameters... ")
323   nn_best <- nn_best_params(x_train, y_train)
324   nn_size <- nn_best$size
325   printf("done. size=%d\n", nn_size)
326   printf("\tgetting accuracy... ")
327   nn_acc <- nn_accuracy(x_train, y_train, x_test, y_test, nn_size)
328   printf("done. accuracy=%.6f", nn_acc)
329   if(nn_acc > nn_best_acc)
330   {
331     printf(" (best so far!)\n")
332     nn_best_acc <- nn_acc
333   }
334   printf("\n")
335
336   #RANDOM FOREST
337   printf("--- on RANDOM FOREST ---\n")
338   printf("\tselecting parameters... ")
339   rf_best <- rf_best_params(x_train, y_train)
340   rf_mtry <- rf_best$mtry
341   rf_ntree <- rf_best$ntree

```

```

343     printf("done. mtry=%d, ntree=%d\n", rf_mtry, rf_ntree)
344     printf("\tgetting accuracy... ")
345     rf_acc <- rf_accuracy(x_train, y_train, x_test, y_test,
346                          rf_mtry, rf_ntree)
347     printf("done. accuracy=%.6f", rf_acc)
348     if(rf_acc > rf_best_acc)
349     {
350         printf(" (best so far!)" )
351         rf_best_acc <- rf_acc
352     }
353     printf("\n")
354
355     #GBM
356     printf("--- on GBM ---\n")
357     printf("\tselecting parameters... ")
358     gbm_best <- gbm_best_params(x_train, y_train)
359     gbm_n_trees <- gbm_best$n.trees
360     gbm_depth <- gbm_best$interaction.depth
361     gbm_shrinkage <- gbm_best$shrinkage
362     gbm_minobs <- gbm_best$n.minobsinnode
363     printf("done.\n")
364     print(gbm_best)
365     printf("\tgetting accuracy... ")
366     gbm_acc <- gbm_accuracy(x_train, y_train, x_test, y_test,
367                            gbm_n_trees, gbm_depth, gbm_shrinkage, gbm_minobs)
368     printf("done. accuracy=%.6f", gbm_acc)
369     if(gbm_acc > gbm_best_acc)
370     {
371         printf(" (best so far!)" )
372         gbm_best_acc <- gbm_acc
373     }
374     printf("\n")
375
376     i <- i + 1
377     printf("\n")
378 }
379
380 printf("FINAL RESULTS:\n")
381 printf("\tknn best accuracy: %.6f\n", knn_best_acc)
382 printf("\tsvm best accuracy: %.6f\n", svm_best_acc)
383 printf("\tnn best accuracy: %.6f\n", nn_best_acc)
384 printf("\trf best accuracy: %.6f\n", rf_best_acc)
385 printf("\tgbm best accuracy: %.6f\n", gbm_best_acc)
386 }

```


5 Saída do código

```
1 FINAL RESULTS:  
2     knn best accuracy: 0.933333  
3     svm best accuracy: 0.933674  
4     nn best accuracy: 0.933674  
5     rf best accuracy: 0.931035  
6     gbm best accuracy: 0.933333
```