

MC884/MO444 - APRENDIZADO DE MÁQUINA

Regressão Logística e LDA

Erik de Godoy Perillo - RA135582

Universidade Estadual de Campinas

23 de setembro de 2016

1 Introdução

O objetivo do trabalho era experimentar com as técnicas de PCA¹, *Logistic Regression* e LDA².

1.1 Implementação

A linguagem de implementação escolhida foi o R. Todo o código utilizado no relatório encontra-se na seção 3. Ao longo do documento, linhas do código serão citadas para referência no mesmo. A função `main` (linha 59 da seção 3) executa todos os itens em ordem, mostrando os resultados.

2 Itens

1. *faça o PCA dos dados (sem a última coluna). Se você quiser que os dados transformados tenham 80% da variância original, quantas dimensões do PCA vc precisa manter? Gere os dados transformados mantendo 80% da variância.*

O PCA é calculado na linha 69 do código da seção 3. Os dados são normalizados antes da função `prcomp` ser chamada, pois já precisaremos deles normalizados no momento em que gerarmos os dados nas novas dimensões. A função `pca_min_pcs` da linha 19 da seção 3 contém o código necessário para determinar o número mínimo de componentes dada uma variância. Para o valor de 80%, é necessário manter 13 das bases com maior variância dos componentes principais. A matriz de transformação dos dados antigos para os nas novas 13 dimensões é obtida na linha 78 da seção 3. Os dados novos são gerados na linha 80.

2. *Treine uma regressão logística no conjunto de treino dos dados originais e nos dados transformados. Qual a taxa de acerto no conjunto de teste nas 2 condições (sem e com PCA)?*

A função que calcula a regressão logística e sua acurácia é dada na função `logit_reg` da linha 32 da seção 3. Os dados de treino e teste são selecionados entre as linhas 84 e 91. Como resultado, obteve-se que a acurácia do modelo obtido pela regressão com todos os dados foi menor que a obtida com as dimensões que mantinham 80% da variância. Os valores foram, respectivamente, 65.58% e 75.72%.

3. *Treine o LDA nos conjuntos de treino com e sem PCA e teste nos respectivos conjuntos de testes. Qual a acurácia nas 2 condições?*

A função que calcula a LDA e sua acurácia é dada na função `lda_reg` da linha 45 da seção 3. Como resultado, obteve-se que a acurácia do modelo obtido pela regressão com todos os dados foi menor que a obtida com as dimensões que mantinham 80% da variância, assim como no item 2. Os valores foram, respectivamente, 67.75% e 78.62%.

4. *Qual a melhor combinação de classificador e PCA ou não?*

A combinação que mostrou o melhor resultado foi o uso de PCA com LDA. O resultado é interessante, ainda mais que na regressão logística, sob o uso de todas as dimensões, o algoritmo demonstrou mensagens de aviso do tipo `algorithm did not converge`. Uma pesquisa para os motivos de tal resultado sugere [1] que isso pode ser pelo fato de o modelo estar

¹do inglês: *Principal Component Analysis*

²do inglês: *Linear Discriminant Analysis*

“perfeito demais” que, devido a algum detalhe de implementação da função em R, faz com que alguns parâmetros fiquem com valores muito pequenos/grandes, sendo assim difíceis de serem representados com precisão pelo computador e então gerando resultados piores que os com PCA.

3 Apêndice: código-fonte

```
1 #package for lda
2 library(MASS)
3
4 #default values:
5 #file path of data
6 data_filepath <- "data1.csv"
7 #minimum variance required
8 min_var <- 0.80
9 #number of lines to use in training
10 train_n_lines <- 200
11 #number of lines to use in test
12 test_n_lines <- 276
13
14 #wrapper for sprintf
15 printf <- function(...) cat(sprintf(...))
16
17 #gets minimum number of principal components to keep in order to conserve
18 #min_var of variance
19 pca_min_pcs <- function(pcs, min_var)
20 {
21   #getting k minimum number of components required for minimum variance
22   pcs_var_cumsum <- cumsum(pcs$sdev^2/sum(pcs$sdev^2))
23   min_pcs <- which(pcs_var_cumsum >= min_var)[1]
24   #printing result
25   printf("\t-Number of components to keep %.2f%% variance: %d\n",
26         min_var*100, min_pcs)
27
28   return(min_pcs)
29 }
30
31 #calculates logistic regression and displays accuracy
32 logit_reg <- function(x_train, y_train, x_test, y_test)
33 {
34   #computing logistic regression
35   lr <- glm(y_train ~ ., data=x_train, family=binomial(link="logit"))
36   #getting predictions
37   pred <- as.matrix(predict(lr, x_test)) >= 1
38   #getting score
39   score = sum(pred == y_test)/length(y_test)
40   #printing accuracy
41   printf("accuracy: %.2f%%\n", score*100)
42 }
43
44 #calculates LDA and displays accuracy
45 lda_reg <- function(x_train, y_train, x_test, y_test)
46 {
47   #computing lda
48   ldar <- lda(y_train ~ ., data=x_train)
49   #getting predictions
50   pred <- predict(ldar, x_test, prior=ldar$prior)
51   pred <- pred$posterior[, 2] > pred$posterior[, 1]
52   #getting score
53   score = sum(pred == y_test)/length(y_test)
54   #printing accuracy
55   printf("accuracy: %.2f%%\n", score*100)
56 }
57
```

```

58 #main method for whole challenge
59 main <- function()
60 {
61   #reading data
62   data <- read.csv(data_filepath)
63   x <- data[, 1:ncol(data)-1]
64
65   #scaling data prior to pca. we would have to do it anyway...
66   x <- scale(x)
67
68   #getting principal components
69   pcs <- prcomp(x, scale=FALSE)
70
71   #1. Faça o PCA dos dados (sem a ultima coluna).
72   #Se voce quiser que os dados transformados tenham 80% da variancia original,
73   #quantas dimensoes do PCA vc precisa manter?
74   #Gere os dados transformados mantendo 80% da variancia.
75   printf("Item 1:\n")
76   min_pcs <- pca_min_pcs(pcs, min_var)
77   #getting transformation matrix
78   transf_mat <- t(pcs$rotation[, 1:min_pcs])
79   #transforming data into k dimensions while keeping percentage of variance
80   transf_x <- as.matrix(x) %*% t(transf_mat)
81
82   #preparing data for regression
83   #train data
84   y <- as.matrix(data[, ncol(data)])
85   y_train <- y[1:train_n_lines, ]
86   x_full_var_train <- x[1:train_n_lines, ]
87   x_part_var_train <- transf_x[1:train_n_lines, ]
88   #test data
89   x_full_var_test <- x[(train_n_lines+1):nrow(x), ]
90   x_part_var_test <- as.matrix(x_full_var_test) %*% t(transf_mat)
91   y_test <- y[(train_n_lines+1):nrow(y), ]
92
93   #2. Treine uma regressao logistica no conjunto de treino dos dados originais
94   #e nos dados transformados.
95   #Qual a taxa de acerto no conjunto de teste nas 2 condicoes (sem e com PCA)?
96   printf("\nItem 2:\n")
97   #logistic regression on all dimensions
98   printf("\tAll dimensions: ")
99   logit_reg(as.data.frame(x_full_var_train), y_train,
100    as.data.frame(x_full_var_test), y_test)
101   #logistic regression on k principal components
102   printf("\tFirst %d dimensions from PCA: ", min_pcs)
103   logit_reg(as.data.frame(x_part_var_train), y_train,
104    as.data.frame(x_part_var_test), y_test)
105
106   #3. Treine o LDA nos conjuntos de treino com e sem PCA e teste nos
107   #respectivos conjuntos de testes. Qual a acuracia nas 2 condicoes?
108   printf("\nItem 3:\n")
109   #lda on all dimensions
110   printf("\tAll dimensions: ")
111   lda_reg(as.data.frame(x_full_var_train), y_train,
112    as.data.frame(x_full_var_test), y_test)
113   #lda on k principal components
114   printf("\tFirst %d dimensions from PCA: ", min_pcs)
115   lda_reg(as.data.frame(x_part_var_train), y_train,
116    as.data.frame(x_part_var_test), y_test)
117
118   printf("\n")
119 }

```

Referências

- [1] *Logistic regression model does not converge*. 2010. URL: <http://stats.stackexchange.com/questions/5354/logistic-regression-model-does-not-converge> (acesso em 19/04/2016).