

Support Vector Machine: Uma aplicação R

Kim Silva

20/05/2022

Vamos aplicar o modelo SVM no conjunto de dados sobre diabetes disponível no R com o comando `data("diabetes392")` do pacote `HTLR`. Este conjunto de dados é originalmente do Instituto Nacional de Diabetes e Doenças Digestivas e Renais. O objetivo do conjunto de dados é prever diagnosticamente se um paciente tem ou não diabetes, com base em certas medidas de diagnóstico incluídas no conjunto de dados. Várias restrições foram colocadas na seleção dessas instâncias de um banco de dados maior. Em particular, todos os pacientes aqui são mulheres com pelo menos 21 anos de ascendência indígena Pima. Diferente da versão original do UCI, o conjunto de dados foi pré-processado para que as linhas com valores ausentes sejam removidas e os recursos sejam dimensionados.

Para mais informações sobre a base de dados consultar no link: <https://search.r-project.org/CRAN/refmans/HTLR/html/diabetes392.html>

```
#install.packages('HTLR')
library(HTLR)
```

```
## Warning: package 'HTLR' was built under R version 4.1.3
```

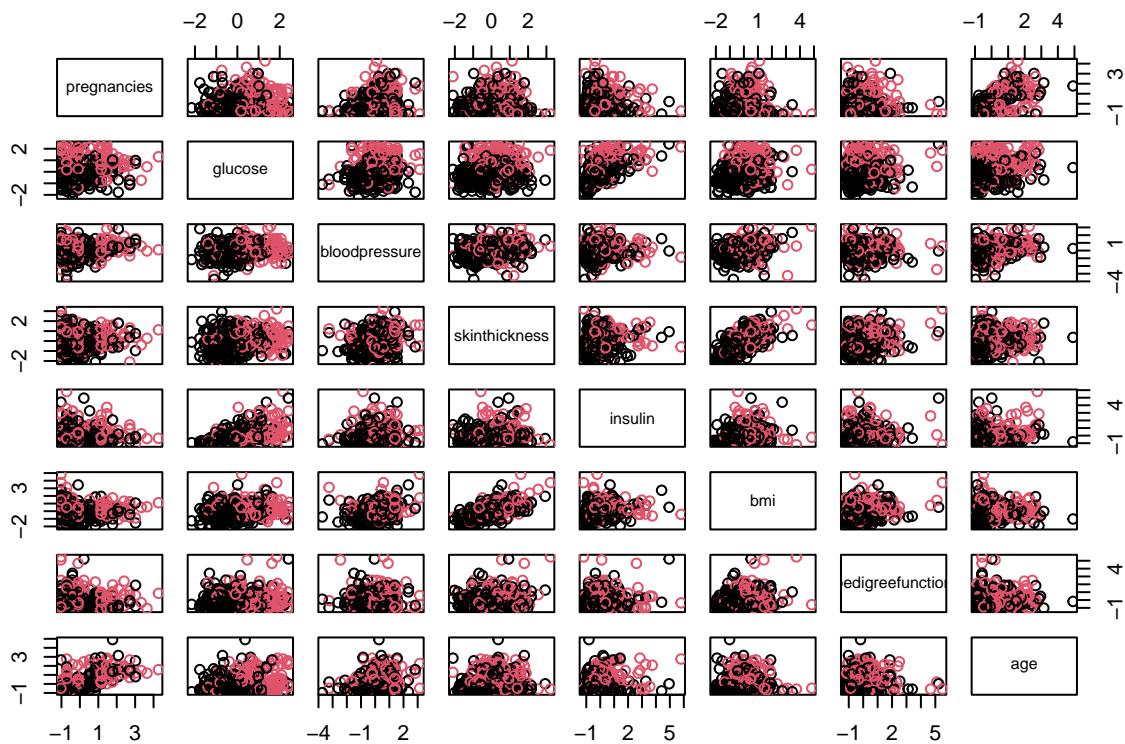
```
data("diabetes392")
dados = as.data.frame(cbind(diabetes392$X,diabetes392$y))
dados$V9 <- as.factor(dados$V9)
summary(dados)
```

```
##      pregnancies      glucose      bloodpressure      skinthickness
## Min.      :-1.0279   Min.      :-2.1590   Min.      :-3.73423   Min.      :-2.10579
## 1st Qu.   :-0.7165   1st Qu.   :-0.7656   1st Qu.   :-0.69328   1st Qu.   :-0.77454
## Median    :-0.4051   Median    :-0.1175   Median    :-0.05308   Median    :-0.01383
## Mean      : 0.0000   Mean      : 0.0000   Mean      : 0.00000   Mean      : 0.00000
## 3rd Qu.   : 0.5290   3rd Qu.   : 0.6601   3rd Qu.   : 0.58712   3rd Qu.   : 0.74689
## Max.      : 4.2657   Max.      : 2.4423   Max.      : 3.14792   Max.      : 3.21921
##      insulin      bmi      pedigreefunction      age
## Min.      :-1.1953   Min.      :-2.11823   Min.      :-1.2679   Min.      :-0.9671
## 1st Qu.   :-0.6673   1st Qu.   :-0.66683   1st Qu.   :-0.7332   1st Qu.   :-0.7710
## Median    :-0.2571   Median      : 0.01619   Median    :-0.2129   Median    :-0.3789
## Mean      : 0.0000   Mean      : 0.00000   Mean      : 0.0000   Mean      : 0.0000
## 3rd Qu.   : 0.2856   3rd Qu.   : 0.57114   3rd Qu.   : 0.4746   3rd Qu.   : 0.5034
## Max.      : 5.8056   Max.      : 4.83999   Max.      : 5.4906   Max.      : 4.9148
## V9
## 1:262
## 2:130
##
##
```

```
##
##
```

Como podemos notar há um desbalanceamento na variável de interesse. A base de dados, conta com 130 pacientes diabéticas e 262 pacientes não diabéticas. Desse modo, o indicado seria utilizar algum método para correção do desbalanceamento das classes.

```
library(graphics)
pairs(dados[, -9], col = dados$V9)
```



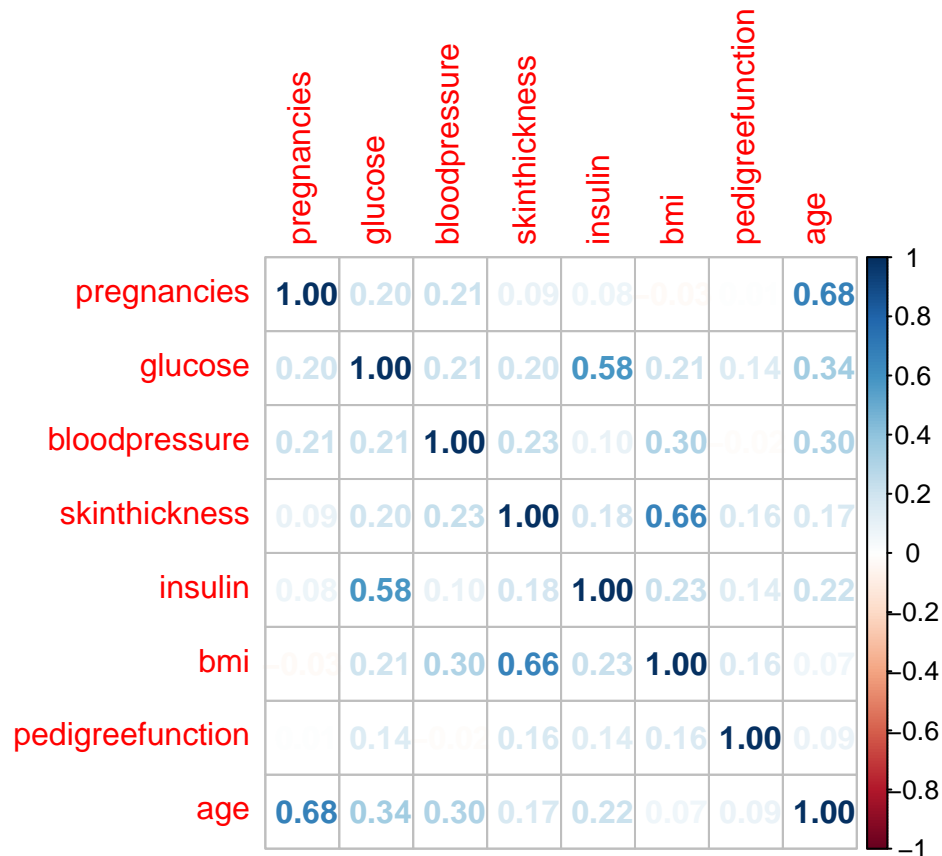
Vemos na imagem acima, o gráfico de dispersão. As observações das pacientes, estão coloridas de acordo com a condição diabética ou não diabética. Podemos perceber, de acordo com o comportamento dos dados, que não conseguimos ver uma distinção clara para classificarmos diabéticas e não diabéticas. Nesse contexto, o modelo SVM pode ser uma boa opção de classificador.

```
library(corrplot)
```

```
## Warning: package 'corrplot' was built under R version 4.1.1
```

```
## corrplot 0.90 loaded
```

```
corrplot(cor(dados[, -9]), method = "number")
```



A correlação entre as variáveis age e pregnancies, skinthickness e bmi apresenta correlação maior que 60%. Desse modo, iremos construir os modelos SVM com todas as variáveis e sem as variáveis age e skinthickness pois apresentam alta correlação com as variáveis pregnancies e bmi, respectivamente.

Support Vector Machine

```
library(e1071)
```

```
## Warning: package 'e1071' was built under R version 4.1.1
```

```
tune.out = tune(svm ,V9~.,data=dados,
               ranges=list(kernel=c("linear", "polynomial",
                                   "radial", "sigmoid"))))
summary(tune.out)
```

```
##
## Parameter tuning of 'svm':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
## kernel
## linear
```

```
##
## - best performance: 0.2242949
##
## - Detailed performance results:
##      kernel      error dispersion
## 1      linear 0.2242949 0.07725619
## 2 polynomial 0.2522436 0.07509171
## 3      radial 0.2446795 0.10957104
## 4      sigmoid 0.2757051 0.08004836
```

A função `tune` realiza o método 10-fold cross validation e como resultado o melhor kernel foi o `linear`. O mesmo método poderia ser utilizado para realizar a escolha do parâmetro `cost`. No entanto optamos por utilizar o parâmetro `cost` igual a cinco.

Iremos utilizar o método 10-fold cross validation para mensurar o desempenho dos modelos.

```
# 10-fold cross validation:
ACC=numeric(10)
MCC=numeric(10)
SEN=numeric(10)
VPP=numeric(10)
F1=numeric(10)
BS=numeric(10)

set.seed(42553)
df = dados[sample(1:nrow(dados)),]
a = round(as.numeric(quantile(row(dados)[,1], probs = seq(0,1,0.1))))
for(i in 1:10){
  # Divisao em treinamento (70%) e teste (30%):
  dados.trein= df[a[i]:a[i+1], ]
  dados.test=df[-(a[i]:a[i+1]),]
  # SVM:
  mod = svm(V9~.,data = dados.trein, kernel = "linear", cost = 5)
  # Desempenho preditivo (teste):
  pred = predict(mod, dados.test)
  mc=table(predict = pred, truth = dados.test$V9)
  VN=mc[1,1]; FP=mc[2,1]; FN=mc[1,2]; VP=mc[2,2]
  ACC[i]=(VP+VN)/(VP+VN+FN+FP)
  MCC[i]=(VP*VN-FP*FN)/(sqrt(VP+FP)*sqrt(VP+FN)*sqrt(VN+FP)*sqrt(VN+FN))
  SEN[i]=VP/(VP+FN)
  VPP[i]=VP/(VP+FP)
  F1[i]=2*VPP[i]*SEN[i]/(VPP[i]+SEN[i])
}
rsvm = c(mean(ACC), mean(MCC), mean(SEN), mean(VPP), mean(F1))
rsvm
```

```
## [1] 0.7297559 0.3732059 0.5174623 0.6145804 0.5454167
```

```
# 10-fold cross validation:
ACC=numeric(10)
MCC=numeric(10)
SEN=numeric(10)
VPP=numeric(10)
F1=numeric(10)
```

```

BS=numeric(10)

for(i in 1:10){
  # Divisao em treinamento (70%) e teste (30%):
  dados.trein = df[a[i]:a[i+1],-c(4,8) ]
  dados.test = df[-(a[i]:a[i+1]),-c(4,8)]
  # SVM:
  mod = svm(V9~.,data = dados.trein, kernel = "linear", cost = 5)
  # Desempenho preditivo (teste):
  pred = predict(mod, dados.test)
  mc=table(predict = pred, truth = dados.test$V9)
  VN=mc[1,1]; FP=mc[2,1]; FN=mc[1,2]; VP=mc[2,2]
  ACC[i]=(VP+VN)/(VP+VN+FN+FP)
  MCC[i]=(VP*VN-FP*FN)/(sqrt(VP+FP)*sqrt(VP+FN)*sqrt(VN+FP)*sqrt(VN+FN))
  SEN[i]=VP/(VP+FN)
  VPP[i]=VP/(VP+FP)
  F1[i]=2*VPP[i]*SEN[i]/(VPP[i]+SEN[i])
}
rsvm1 = c(mean(ACC), mean(MCC), mean(SEN), mean(VPP), mean(F1))
rsvm1

```

```
## [1] 0.7343054 0.3753489 0.4717926 0.6515379 0.5225558
```

```

medidas <- data.frame(Medidas = c('ACC', 'MCC', 'SEN', 'VPP', 'F1'),
  SVM = rsvm, SVM_2 = rsvm1)

library(knitr)
kable(medidas)

```

Medidas	SVM	SVM_2
ACC	0.7297559	0.7343054
MCC	0.3732059	0.3753489
SEN	0.5174623	0.4717926
VPP	0.6145804	0.6515379
F1	0.5454167	0.5225558

A tabela acima mostra o desempenho dos modelos. O modelo SVM_2, que remove as variáveis **age** e **skinthickness**, foi o modelo com maiores Acurácia (ACC), Coeficiente de Correlação de Matthews (MCC) e Valor Preditivo Positivo (VPP). Por outro lado, o modelo ajustado com todas as variáveis tem uma maiores sensibilidade (SEN) e F1-scores. Desse modo, concluímos que para esta base de dados, ajustar o SVM com a remoção das variáveis **age** e **skinthickness** é melhor para a performance do modelo.