

Feature Selection

Karen Mazidi

Look for correlations in Pima data

The `findCorrelation()` function suggests that we could remove column 6, mass, because it correlates with triceps. And that we could remove column 2, glucose, because it correlates with insulin.

```
library(caret)

## Warning: package 'caret' was built under R version 3.4.3
## Loading required package: lattice
## Loading required package: ggplot2
library(mlbench)
data("PimaIndiansDiabetes2")
df <- PimaIndiansDiabetes2[complete.cases(PimaIndiansDiabetes2[]),]
corMatrix <- cor(df[,1:7])
findCorrelation(corMatrix, cutoff=0.5, verbose=TRUE)

## Compare row 6 and column 4 with corr 0.664
## Means: 0.265 vs 0.187 so flagging column 6
## Compare row 2 and column 5 with corr 0.581
## Means: 0.266 vs 0.161 so flagging column 2
## All correlations <= 0.5

## [1] 6 2
```

Remove the highly correlated columns

```
df <- df[, -c(2,6)]
```

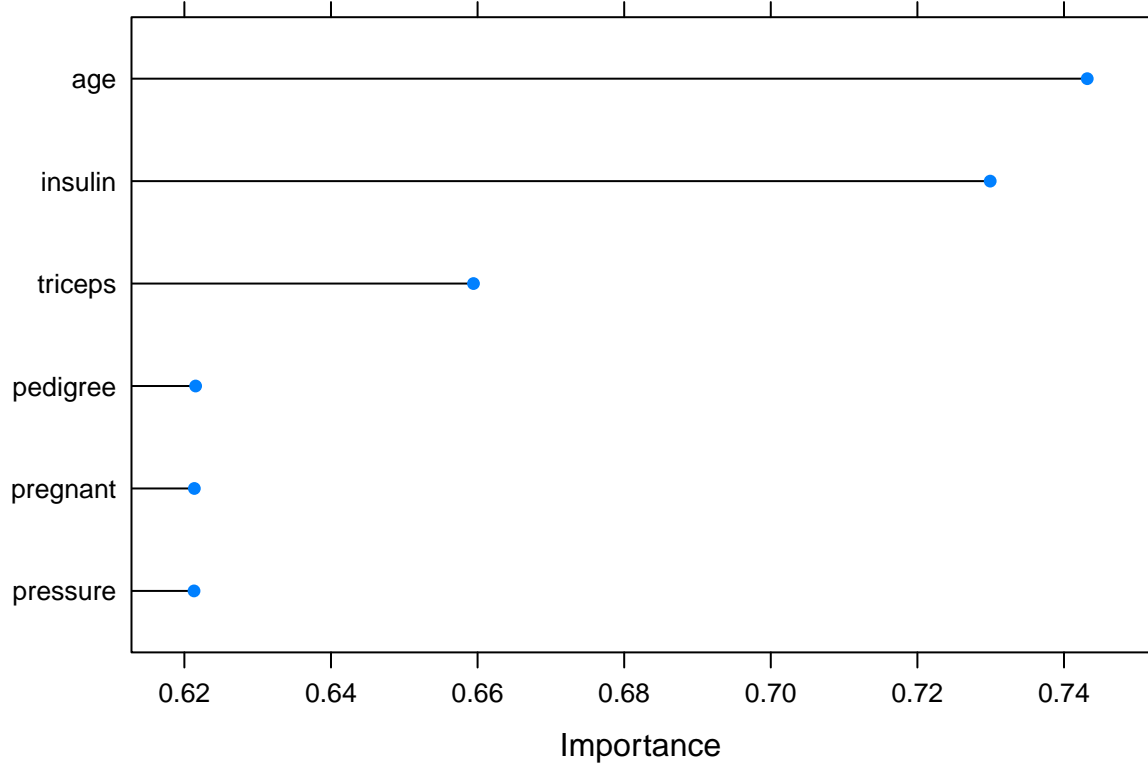
Rank features

The `varImp()` function ranks variables by importance. It requires a model which we trained on method `knn`, using control parameters stored in variable `ctrl`.

```
ctrl <- trainControl(method="repeatedcv", repeats=5)
model <- train(diabetes~., data=df, method="knn", preProcess="scale", trControl=ctrl)
importance <- varImp(model, scale=FALSE)
importance
```

```
## ROC curve variable importance
##
##      Importance
## age          0.7432
## insulin      0.7299
## triceps      0.6594
## pedigree     0.6215
## pregnant     0.6214
## pressure     0.6213
```

```
plot(importance)
```



###

Recursive feature selection

We start with the data set including all columns.

```
df <- PimaIndiansDiabetes2[complete.cases(PimaIndiansDiabetes2[]),]
ctrl <- rfeControl(functions=rfFuncs, method="cv", number=10)
rfe_out <- rfe(df[,1:7], df[,8], sizes=c(1:7), rfeControl=ctrl)
rfe_out
```

```
##
## Recursive feature selection
##
## Outer resampling method: Cross-Validated (10 fold)
##
## Resampling performance over subset size:
##
## Variables  RMSE Rsquared  MAE RMSESD RsquaredSD  MAESD Selected
##          1 7.399   0.4865 5.232  1.312    0.1649 0.7959
##          2 7.523   0.4663 5.434  1.403    0.1841 0.9263
##          3 7.471   0.4697 5.372  1.395    0.1788 1.0034
##          4 7.286   0.4963 5.256  1.247    0.1592 0.8430
##          5 7.294   0.5007 5.298  1.242    0.1518 0.8417
##          6 7.132   0.5114 5.133  1.286    0.1522 0.8548
##          7 7.030   0.5293 5.058  1.253    0.1606 0.7969      *
##
## The top 5 variables (out of 7):
##    pregnant, glucose, insulin, triceps, mass
```