

ML_With_Iris.R

dieudonneouedraogo

Sun Nov 11 20:07:01 2018

```
#install.packages("caret")
#install.packages("caret", dependencies=c("Depends", "Suggests"))
library(caret)

## Loading required package: lattice
## Loading required package: ggplot2

data(iris)
# rename the dataset
dataset <- iris
validation_index <- createDataPartition(dataset$Species, p=0.80, list=FALSE)
# select 20% of the data for validation
validation <- dataset[-validation_index,]
# use the remaining 80% of data to training and testing the models
dataset <- dataset[validation_index,]

# dimensions of dataset
dim(dataset)

## [1] 120 5

# list types for each attribute
sapply(dataset, class)

## Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## "numeric" "numeric" "numeric" "numeric" "factor"

# list the levels for the class
levels(dataset$Species)

## [1] "setosa" "versicolor" "virginica"

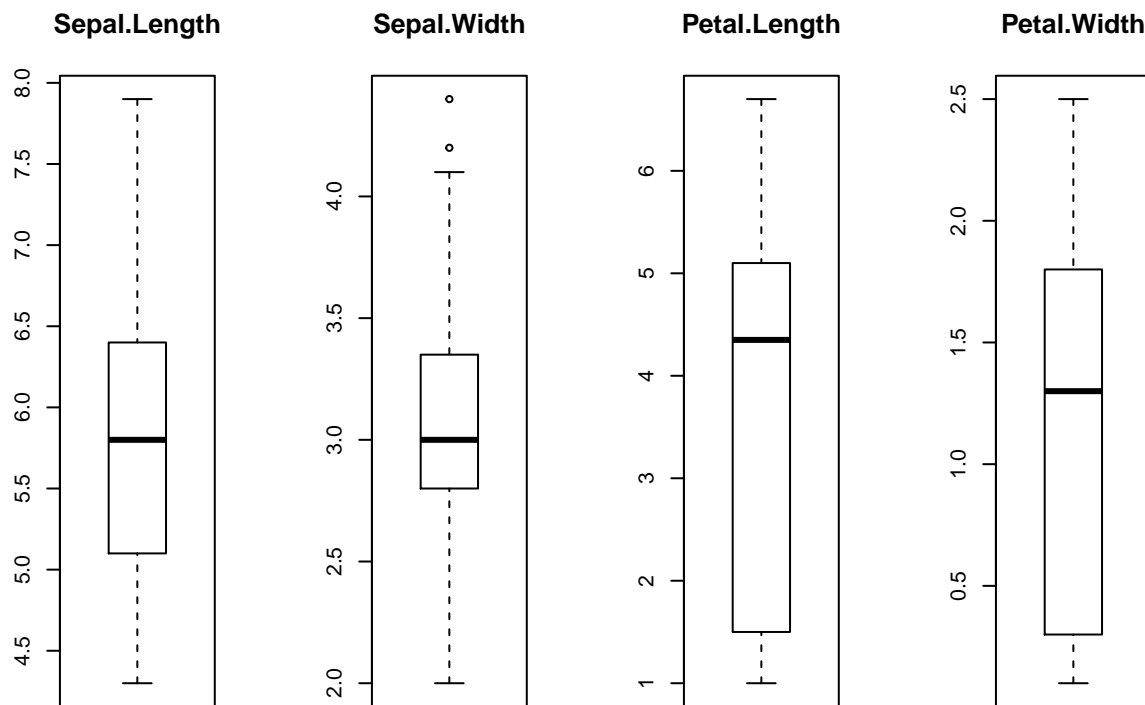
# summarize the class distribution
percentage <- prop.table(table(dataset$Species)) * 100
cbind(freq=table(dataset$Species), percentage=percentage)

##          freq percentage
## setosa      40  33.33333
## versicolor  40  33.33333
## virginica   40  33.33333

# split input and output
x <- dataset[,1:4]
y <- dataset[,5]

# boxplot for each attribute on one image
par(mfrow=c(1,4))
for(i in 1:4) {
  boxplot(x[,i], main=names(iris)[i])
}
```

3



```
# barplot for class breakdown
plot(y)
# scatterplot matrix
featurePlot(x=x, y=y, plot="ellipse")
# box and whisker plots for each attribute
featurePlot(x=x, y=y, plot="box")

# density plots for each attribute by class value
scales <- list(x=list(relation="free"), y=list(relation="free"))
featurePlot(x=x, y=y, plot="density", scales=scales)

# Run algorithms using 10-fold cross validation
control <- trainControl(method="cv", number=10)
metric <- "Accuracy"

# a) linear algorithms
set.seed(7)
fit.lda <- train(Species~., data=dataset, method="lda", metric=metric, trControl=control)
# b) nonlinear algorithms
# CART
set.seed(7)
fit.cart <- train(Species~., data=dataset, method="rpart", metric=metric, trControl=control)
# kNN
set.seed(7)
fit.knn <- train(Species~., data=dataset, method="knn", metric=metric, trControl=control)
# c) advanced algorithms
# SVM
set.seed(7)
```

```

fit.svm <- train(Species~., data=dataset, method="svmRadial", metric=metric, trControl=control)
# Random Forest
set.seed(7)
fit.rf <- train(Species~., data=dataset, method="rf", metric=metric, trControl=control)

# summarize accuracy of models
results <- resamples(list(lda=fit.lda, cart=fit.cart, knn=fit.knn, svm=fit.svm, rf=fit.rf))
summary(results)

##
## Call:
## summary.resamples(object = results)
##
## Models: lda, cart, knn, svm, rf
## Number of resamples: 10
##
## Accuracy
##      Min.   1st Qu.   Median     Mean 3rd Qu.  Max. NA's
## lda  0.9166667 1.0000000 1.0000000 0.9833333      1      1      0
## cart 0.8333333 0.9166667 1.0000000 0.9500000      1      1      0
## knn  0.8333333 0.9166667 0.9583333 0.9500000      1      1      0
## svm  0.7500000 0.9375000 1.0000000 0.9583333      1      1      0
## rf   0.8333333 0.9166667 1.0000000 0.9500000      1      1      0
##
## Kappa
##      Min. 1st Qu. Median     Mean 3rd Qu.  Max. NA's
## lda  0.875 1.00000 1.0000 0.9750      1      1      0
## cart 0.750 0.87500 1.0000 0.9250      1      1      0
## knn  0.750 0.87500 0.9375 0.9250      1      1      0
## svm  0.625 0.90625 1.0000 0.9375      1      1      0
## rf   0.750 0.87500 1.0000 0.9250      1      1      0

# compare accuracy of models
dotplot(results)

# summarize Best Model
print(fit.lda)

## Linear Discriminant Analysis
##
## 120 samples
## 4 predictor
## 3 classes: 'setosa', 'versicolor', 'virginica'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 108, 108, 108, 108, 108, 108, ...
## Resampling results:
##
## Accuracy  Kappa
## 0.9833333 0.975

# estimate skill of LDA on the validation dataset
predictions <- predict(fit.lda, validation)
confusionMatrix(predictions, validation$Species)

```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  setosa versicolor virginica
##   setosa      10         0         0
##   versicolor   0         10        1
##   virginica    0         0         9
##
## Overall Statistics
##
##           Accuracy : 0.9667
##           95% CI : (0.8278, 0.9992)
##   No Information Rate : 0.3333
##   P-Value [Acc > NIR] : 2.963e-13
##
##           Kappa : 0.95
##   McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: setosa Class: versicolor Class: virginica
## Sensitivity           1.0000           1.0000           0.9000
## Specificity           1.0000           0.9500           1.0000
## Pos Pred Value        1.0000           0.9091           1.0000
## Neg Pred Value        1.0000           1.0000           0.9524
## Prevalence            0.3333           0.3333           0.3333
## Detection Rate        0.3333           0.3333           0.3000
## Detection Prevalence  0.3333           0.3667           0.3000
## Balanced Accuracy     1.0000           0.9750           0.9500

```

