

ReadMe

Machine Learning Practice

```
#install.packages("caret")  
library(caret)
```

```
## Warning: package 'caret' was built under R version 4.0.5
```

```
## Loading required package: lattice
```

```
## Loading required package: ggplot2
```

```
library(ggplot2)
```

```
# attach the iris dataset to the environment  
data(iris)  
# rename the dataset  
dataset <- iris
```

```
# create a list of 80% of the rows in the original dataset we can use for training  
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,]
```

```
# We can get a quick idea of how many instances (rows) and how many attributes (columns) the data contains
```

```
# 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"
```

```
# take a peek at the first 5 rows of the data  
head(dataset)
```

```
## Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1          5.1          3.5          1.4          0.2 setosa
## 2          4.9          3.0          1.4          0.2 setosa
## 4          4.6          3.1          1.5          0.2 setosa
## 5          5.0          3.6          1.4          0.2 setosa
## 6          5.4          3.9          1.7          0.4 setosa
## 7          4.6          3.4          1.4          0.3 setosa
```

```
# 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
```

```
summary(dataset)
```

```
## Sepal.Length Sepal.Width Petal.Length Petal.Width
## Min. :4.300 Min. :2.200 Min. :1.000 Min. :0.100
## 1st Qu.:5.100 1st Qu.:2.800 1st Qu.:1.500 1st Qu.:0.300
## Median :5.800 Median :3.000 Median :4.500 Median :1.400
## Mean :5.848 Mean :3.069 Mean :3.772 Mean :1.213
## 3rd Qu.:6.400 3rd Qu.:3.300 3rd Qu.:5.100 3rd Qu.:1.800
## Max. :7.900 Max. :4.400 Max. :6.700 Max. :2.500
## Species
## setosa :40
## versicolor:40
## virginica :40
##
##
##
```

Testing and Machine Learning bit

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

Let's evaluate 5 different algorithms:

Linear Discriminant Analysis (LDA) Classification and Regression Trees (CART). k-Nearest Neighbors (kNN). Support Vector Machines (SVM) with a linear kernel. Random Forest (RF)

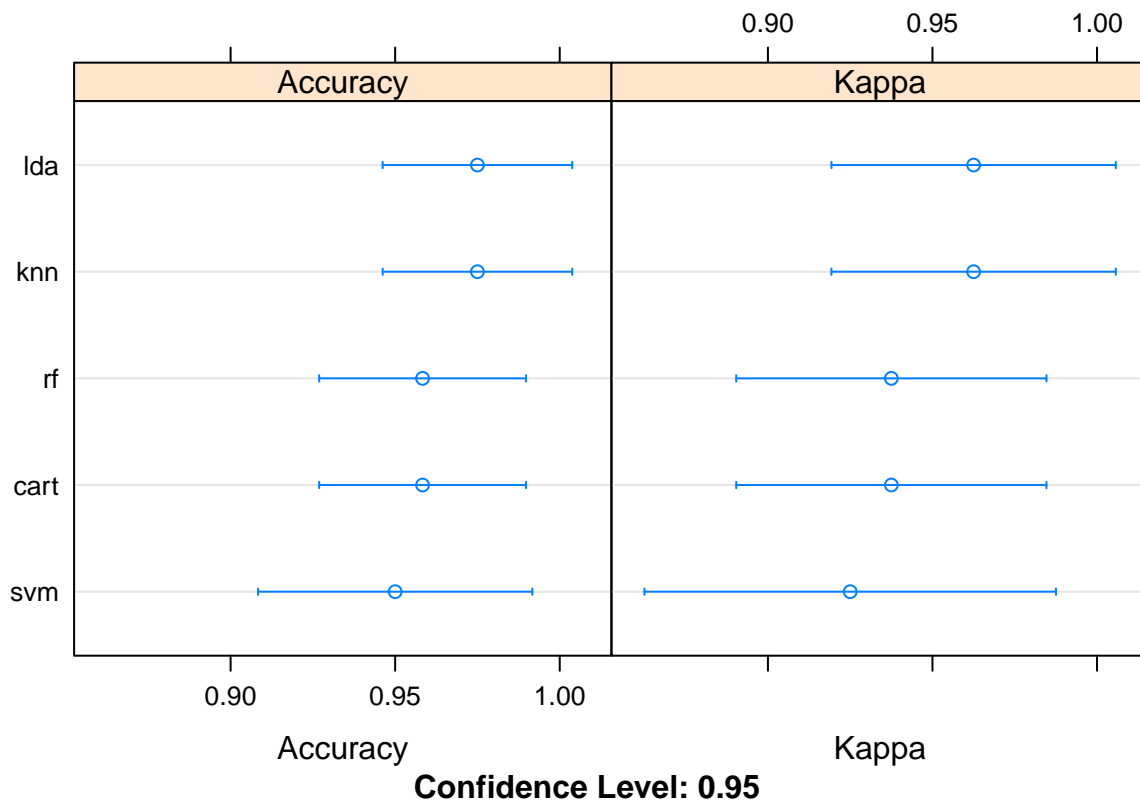
```
#install.packages("e1071")
```

```
# 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 0.9375000 1.0000000 0.9750000      1      1      0
## cart 0.9166667 0.9166667 0.9583333 0.9583333      1      1      0
## knn  0.9166667 0.9375000 1.0000000 0.9750000      1      1      0
## svm  0.8333333 0.9166667 0.9583333 0.9500000      1      1      0
## rf   0.9166667 0.9166667 0.9583333 0.9583333      1      1      0
##
## Kappa
##      Min. 1st Qu. Median     Mean 3rd Qu.  Max. NA's
## lda  0.875 0.90625 1.0000 0.9625      1      1      0
## cart 0.875 0.87500 0.9375 0.9375      1      1      0
## knn  0.875 0.90625 1.0000 0.9625      1      1      0
## svm  0.750 0.87500 0.9375 0.9250      1      1      0
## rf   0.875 0.87500 0.9375 0.9375      1      1      0
```

```
dotplot(results)
```



```
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
## 0.975 0.9625
```

Predict

```
# 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

```