

INFS4203/7203 Data Mining
The University of Queensland, Australia
Semester 2, 2018

Tutorial Week 10: Classification with R

Chandra Prasetyo Utomo
c.utomo@uq.edu.au

Objectives

1. To be able to divide dataset into training and test data.
2. To be able to implement decision tree algorithm.
3. To be able to implement K-NN algorithm.
4. To be able to evaluate classification results.

Outline

1. Experimental Data Creation (*10 minutes*)
2. Decision Tree Algorithm Implementation (*15 minutes*)
3. K-NN Algorithm Implementation (*15 minutes*)
4. Evaluation Computation (*10 minutes*)

Part 1

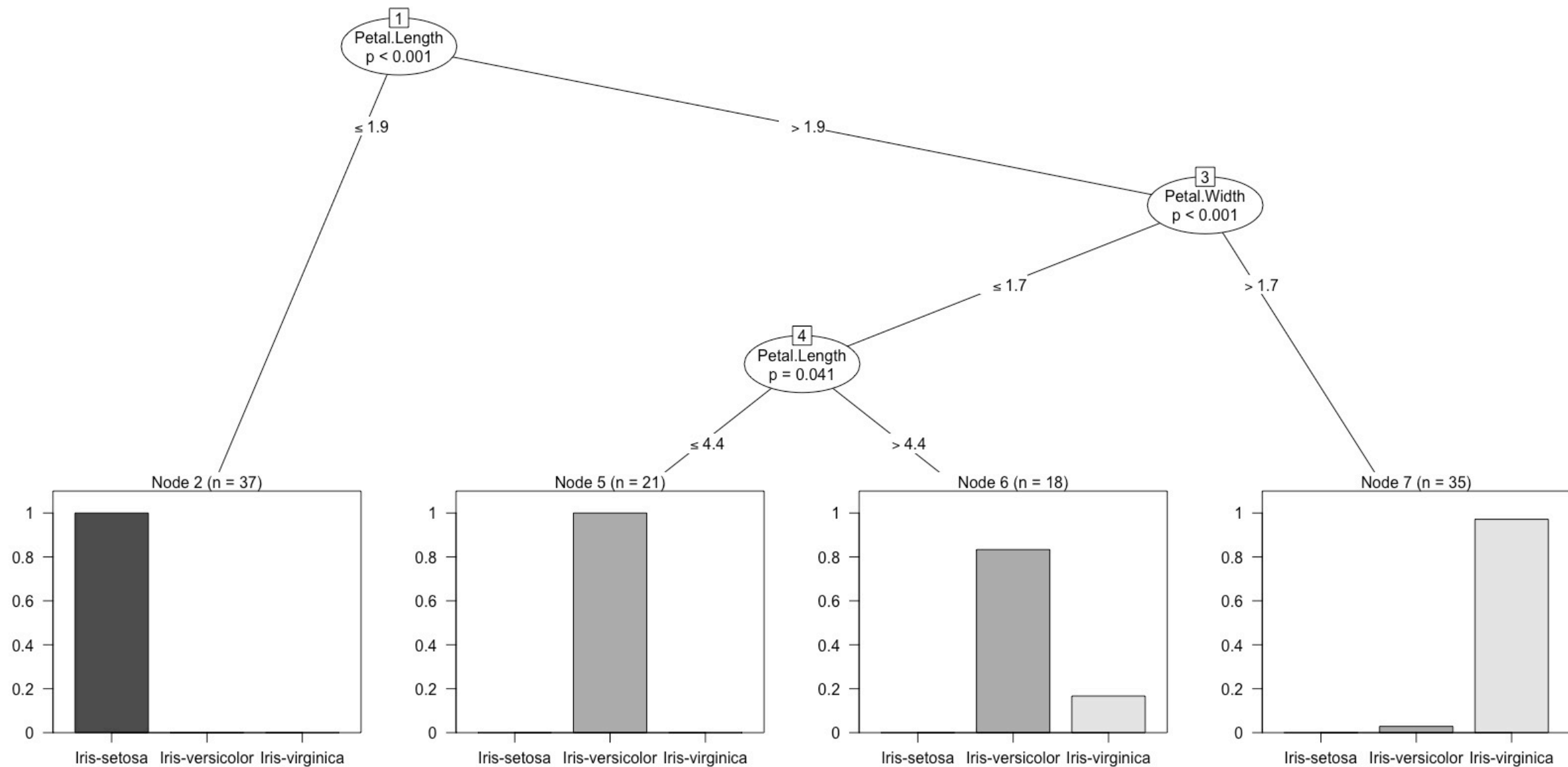
Experimental Data

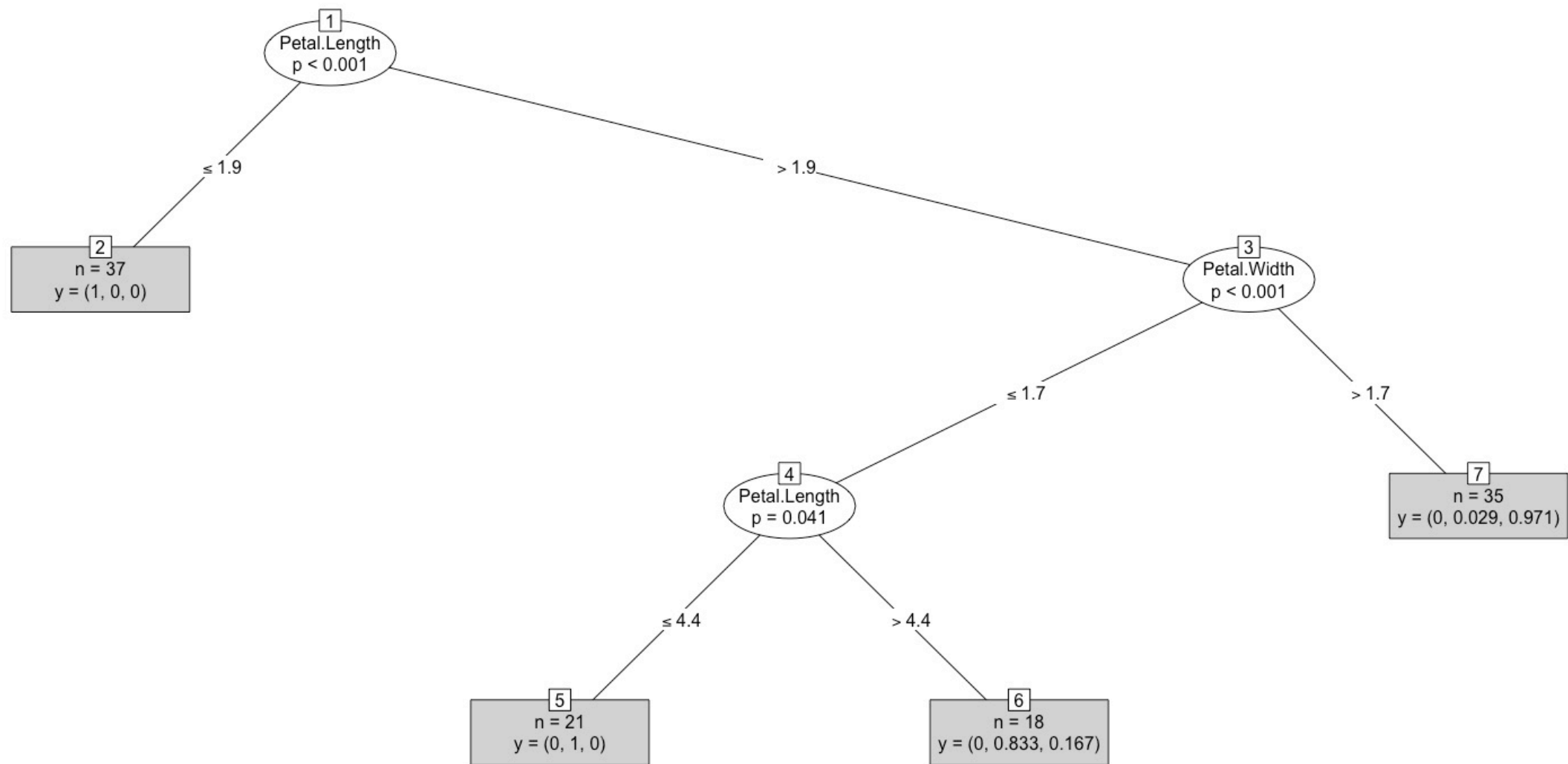
```
1 ## Practical Tutorial 4: Classification ##
2 ## Part 4A: Create Experimental Data ##
3
4 # Extract Data
5 iris <- read.table("./data/iris.data", sep = ',')
6
7 # Assign name to variables
8 names(iris) <- c("Sepal.Length", "Sepal.Width", "Petal.Length", "Petal.Width", "Species")
9
10 # For reproducible result
11 set.seed(2018)
12
13 # Set training and test ratio
14 m = nrow(iris)
15 training_percentage = 0.7
16 test_percentage = 0.3
17
18 # Sample random index
19 ind <- sample(2, m, replace = TRUE, prob = c(training_percentage, test_percentage))
20
21 # Select training and test data
22 training_data = iris[ind == 1, ]
23 test_data = iris[ind == 2, ]
24
25 # Save datasets to files
26 saveRDS(training_data, file="./data/training_data.Rda")
27 saveRDS(test_data, file="./data/test_data.Rda")
```

Part 2

Decision Tree Algorithm

```
1 ## Practical Tutorial 4: Classification ##
2 ## Part 4B: Decision Tree ##
3
4 # load training and test data
5 training_data <- readRDS(file="./data/training_data.Rda")
6 test_data <- readRDS(file="./data/test_data.Rda")
7
8 # divide features and labels
9 training_features <- training_data[,1:4]
10 training_labels <- training_data[,5]
11 test_features <- test_data[,1:4]
12 test_labels <- test_data[,5]
13
14 # install and import "party" library
15 install.packages("party")
16 library(party)
17
18 # specify target (class) and predictors (features)
19 myFormula <- Species ~ Sepal.Length + Sepal.Width +
20   Petal.Length + Petal.Width
21
22 # generate classification tree
23 iris_ctree <- ctree(myFormula, data = training_data)
24
25 # visualise the tree
26 plot(iris_ctree)
27 plot(iris_ctree, type="simple")
28
29 # predict test labels
30 ctree_pred <- predict(iris_ctree, newdata = test_features)
31 saveRDS(ctree_pred, file="./data/ctree_pred.Rda")
```





Part 3

K-NN Algorithm

```
1 ## Practical Tutorial 4: Classification ##
2 ## Part 4C: K-NN ##
3
4 # load training and test data
5 training_data <- readRDS(file="./data/training_data.Rda")
6 test_data <- readRDS(file="./data/test_data.Rda")
7
8 # divide features and labels
9 training_features <- training_data[,1:4]
10 training_labels <- training_data[,5]
11 test_features <- test_data[,1:4]
12 test_labels <- test_data[,5]
13
14 # install and import "class" library
15 install.packages("class")
16 library(class)
17
18 # classify using K-NN
19 knn_pred <- knn(train = training_features,
20                test = test_features,
21                cl = training_labels,
22                k = 3)
23
24 # save actual and predicted labels
25 saveRDS(test_labels, file="./data/test_labels.Rda")
26 saveRDS(knn_pred, file="./data/knn_pred.Rda")
```

Part 4

Evaluation Computation

```
1 |## Practical Tutorial 4: Classification ##
2 |## Part 4D: Evaluation ##
3 |
4 |# load actual and predicted labels
5 |test_labels <- readRDS(file="./data/test_labels.Rda")
6 |ctree_pred <- readRDS(file="./data/ctree_pred.Rda")
7 |knn_pred <- readRDS(file="./data/knn_pred.Rda")
8 |
9 |# create the confusion matrix
10 |cm = as.matrix(table(Actual = test_labels, Predicted = ctree_pred))
11 |
12 |n = sum(cm) # number of instances
13 |nc = nrow(cm) # number of classes
14 |diag = diag(cm) # number of correctly classified instances per class
15 |rowsums = apply(cm, 1, sum) # number of instances per class
16 |colsums = apply(cm, 2, sum) # number of predictions per class
17 |
18 |# compute accuracy, precision, recall, and f1
19 |accuracy = sum(diag) / n
20 |precision = diag / colsums
21 |recall = diag / rowsums
22 |f1 = 2 * precision * recall / (precision + recall)
23 |
24 |results <- data.frame(precision, recall, f1)
25 |accuracy
26 |results
```

> cm

Actual	Predicted		
	Iris-setosa	Iris-versicolor	Iris-virginica
Iris-setosa	13	0	0
Iris-versicolor	0	13	0
Iris-virginica	0	2	11

> accuracy

[1] 0.9487179

> results

	precision	recall	f1
Iris-setosa	1.0000000	1.0000000	1.0000000
Iris-versicolor	0.8666667	1.0000000	0.9285714
Iris-virginica	1.0000000	0.8461538	0.9166667

References:

- R and Data Mining. Yangchan Zhao. Academic Press 2012.
<https://www.sciencedirect.com/book/9780123969637/r-and-data-mining>
- Computing Classification Evaluation Metrics in R. *by Said Bleik, Shaheen Gauher, Data Scientists at Microsoft.*
http://blog.revolutionanalytics.com/2016/03/com_class_eval_metrics_r.html