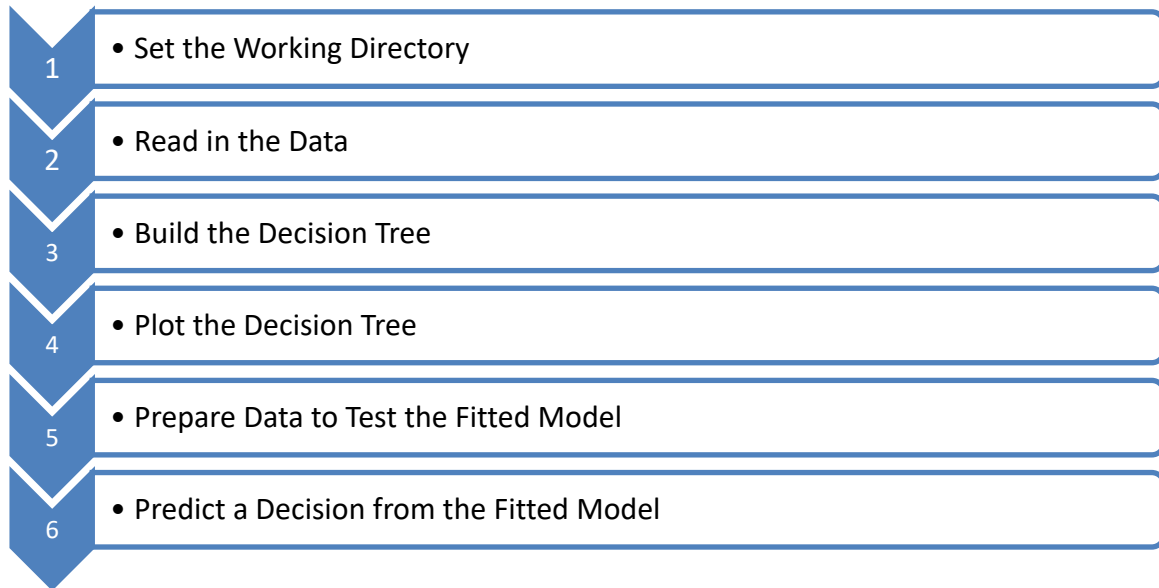


## Lab Exercise 9: Decision Trees

<b>Purpose:</b>	<p>This lab is designed to investigate and practice Decision Tree (DT) models covered in the course work. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Use R functions for Decision Tree models</li><li>• Predict the outcome of an attribute based on the model</li></ul>
<b>Tasks:</b>	<p>Tasks you will complete in this lab exercise include:</p> <ul style="list-style-type: none"><li>• Use the R –Studio environment to code Decision Tree Models</li><li>• Build a Decision Tree Model based on data whose schema is composed of attributes</li><li>• Predict the outcome of one attribute based on the model</li></ul>

## Workflow Overview



## LAB Instructions

Step	Action
1	Log in with GPADMIN credentials on to R-Studio.
2	<p><b><u>Set the Working Directory:</u></b></p> <ol style="list-style-type: none"> <li>Execute the command:</li> </ol> <pre>setwd("~/LAB09") &gt; library("rpart") &gt; library("rpart.plot")</pre>
3	<p><b><u>Read in the Data:</u></b></p> <ul style="list-style-type: none"> <li>Use a data table with columns for data attributes : Play, Outlook, Temperature, Humidity and Windy</li> <li>A Decision Tree allows for predicting the values of the attribute Play, given that you know the values for attributes like Outlook, Humidity and Windy.</li> </ul> <ol style="list-style-type: none"> <li>Read in the data from the "Dtdata.csv" file in the working directory and display the contents:</li> </ol> <pre>&gt; #Read the data &gt; play_decision &gt; play_decision&lt;- read.table("DTdata.csv",header=TRUE,sep=",") &gt; play_decision</pre> <ol style="list-style-type: none"> <li>How many observations did you read in? 10</li> <li>How many variables (attributes) did you read in? 5 – Play, Outlook, Temperature, Humidity, and Wind</li> <li>Use the command "summary" for a detailed list of the table object you read in</li> </ol> <pre>summary(play_decision)</pre> <ol style="list-style-type: none"> <li>Review the results. (The Summary is located in the console window.)</li> </ol> <pre>&gt; summary(play_decision)   Play      Outlook Temperature  Humidity   Wind no :3   overcast:2   cool:5     high :4   Mode :logical yes:7   rainy  :4   hot :2     normal:6  FALSE:7       sunny  :4   mild:3                      TRUE :3  NA's :0</pre>

Step	Action
4	<p><b><u>Build the Decision Tree:</u></b></p> <p>Use the “rpart” package in R for classification by Decision Trees. The RPart Programs build classification or regression models of a very general structure using a two stage procedure; the resulting models can be represented as binary trees.</p> <p>1. Use the following rpart commands to grow a Decision Tree:</p> <pre>rpart (formula, data=, method=, control=)</pre> <div data-bbox="349 537 1507 1350" style="border: 1px solid black; padding: 10px;"> <ul style="list-style-type: none"> <li>• <b>formula</b> is in the format: outcome ~ predictor1+predictor2+predictor3+ect.</li> <li>• <b>data=</b> specifies the dataframe</li> <li>• <b>method=</b> "class" for a classification tree "anova" for a regression tree</li> <li>• <b>control=</b> optional parameters for controlling tree growth. For example, control=rpart.control(minsplit=30, cp=0.001) requires that the minimum number of observations in a node be 30 before attempting a split and that a split must decrease the overall lack of fit by a factor of 0.001 (cost complexity factor) before being attempted.</li> <li>• <b>parms=</b> Optional parameters for the splitting function. Anova splitting has no parameters. Poisson splitting has a single parameter, the coefficient of variation of the prior distribution on the rates. The default value is 1. Exponential splitting has the same parameter as Poisson. For classification splitting, the list can contain any of: the vector of prior probabilities (component prior), the loss matrix (component loss) or the splitting index (component split). The priors must be positive and sum to 1. The loss matrix must have zeros on the diagonal and positive off-diagonal elements. The splitting index can be gini or information. The default priors are proportional to the data counts, the losses default to 1, and the split defaults to gini.</li> </ul> </div> <p>The "Play" attribute is the outcome that will be predicted.</p> <p>2. Use the command:</p> <pre>&gt; fit &lt;- rpart(Play ~ Outlook + Temperature + Humidity + Wind, method="class", data=play_decision, + control=rpart.control(minsplit=1) + parms=list(split='information')</pre> <p>3. You can now display “fit” and review the results:</p> <pre>&gt; summary(fit)</pre> <p>Note that the leaf nodes information includes both the class label and the class probabilities (P(no), P(yes))</p>

Step	Action
5	<p><b><u>Plot the Decision Tree:</u></b></p> <ol style="list-style-type: none"> <li>1. Review the arguments for <code>rpart.plot</code> function. Type in:  <code>&gt; ?rpart.plot</code></li> </ol> <p>We will use the arguments “type” and “extra” in our plot.</p> <ol style="list-style-type: none"> <li>2. Type in the following :  <code>&gt; rpart.plot(fit, type=4, extra=1)</code></li> <li>3. Review the Decision Tree plot on the graphics window.  <b>***Screenshot</b></li> </ol> <pre> graph TD     Root["yes 3 7"] --&gt; Temperat=mld  Left["no 2 1"]     Root --&gt; col,hot  Right["yes 1 6"]     Left --&gt; Outlook=snn  LeftLeft["no 2 0"]     Left --&gt; ovr  LeftRight["yes 0 1"]     Right --&gt; Wind&gt;=0.5  RightLeft["no 1 0"]     Right --&gt; &lt;0.5  RightRight["yes 0 6"] </pre>

Step	Action																		
6	<p><b><u>Prepare Data to Test the Fitted Model:</u></b></p> <p>You must use “fit” for a new set of data to create predictions from the DT:</p> <table><tr><td>Play Decision</td><td>Outlook</td><td>Temperature</td><td>Humidity</td><td>Wind</td></tr><tr><td>?</td><td>rainy</td><td>mild</td><td>high</td><td>FALSE</td></tr></table> <p>1. “newdata” is a data frame object and can be built for our test data. Type in the following statement:</p> <pre>newdata &lt;- data.frame(Outlook="rainy",Temperature="mild",Humidity="high",Wind=F ALSE)</pre> <p>2. Review the “newdata” displaying the dataframe</p> <pre>&gt; newdata</pre> <p>3. The data displayed as follows:</p> <table><tr><td>Outlook</td><td>Temperature</td><td>Humidity</td><td>Wind</td></tr><tr><td>1 rainy</td><td>mild</td><td>high</td><td>FALSE</td></tr></table>	Play Decision	Outlook	Temperature	Humidity	Wind	?	rainy	mild	high	FALSE	Outlook	Temperature	Humidity	Wind	1 rainy	mild	high	FALSE
Play Decision	Outlook	Temperature	Humidity	Wind															
?	rainy	mild	high	FALSE															
Outlook	Temperature	Humidity	Wind																
1 rainy	mild	high	FALSE																

7

**Predict a Decision from the Fitted Model:**

The “predict” function is used to generate predictions from a fitted rpart object.

- “type” is a character string denoting the type of the predicted value
- Use both “prob” and “class” to predict from a Decision Tree model

```
predict(object, newdata = list(),  
        type = c("vector", "prob", "class", "matrix"))
```

1. The **type="prob"** gives the class probabilities for the prediction for newdata Type in  
**> predict(fit,newdata=newdata,type="prob")**

2. Repeat the prediction with type="class"

```
> predict(fit,newdata=newdata,type="class")
```

Review the results.

3. What is the prediction for the “newdata”?

```
> predict(fit,newdata=newdata,type=c("class"))  
1  
no  
Levels: no yes
```

*End of Lab Exercise*

