

 Fork Settings

Find file Copy path

7c776e2 Oct 21, 2017

1 contributor

[illegible]

```

62     font-size: 12px;
63 }
64 .table th:not([align]) {
65     text-align: left;
66 }
67 </style>
68
69
70 </head>
71
72 <body>
73
74 <style type="text/css">
75 .main-container {
76     max-width: 940px;
77     margin-left: auto;
78     margin-right: auto;
79 }
80 code {
81     color: inherit;
82     background-color: rgba(0, 0, 0, 0.04);
83 }
84 img {
85     max-width: 100%;
86     height: auto;
87 }
88 .tabbed-pane {
89     padding-top: 12px;
90 }
91 button.code-folding-btn:focus {
92     outline: none;
93 }
94 </style>
95
96
97
98 <div class="container-fluid main-container">
99
100 <!-- tabsets -->
101 <script>
102 $(document).ready(function () {
103     window.buildTabsets("TOC");
104 });
105 </script>
106
107 <!-- code folding -->
108
109
110
111
112
113
114 <div class="fluid-row" id="header">
115
116
117
118 <h1 class="title toc-ignore">post01-Wei-Li</h1>
119
120 </div>
121
122
123 <pre><code>                                Trees in Statistics and R language
124     Tree is a widely used data structure in computer science. We have binary search tree, heap, B+ tree, trie and a lot of mutati
125     When I looked into the manual of R language, which is one of the best programming language for statistical computing and also
126     We will first load the package with library call.
127 </code></pre>
128 <pre class="r"><code>library(rpart)</code></pre>
129 <p>Then we can grow a classification tree with rpart method and the built-in data frame kyphosis</p>
130 <pre class="r"><code>fit <math>\leftarrow</math> rpart(Kyphosis ~ Age + Number + Start,
131     method="class", data=kyphosis)</code></pre>
132 <p>With printcp, we can display the result.</p>
133 <pre class="r"><code>printcp(fit)</code></pre>
134 <pre><code>##
135 ## Classification tree:
136 ## rpart(formula = Kyphosis ~ Age + Number + Start, data = kyphosis,
137 ##     method = "class", data=kyphosis)
138 ##
139 ## Variables actually used in tree construction:
140 ## [1] Age    Start

```

```

141 ##
142 ## Root node error: 17/81 = 0.20988
143 ##
144 ## n= 81
145 ##
146 ##      CP nsplit rel error  xerror   xstd
147 ## 1 0.176471      0  1.00000 1.00000 0.21559
148 ## 2 0.019608      1  0.82353 0.82353 0.20018
149 ## 3 0.010000      4  0.76471 0.82353 0.20018</code></pre>
150 <p>Create simple plot with plotcp</p>
151 <pre class="r"><code>pdf(file = &quot;../images/cp_plot.pdf&quot;);
152 plotcp(fit)
153 dev.off()</code></pre>
154 <pre><code>## quartz_off_screen
155 ##      2</code></pre>
156 <p>Summarize the data with summary</p>
157 <pre class="r"><code>sink(&quot;../output/kyphosis-summary&quot;);
158 summary(fit)</code></pre>
159 <pre><code>## Call:
160 ## rpart(formula = Kyphosis ~ Age + Number + Start, data = kyphosis,
161 ##       method = &quot;class&quot;);
162 ##      n= 81
163 ##
164 ##      CP nsplit rel error  xerror   xstd
165 ## 1 0.17647059      0 1.0000000 1.0000000 0.2155872
166 ## 2 0.01960784      1 0.8235294 0.8235294 0.2001751
167 ## 3 0.01000000      4 0.7647059 0.8235294 0.2001751
168 ##
169 ## Variable importance
170 ##      Start   Age Number
171 ##      64    24    12
172 ##
173 ## Node number 1: 81 observations,   complexity param=0.1764706
174 ##   predicted class=absent   expected loss=0.2098765   P(node) =1
175 ##   class counts:    64    17
176 ##   probabilities: 0.790 0.210
177 ##   left son=2 (62 obs) right son=3 (19 obs)
178 ##   Primary splits:
179 ##     Start &lt; 8.5 to the right, improve=6.762330, (0 missing)
180 ##     Number &lt; 5.5 to the left, improve=2.866795, (0 missing)
181 ##     Age &lt; 39.5 to the left, improve=2.250212, (0 missing)
182 ##   Surrogate splits:
183 ##     Number &lt; 6.5 to the left, agree=0.802, adj=0.158, (0 split)
184 ##
185 ## Node number 2: 62 observations,   complexity param=0.01960784
186 ##   predicted class=absent   expected loss=0.09677419   P(node) =0.7654321
187 ##   class counts:    56     6
188 ##   probabilities: 0.903 0.097
189 ##   left son=4 (29 obs) right son=5 (33 obs)
190 ##   Primary splits:
191 ##     Start &lt; 14.5 to the right, improve=1.0205280, (0 missing)
192 ##     Age &lt; 55 to the left, improve=0.6848635, (0 missing)
193 ##     Number &lt; 4.5 to the left, improve=0.2975332, (0 missing)
194 ##   Surrogate splits:
195 ##     Number &lt; 3.5 to the left, agree=0.645, adj=0.241, (0 split)
196 ##     Age &lt; 16 to the left, agree=0.597, adj=0.138, (0 split)
197 ##
198 ## Node number 3: 19 observations
199 ##   predicted class=present expected loss=0.4210526   P(node) =0.2345679
200 ##   class counts:     8    11
201 ##   probabilities: 0.421 0.579
202 ##
203 ## Node number 4: 29 observations
204 ##   predicted class=absent   expected loss=0   P(node) =0.3580247
205 ##   class counts:    29     0
206 ##   probabilities: 1.000 0.000
207 ##
208 ## Node number 5: 33 observations,   complexity param=0.01960784
209 ##   predicted class=absent   expected loss=0.1818182   P(node) =0.4074074
210 ##   class counts:    27     6
211 ##   probabilities: 0.818 0.182
212 ##   left son=10 (12 obs) right son=11 (21 obs)
213 ##   Primary splits:
214 ##     Age &lt; 55 to the left, improve=1.2467530, (0 missing)
215 ##     Start &lt; 12.5 to the right, improve=0.2887701, (0 missing)
216 ##     Number &lt; 3.5 to the right, improve=0.1753247, (0 missing)
217 ##   Surrogate splits:
218 ##     Start &lt; 9.5 to the left, agree=0.758, adj=0.333, (0 split)
219 ##     Number &lt; 5.5 to the right, agree=0.607, adj=0.167, (0 split)

```

```

217 ##      Number <math>0.5</math> to the right, agree=0.097, adj=0.107, (0 split)
218 ##
219 ## Node number 10: 12 observations
220 ##      predicted class=absent      expected loss=0      P(node) =0.1481481
221 ##      class counts:      12      0
222 ##      probabilities: 1.000 0.000
223 ##
224 ## Node number 11: 21 observations,      complexity param=0.01960784
225 ##      predicted class=absent      expected loss=0.2857143      P(node) =0.2592593
226 ##      class counts:      15      6
227 ##      probabilities: 0.714 0.286
228 ##      left son=22 (14 obs) right son=23 (7 obs)
229 ##      Primary splits:
230 ##      Age      <math>111</math> to the right, improve=1.71428600, (0 missing)
231 ##      Start    <math>12.5</math> to the right, improve=0.79365080, (0 missing)
232 ##      Number   <math>3.5</math> to the right, improve=0.07142857, (0 missing)
233 ##
234 ## Node number 22: 14 observations
235 ##      predicted class=absent      expected loss=0.1428571      P(node) =0.1728395
236 ##      class counts:      12      2
237 ##      probabilities: 0.857 0.143
238 ##
239 ## Node number 23: 7 observations
240 ##      predicted class=present      expected loss=0.4285714      P(node) =0.08641975
241 ##      class counts:      3      4
242 ##      probabilities: 0.429 0.571</code></pre>
243 <pre class="r"><code>sink()</code></pre>
244 <p>Create tree-structured plot that better visualises hierarchical data</p>
245 <pre class="r"><code>pdf(file = &quot;../images/tree_plot.dpf&quot;);
246 plot(fit, uniform=TRUE,
247      main=&quot;Classification Tree for Kyphosis&quot;);
248 text(fit, use.n=TRUE, all=TRUE, cex=.8)
249 dev.off()</code></pre>
250 <pre><code>## quartz_off_screen
251 ##      2</code></pre>
252 <p>To avoid overfitting the data, we can prune the tree to minimize the cross-validated error.</p>
253 <pre class="r"><code>pfit<math>-</math> prune(fit, cp= fit$cp.table[which.min(fit$cp.table[,&quot;xerror&quot;]),&quot;CP&quot;])</code></pre>
254 <pre class="r"><code>printcp(pfit)</code></pre>
255 <pre><code>##
256 ## Classification tree:
257 ## rpart(formula = Kyphosis ~ Age + Number + Start, data = kyphosis,
258 ##      method = &quot;class&quot;);
259 ##
260 ## Variables actually used in tree construction:
261 ## [1] Start
262 ##
263 ## Root node error: 17/81 = 0.20988
264 ##
265 ## n= 81
266 ##
267 ##      CP nsplit rel error  xerror   xstd
268 ## 1 0.176471      0  1.00000 1.00000 0.21559
269 ## 2 0.019608      1  0.82353 0.82353 0.20018</code></pre>
270 <pre class="r"><code>plotcp(pfit)</code></pre>
271
272 <p>Reference: <a href="https://www.statmethods.net/advstats/cart.html" class="uri">https://www.statmethods.net/advstats/cart.ht
273
274
275
276
277
278
279 </div>
280
281 <script>
282
283 // add bootstrap table styles to pandoc tables
284 function bootstrapStylePandocTables() {
285   $('tr.header').parent('thead').parent('table').addClass('table table-condensed');
286 }
287 $(document).ready(function () {
288   bootstrapStylePandocTables();
289 });
290
291
292 </script>
293
294 <!-- dynamically load mathjax for compatibility with self-contained -->
295 <script>
296   (function () {
297     var script = document.createElement("script");

```

```
298     script.type = "text/javascript";
299     script.src   = "https://mathjax.rstudio.com/latest/MathJax.js?config=TeX-AMS-MML_HTMLorMML";
300     document.getElementsByTagName("head")[0].appendChild(script);
301   }());
302 </script>
303
304 </body>
305 </html>
```

