

```

1 local help= []
2 NB:
3 (c)2022 Tim Menzies, timm@ieee.org
4
5 OPTIONS:
6 -b --Bins      max number of bins      = 16
7 -k --k         handle rare classes     = 1
8 -m --m         handle rare attributes  = 2
9 -p --p         distance coefficient     = 2
10 -S --small     small leaf size         = .5
11 -w --wait      wait before classifying = 5
12
13 OPTIONS (other):
14 -f --file      file                    = ../.data/auto93.csv
15 -g --go        start-up goal           = nothing
16 -h --help      show help               = false
17 -s --seed      seed                    = 10019]
18
19 -- 1 2 3 4 5
20
21 local _ = require"lib"
22 local argmax,big      = _:argmax, _:big
23 local cli, csv, demos, is, normpdf = _:cli, _:csv, _:demos, _:is, _:normpdf
24 local oo, push, read, rnd, same, str = _:oo, _:push, _:read, _:rnd, _:same, _:str
25
26 local THE={}
27 help:gsub("[^-][^%s+][^n]%%s{*(%s)+}", function(key,x) THE[key] = read(x) end)
28
29 local NB, NUM, SYM, COLS, ROW, ROWS= is"NB", is"NUM", is"SYM", is"COLS", is"ROW", is"ROWS"
30 local FEW, RANGE, TREE = is"FEW", is"RANGE", is"TREE"
31
32 -- 1 2 3 4 5
33
34 function RANGE.new(i, xlo, xhi, ys) i.xlo, i.xhi, i.ys = xlo, xhi, ys end
35 function RANGE.add(i,x,y)
36 if x < i.xlo then i.xlo = x end -- works for string or num
37 if x > i.xhi then i.xhi = x end -- works for string or num
38 i.ys:add(y) end
39
40 function RANGE._:tostring(i)
41 local x, lo, hi = i.ys.txt, i.xlo, i.xhi
42 if lo == hi then return fmt("%s==%s",x, lo)
43 elseif hi == big then return fmt("%s>%s",x, lo)
44 elseif lo == -big then return fmt("%s<=%s",x, hi)
45 else return fmt("%s<%s<=%s",lo,x,hi) end end
46
47 -- 2 3 4 5 6 7 8
48
49 function FEW.new(i) i.n,i.t,i.ok=0, {}, true end
50 function FEW.has(i) i.t=i.ok and i.t or sort(i.t); i.ok=true; return i.t end
51 function FEW.add(i,x)
52 if x=="M" then return x end
53 i.n=i.n+1
54 if #i.t < THE.some then i.ok=false; push(i.t,x)
55 elseif rand() < THE.some/i.n then i.ok=false; i.t[rand(#i.t)]x end end
56
57 function NUM.new(i) i.n,i.mu,i.m2,i.mu,i.lo,i.hi,i.few=0,0,0,0,big,-big,FEW(i) end
58 function NUM.mid(i,p) return rnd(i.mu,p) end
59 function NUM.like(i,x,...) return normpdf(x, i.mu, i.sd) end
60 function NUM.bin(x)
61 b=(i.hi - i.lo)/THE.bins; return i.lo==i.hi and 1 or math.floor(x/b+.5)*b end
62
63 function NUM.add(i, NUM, v, number)
64 if v=="?" then return v end
65 i.few:add(v)
66 i.n = i.n + 1
67 local d = v - i.mu
68 i.mu = i.mu + d/i.n
69 i.m2 = i.m2 + d*(v - i.mu)
70 i.sd = i.n<2 and 0 or (i.m2/(i.n-1))^0.5
71 i.lo = math.min(v, i.lo)
72 i.hi = math.max(v, i.hi) end
73
74 function NUM.merge(i,j, k)
75 local k = NUM(i.at, i.txt)
76 for _,n in pairs(i.few.t) do k:add(x) end
77 for _,n in pairs(j.few.t) do k:add(x) end
78 return k end
79
80 function NUM.mergeRanges(i,b4,min)
81 local t,j, a,b,c,A,B,C = {},1,
82 while j < #b4 do
83 a, b = b4[j], b4[j+1]
84 if b then
85 A,B = a.ys, b.ys
86 C = A:merge(B)
87 if A.n<min or B.n<min or C:div() <= (A.n*A:div() + B.n*B:div())/C.n then
88 j = j + 1
89 a = RANGE(a.xlo, b.xhi, C) end end
90 t[#t+1] = a
91 j = j + 1 end
92 if #t < b4 then return i:mergeRanges(t,min) end
93 for j=2,#t do t[j].xlo = t[j-1].xhi end
94 t[1].xlo, t[#t].xhi = -big, big
95 return t end
96
97 function SYM.new(i) i.n,i.syms,i.most,i.mode = 0, {}, 0, nil end
98 function SYM.mid(i,...) return i.mode end
99 function SYM.like(i,x,prior) return ((i.syms[x] or 0)+THE.m*prior)/(i.n+THE.m) end
100 function SYM.bin(x) return x end
101 function SYM.add(i,v,inc)
102 if v=="?" then return v end
103 inc=inc or 1
104 i.n = i.n + inc
105 i.syms[v] = inc + (i.syms[v] or 0)
106 if i.syms[v] > i.most then i.most,i.mode = i.syms[v],v end end
107
108 function SYM.merge(i,j, k)
109 local k = SYM(i.at, i.txt)
110 for x,n in pairs(i.has) do k:add(x,n) end
111 for x,n in pairs(j.has) do k:add(x,n) end
112 return k end
113
114 function SYM.mergeRanges(i,t,...) return t end
115
116 -- 2 3 4 5
117
118 local function usep(x) return not x:find"$" end
119 local function nump(x) return x:find"^[A-Z]" end
120 local function goalp(x) return x:find"[!-]$" end
121 local function klassp(x) return x:find"$" end
122 local function new(at,txt, i)
123 txt = txt or ""
124 i = (nump(txt) and NUM or SYM)()
125 i.txt, i.usep, i.at, i.w = txt, usep(txt), at or 0, txt:find"$" and -1 or 1
126 return i end
127
128 function COLS.new(i,t, col)
129 i.all, i.xs, i.ys, i.names = {}, {}, {}, {}
130 for at,x in pairs(t) do
131 col = push(i.all, new(at,x))
132 if col.usep then
133 if klassp(col.txt) then i.klass=col end
134 push(goalp(col.txt) and i.ys or i.xs, col) end end end
135
136 function COLS.add(i,t)
137 for _,cols in pairs(i.xs,i.ys) do
138 for _,col in pairs(cols) do col:add(t[col.at]) end end
139 return t end
140
141 -- 1 2 3 4 5
142
143 function ROW.new(i,of,cells) i.of,i.cells,i.eval=of,cells,false end
144 function ROW.klass(i) return i.cells[i.of.cols.klass.at] end
145 function ROW.within(i,range, lo,hi,at,v)
146 lo, hi, at = range.xlo, range.xhi, range.ys.at
147 v = i.i[llist[at]]
148 return v=="?" or (lo==hi and v==lo) or (lo<v and v<=hi) end
149
150 function ROW.b4(i,j,at, x,y)
151 x, y = i.cells[at], j.cells[at]
152 x = x=="?" and -big or x
153 y = y=="?" and -big or y
154 return x < y end
155
156 -- 1 2 3 4 5
157
158 local function data(src, fun)
159 if type(src)=="string" then for _,t in pairs(src) do fun(t) end
160 else for t in csv(src) do fun(t) end end end
161
162 function ROWS.new(i,t) i.cols=COLS(t); i.rows={} end
163 function ROWS.add(i,t)
164 t=t.cells and t or ROW(i,t)
165 i.cols:add(t.cells)
166 return push(i.rows, t) end
167
168 function ROWS.mid(i, cols, p, t)
169 t={};for _,col in pairs(cols or i.cols.ys) do t[col.txt]=col:mid(p) end;return t end
170
171 function ROWS.clone(i,t, j)
172 j= ROWS(i.cols.names);for _,row in pairs(t or {}) do j:add(row) end; return j end
173
174 function ROWS.like(i,t, nklasses, nrows, prior,like,inc,x)
175 prior = (#i.rows + THE.k) / (nrows + THE.k * nklasses)
176 like = math.log(prior)
177 for _,col in pairs(i.cols.xs) do
178 x = t.cells[col.at]
179 if x and x ~= "?" then
180 inc = col:like(x,prior)
181 like = like + math.log(inc) end end
182 return like end
183
184 -- 1 2 3
185
186 -- (0) Use rowl to initial our 'overall' knowledge of all rows.
187 -- After that (1) add row to 'overall' and (2) ROWS about this row's klass.
188 -- (3) After 'wait' row, classify row BEFORE updating training knowledge
189 function NB.new(i,src,report, row)
190 report = report or print
191 i.overall,i.dict, i.list = nil, {}, {}
192 data(src, function(row, k)
193 if not i.overall then i.overall = ROWS(row) else -- (0) eat rowl
194 row = i.overall:add(row) -- add to overall
195 if #i.overall.rows > THE.wait then report(row:klass(i), i:guess(row)) end
196 i:train(row) end end) add tp row's klass
197
198 function NB:train(i,row) i:_known(row:klass()):add(row) end
199 function NB:_known(i,k)
200 i:dict[k] = i:dict[k] or push(i.list, i:overall:clone()) -- klass is known
201 i:dict[k].txt = k -- each klass knows its name
202 return i:dict[k] end
203
204 function NB:guess(i,row)
205 return argmax(i:dict,
206 function(klass) return klass:like(row,#i.list,#i.overall.rows) end) end
207
208 function TREE.new(i,rows,gaurd)
209 -- i.gaurd, i.kids = gaurd, {}
210 of = listOfRows[i][1].of
211 best = sort(map(of:cols.x,
212 function(col) i:bins(col,listOfRows) end),lt"div")[1]
213
214 i.kids = map(best.ranges, function(range)
215 function(listOfRowsl) = {}
216 -- local function within(row) return row:within(best) end
217 -- local function withins(rows) return map(rows, within) end
218 -- map(listOfRows, function(rows) return withins(rows) end) end
219 ttmp= map(rows,withins)
220 -- if #tmp > stop then
221 -- end)
222
223 function Tree.new(i,rows,gaurd)
224 i.gaurd, i.kids, labels = gaurd, {}, {}
225 xcols,rows = nil, {}
226 for label,rows0 in pairs(rowsa) do
227 for _,row in pairs(rows0) do
228 labels[row.id] = label
229 xcols = push(rows,row).of.cols.xs end end
230 for _,xcol in pairs(of.cols.xs) do
231 i:bins(rows, xcol, SYM, function(row) return labels[row.id] end) end end end
232
233 function TREE.bins(i,rows,xcol,yklass,y)
234 local n,list, dict = 0, {}, {}
235 for _,row in pairs(rows) do
236 local v = row.cells[xcol.at]
237 if v ~= "?" then
238 n = n + 1
239 local pos = xcol:bin(v)
240 dict[pos] = dict[pos] or push(list, RANGE(v,v, yklass(xcol.at, xcol.txt)))
241
242 dict[pos]:add(v, y(row)) end end
243 list = xcol:mergeRanges(sort(list, lt"xlo"),n*THE.min)
244 return div = sum(list,function(z) return z.ys:div(i)*z.ys.n/n end) end

```

```

244 --      -|_  (7_ > -|_ > -----
245 --
246
247 local no,go = {},{}
248 function go.the() print(1); print(1); return type(1)=="number" and 1==2 end
249
250 function go.argmax( t,fun)
251   fun=function(x) return -x end
252   t={50,40,0,40,50}
253   return 3 == argmax(t,fun) end
254
255 function go.num(n) n=NUM(); for x=1,100 do n:add(x) end; return n.mu==50.5 end
256
257 function go.sym(s)
258   s=SYM(); for _,x in pairs{"a","a","a","a","b","b","c"} do s:add(x) end
259   return s.mode=="a" end
260
261 function go.csv( n,s)
262   n,s=0,0; for row in csv(1) do n=n+1; if n>1 then s=s+row[1] end end
263   return rnd(s/n,3) == 5.441 end
264
265 function go.rows( rows)
266   data(1,1,function(t) if rows then rows:add(t) else rows=ROWS(t) end end)
267   return rnd(rows.cols.ys[1].sd,0)==847 end
268
269 function go.nb()
270   return 268 == #NB("./data/diabetes.csv").dict["positive"].rows end
271
272 local function _classify(file)
273   local Abcd=require"abcd"
274   local abcd=Abcd()
275   NB(file, function(got,want) abcd:add(got,want) end)
276   abcd:pretty(abcd:report())
277   return true end
278
279 function go.soybean() return _classify("./data/soybean.csv") end
280 function go.diabetes() return _classify("./data/diabetes.csv") end
281
282 --      > -|_  C|  | -|_ -----
283 --
284 if pcall(debug.getlocal, 4, 1)
285 then return (ROW=ROW, ROWS=ROWS, NUM=NUM, SYM=SYM, THE=THE,lib=lib)
286 else THE = cli(THE,help)
287   demos(THE,go) end
288
289 ---
290 ---
291 ---
292 ---
293 ---
294 ---
295 ---
296 ---
297 ---

```