

```

1  -- If you understand "it", can you write "it" shorter? Lets try.
2  -- E.G. how short can we write a multi-objective semi-supervised learner?<br>
3  -- <span id="forkongithub"><a href="https://github.com/timm/15">Fork me on GitHub</a></span>
4
5  -- _Mother Teresa_: The more you have, the more you are occupied.
6  -- The less you have, the more free you are.<p>
7  -- _Ken Thompson_: One of my most productive days was throwing
8  -- away 1,000 lines of code.<p>
9  -- _William of Occam_: It is vain to do with more what can be done with less. <p>
10 -- _King width=200 align=left src=cup.png>
11 -- _Leonardo da Vinci_: Simplicity is the ultimate sophistication.<p>
12 -- _Edeger D. Dijkstra_: Simplicity is prerequisite for reliability.<p>
13 -- _Dieter Rams_: Less, but better.<p>
14 -- _timtm_: less, plz <p>
15 -- My heroes: [Jack Diederich] (https://www.youtube.com/watch?v=o9pEgZHorH0)
16 -- [Hilary Mason] (https://boingboing.net/2017/06/30/next-level-regex.html)<p>
17 local help:lua is a multi-objective semi-supervised learner.
18 (c)2022 Tim Menzies, timm@ieee.org
19
20 OPTIONS:
21 -b --Bins max number of bins = 16
22 -F --Few only keep a "Few" numbers = 256
23 -k --k handle rare classes = 1
24 -n --n handle rare attributes = 2
25 -p --p distance coefficient = 2
26 -S --small small leaf size = .5
27 -w --wait wait before classifying =
28
29 OPTIONS (other):
30 -f --file file = .././data/auto93.csv
31 -g --goal start-up goal =
32 -h --help show help = false
33 -s --seed seed = 10019]
34
35 -----
36 -- ## Names
37 local _ = require"lib"
38 local argmax,big = _._argmax, _._big
39 local cli, csv, demos, klass, normpdf = _._cli, _._csv, _._demos, _._klass, _._normpdf
40 local o, push, read, rnd, same, str = _._o, _._push, _._read, _._rnd, _._same, _._str
41
42 -- 'THE' settings is parsed from 'help'
43 -- string lines that contain two dashesnbsps;""---".
44 local THE={}
45 help:gsuib("[-][-]([%s+])[%s+]([%s+])", function(key,x) THE[key] = read(x) end)
46 -- ## Our classes
47
48 -- ROWS use COLS to make either NUMs or SYMs.
49 -- ROWS holds data in ROWs, and summarizes columns in NUMs and SYMs.
50 -- NUMs use SOMEs to store at most 'THE.Few' samples per numeric columns.
51 -- RANGE objects track what 'y' values are seen between 'xlo' and 'xhi'.
52 local ROWS, COLS, NUM, SYM = klass"ROWS", klass"COLS", klass"NUM", klass"SYM"
53 local ROW = klass"ROW"
54 local SOME = klass"SOME"
55 local RANGE = klass"RANGE"
56
57 -- ## class COLS
58
59 -- The system uses COLS to make lists of NUMs or SYMs.
60 -- Independent and dependent columns are stored in 'xs' or 'ys' lists.<p>
61 -- This code reads data from csv files or lists of tables.
62 -- The first row of data is a list of column names.
63 -- Those can contain some special symbols.
64 local is={}
65 -- For example, we want to ignore columns whose name ends with "':".
66 -- Number columns start with an upper case letter
67 -- Dependent variables end in "':", "':", or "':"
68 -- For classes or goals to be minimized or maximized.
69 function is.use(x) return not x:find"'" end
70 function is.num(x) return x:find"^[A-Z]" end
71 function is.goal(x) return x:find"^[+-$]" end
72 function is.klass(x) return x:find"^[S]" end
73 function is.dislike(x) return x:find"^[S]" end
74
75 -- For example, 'COLS' would read line1 of this data, skipping column #2, making
76 -- NUMerics out of columns #1, #2, and a SYMbollic out of column1 #4. Further,
77 -- any goal cols in this case, "the" "y" get stored in a 'ys' list
78 -- and all other columns are stored in 'xs'.
79
80 -- Clndrs, Hp:, Lbs:, origin
81 -- 8, 193, 4732, 1
82 -- 8, 215, 4615, 1
83 -- 8, 200, 4376, 1
84 -- ...
85
86 -- _COLS('t':[string])_<br>constructor.
87 function COLS.new(i,t, new,is)
88 i,xs,i.ys,i.names = {},{},{},{}
89 for at,txt in pairs(t) do
90 new = (is.nump(txt) and NUM or SYM) (at,txt)
91 new.usep, new.w = is.use(txt), is.dislike(txt) and -1 or 1
92 if new.usep then
93 if is.klass(new.txt) then i.klass=new end
94 push(is.goal(new.txt) and i.ys or i.xs, new) end end end
95
96 -- _add('t':ROW)_ :S_<br>update column summaries.
97 function COLS.add(i,t)
98 for _,cols in pairs(i.xs,i.ys) do
99 for _,col in pairs(cols) do col:add(t.cells[col.at]) end end
100 return t end
101
102 -- ## class SOME
103 -- NUMs use SOMEs to store at most 'THE.Few' samples per numeric columns.
104
105 -- _SOME()_<br>constructor.
106 function SOME.new(i) i.n,i.t,i.ok=0,{},true end
107 -- _add('x:any')_<br>
108 -- Technically, this is a reservoir sampler; i.e. given a small cache,
109 -- sample at an equal (but low) probability across a much larger space.
110 -- Note one trick: if we full, make space by replacing anything at all.
111 function SOME.add(i,x)
112 if x=="?" then return x end
113 i.n=i.n+1
114 if #i.t < THE.some then i.ok=false; push(i.t,x)
115 elseif rand() < THE.some/i.n then i.ok=false; i.t[rand(#i.t)]=x end end
116 -- _all()_ :table<br>If we ask for 'all()', then make sure they are sorted.
117 function SOME.all(i) i.t=i.ok and i.t or sort(i.t); i.ok=true; return i.t end
118
119 -- ## class NUM
120 -- Summarize a stream of numbers, maintaining its 'lo', 'hi', 'mu' (mean).
121
122 -- _NUM()_<br>constructor.
123 function NUM.new(i) i.n,i.mu,i.m2,i.w,i.lo,i.hi,i.some=0,0,0,1,big,-big,SOME(i) end
124 -- _add('v':number)_<br>add 'v' to the summary, updating 'mu,sd,lo,hi'.
125
126 function NUM.add(i, v)
127 if v=="?" then return v end
128 i.some=add(v)
129 i.n = i.n + 1
130 local d = v - i.mu
131 i.mu = i.mu + d/i.n
132 i.m2 = i.m2 + d*(v - i.mu)
133 i.sd = i.n<2 and 0 or (i.m2/(i.n-1))^0.5
134 i.lo = math.min(v, i.lo)
135 i.hi = math.max(v, i.hi) end
136 -- _bin('x':number)_ :number<br>return x, rounded into 'THE.bins' values.
137 function NUM.bin(x)
138 b=(i.hi - i.lo)/THE.bins; return i.lo==i.hi and 1 or math.floor(x/b+.5)*b end
139 -- _like('x':number)_ :number<br>return likelihood.
140 function NUM.like(i,x,...) return normpdf(x, i.mu, i.sd) end
141 -- _bin('x':number)_ :NUM<br>Return a NUM that combines self and 'j'.
142 -- _merge('j':NUM)_ :NUM<br>Return a NUM that combines self and 'j'.
143 function NUM.merge(i,j,k)
144 local k = NUM(i.at, i.txt)
145 for _,n in pairs(i.some.t) do k:add(x) end
146 for _,n in pairs(j.some.t) do k:add(x) end
147 return k end
148 -- _mid('p':p=2)_ :number<br>report rounded middle.
149 function NUM.mid(i,p) return rnd(i.mu,p) end
150
151 -- ## class SYM
152 -- Summarize a stream of symbols, maintaining the 'mode' and frequencies counts
153 -- of symbols seen so far.
154
155 -- _SYM()_<br>constructor.
156 function SYM.new(i) i.n,i.syms,i.most,i.mode = 0,{},0,nil end
157 -- _add('v:any')_ :inc=1<br>add 'v', 'inc' number of times.
158 function SYM.add(i,v,inc)
159 if v=="?" then return v end
160 inc=inc or 1
161 i.n = i.n + inc
162 i.syms[v] = inc + (i.syms[v] or 0)
163 if i.syms[v] > i.most then i.most,i.mode = i.syms[v],v end end
164 -- _bin('x':any)_ :any<br>discrediting a symbol just returns that symbol.
165 function SYM.bin(x) return x end
166 -- _like('x':any)_ :number<br>report likelihood.
167 function SYM.like(i,x,prior) return ((i.syms[x] or 0)+THE.m*prior)/(i.n+THE.m) end
168 -- _merge('j':SYM)_ :SYM<br>Return a SYM that combines self and 'j'.
169 function SYM.merge(i,j,k)
170 local k = SYM(i.at, i.txt)
171 for x,n in pairs(i.has) do k:add(x,n) end
172 for x,n in pairs(j.has) do k:add(x,n) end
173 return k end
174 -- _mid()_ :any<br>return the most commonly seen symbol.
175 function SYM.mid(i,...) return i.mode end
176
177 -- ## class ROW
178 -- Store data from one record. Track if we ever use this ROW's dependent variables.
179 -- Hold a pointer back to a ROWs (so we can access attribute positions and
180 -- weightings).
181
182 -- _ROW()_<br>constructor.
183 function ROW.new(i,of,t) i.of,i.cells,i.evald = of,t,false end
184 -- _klass_<br>return the klass variable of this ROW.
185 function ROW.klass(i) return i.cells[i.of.cols.klass.at] end
186 -- _within('range':RANGE)_ :boolean<br>True if ROW is in 'range'.
187 function ROW.within(i,range)
188 local lo, hi, at = range.xlo, range.xhi, range.ys.at
189 local v = i.cells[at]
190 return v=="?" or (lo==hi and v==lo) or (lo<v and v<=hi) end
191
192 -- ## class ROWS
193
194 -- Store some ROWs. Keep a summary of the columns inside 'cols'. Keep the
195 -- dependent and independent summaries seepage in 'cols.ys' and 'cols.xs'.
196
197 -- _ROWS('t':[string])_<br>constructor.
198 function ROWS.new(i,t) i.cols=COLS(t); i.rows={} end
199 -- _add('t':ROW)_ :ROW<br>update with a table or ROW.
200 function ROWS.add(i,t) return push(i.rows, i.cols:add(t.cells and t or ROW(i,t))) end
201 -- _mid('cols':(NUM|SYM), 'p'=2)_ :table<br>
202 -- returns 'mid's of some columns; round numerics to 'p' decimal places.
203 function ROWS.mid(i, cols, p)
204 t={};for _,col in pairs(cols or i.cols.ys) do t[col.txt]=col:mid(p) end;return t end
205 -- _clone('data':(table|ROW))_ :ROWS<br>copy this structure, maybe add data
206
207 function ROWS.clone(i,data, j)
208 j= ROWS(i.cols.names);for _,row in pairs(data or {}) do j:add(row) end; return j end
209 -- _like('row':ROWS, 'nklasses':int, 'nrows':int)_ :number<br>
210 -- How likely is it that 'row' could live here?
211 function ROWS.like(i,row, nklasses, nrows, prior,like,inc,x)
212 prior = (#i.rows + THE.k) / (nrows + THE.k * nklasses)
213 like = math.log(prior)
214 for _,col in pairs(i.cols.xs) do
215 x = row.cells[col.at]
216 if x and x ~= "?" then
217 inc = col:like(x,prior)
218 like = like + math.log(inc) end end
219 return like end
220
221 -- _doRows('src':(string|table), 'fun':function(table|ROW))_<br>
222 -- helper function for reading from tables or files. Case argl of ...
223 -- ... _table_ : call function for all items in table.
224 -- ... _string_ : call function on all rows from a file.
225 -- ... _nil_ : call function of all rows from standard input.
226 local function doRows(src, fun)
227 if type(src)=="table" then for _,t in pairs(src) do fun(t) end
228 else for t in csv(src) do fun(t) end end end
229
230 -- ## class NB
231 -- (0) Use row1 to initial our 'overall' knowledge of all rows.
232 -- After that (1) add row t 'overall' and (2) ROWs about this row's klass.
233 -- (3) After 'wait' rows, classify row BEFORE updating training knowledge
234 function NB.new(i,src,report, row)
235 report = report or print
236 i.overall, i.dict, i.list = nil, {}, {}
237 doRows(src, function(row, k)
238 if not i.overall then i.overall = ROWS(row) else -- (0) eat row1
239 row = i.overall:like(row)
240 if #i.overall.rows > THE.wait then report(row:klass() , i:guess(row)) end
241 i:train(row) end end) end -- add tp rows's klass
242
243 function NB.train(i,row, k)
244 k=row:klass()
245 i.dict[k] = i.dict[k] or push(i.list,i.overall:clone()) -- klass is known
246 i.dict[k].txt = k -- each klass knows its name
247 i.dict[k]:add(row) end -- update klass with row
248
249 function NB.guess(i,row)
250
251 return argmax(i.dict,
252 function(klass) return klass:like(row,#i.list,#i.overall.rows) end) end
253
254 -- function TREE.new(i,listOfRows,gaud)
255 -- i.gaud,i.kids = gaud, {}
256 -- of = listOfRows[1][1].of
257 -- best = sort(map(of.cols.x,
258 -- function(col) i:bins(col,listOfRows) end),lt"div")[1]
259 -- i.kids = map(best.ranges, function(range)
260 -- listOfRows= {}
261 -- local function within(row) return row:within(best) end
262 -- local function within(rows) return map(rows, within) end
263 -- map(listOfRows, function(rows) return within(rows) end)
264 -- tmp = map(rows,within)
265 -- if #tmp > stop then
266 -- end)
267
268 -----
269 -- ## class TREE
270 -- function decisionTree(listOfRows)
271 -- function tree(rows, xcols, yklass,y, gaud)
272 -- local function ranges(rows,col) return i:ranges(rows,col,yklass,y) end
273 -- i.gaud = gaud
274 -- ranges = sort(map(xcols, xranges),lt"div")[1].ranges
275 -- for _,row in pairs(rows) do
276 -- for _,range in pairs(ranges) do
277 -- if row:within(range) then push(range.rows,row) end; break end end
278 -- i.kids = map(ranges,
279 -- function(range) return TREE(range.rows,xcols,yklass,y,range) end)
280 end
281
282 -- labels , all, xcols = {},{}
283 -- for label,rows in pairs(listOfRows) do
284 -- for _,row in pairs(rows) do
285 -- xcols = row.of.cols.xs
286 -- labels[ push(all,row.id) ] = label end end
287 -- return TREE(all, xcols, SYM, function(row) return labels[row.id] end) end
288
289 local _ranges, _merge
290 function _ranges(i,rows,xcol,yklass,y)
291 local n,list, dict = 0,{}, {}
292 for _,row in pairs(rows) do
293 local v = row.cells[xcol.at]
294 if v ~= "?" then
295 n = n + 1
296 local pos = xcol:bin(v)
297 dict[pos] = dict[pos] or push(list, RANGE(v,v, yklass(xcol.at, xcol.txt)))
298 dict[pos]:add(v, y(row)) end end
299 list = sort(list, lt"lo")
300 list = xcol.is=="NUM" and _merge(list, n"THE.min) or list
301 return (ranges = list,
302 div = sun(list,function(z) return z.ys:div() * z.ys.n/n end)) end
303
304 function _merge(b4,min)
305 local j,t a,b,c,ay,by,cy = 1,{}
306 while j <= #b4 do
307 a, b = b4[j], b4[j+1]
308 if b then
309 ay,by,cy = a.ys, b.ys, a.ys:merge(b.ys)
310 if ay.n<min or by.n<min or cy:div() <= (ay.n*ay:div()+by.n*by:div())/cy.n
311 then a = RANGE(a.xlo, b.xhi, cy)
312 j = j + 1 end end -- skip one, since it has just been merged
313 t[#t+1] = a
314 j = j + 1 end
315 if #t <= #b4 then return _merge(t,min) end
316 for j=2,#t do t[j].xlo = t[j-1].xhi end
317 t[1].xlo, t[#t].xhi = -big, big
318 return t end
319
320 -----
321 -- ## class RANGE
322 function RANGE.new(i, xlo, xhi, ys) i.xlo,i.xhi,i.ys,i.rows = xlo,xhi,ys,{} end
323 function RANGE.add(i,x,y)
324 if x < i.xlo then i.xlo = x end -- works for string or num
325 if x > i.xhi then i.xhi = x end -- works for string or num
326 i.ys:add(y) end
327
328 function RANGE._tostring(i)
329 local x, lo, hi = i.ys.txt, i.xlo, i.xhi
330 if lo == hi then return fmt("%s<=%s",x, lo)
331 elseif hi == big then return fmt("%s<=%s",x, lo)
332 elseif lo == -big then return fmt("%s<=%s",x, hi)
333 else return fmt("%s<=%s<=%s",lo,x,hi) end end

```

```

329 -----
330 -- ## Tests
331 local no,go = {},{}
332 function go.the() print(1); print(THE); return type(THE.p)=="number" and THE.p==2 end
333
334 function go.argmax( t,fun)
335   fun=function(x) return -x end
336   t={50,40,0,40,50}
337   return 3 == argmax(t,fun) end
338
339 function go.num(n) n=NUM(); for x=1,100 do n:add(x) end; return n.mu==50.5 end
340
341 function go.sym(s)
342   s=SYM(); for _,x in pairs{"a","a","a","a","b","b","c"} do s:add(x) end
343   return s.mode=="a" end
344
345 function go.csv( n,s)
346   n,s=0,0; for row in csv(THE.file) do n=n+1; if n>1 then s=s+row[1] end end
347   return rnd(s/n,3) == 5.441 end
348
349 function go.rows( rows)
350   doRows(THE.file,function(t) if rows then rows:add(t) else rows=ROWS(t) end end)
351   return rnd(rows.cols.ys[1].sd,0)==847 end
352
353 function go.nb()
354   return 268 == #NB("../data/diabetes.csv").dict["positive"].rows end
355
356 local function _classify(file)
357   local Abcd=require"abcd"
358   local abcd=Abcd()
359   NB(file, function(got,want) abcd:add(got,want) end)
360   abcd:pretty(abcd:report())
361   return true end
362
363 function go.soybean() return _classify("../data/soybean.csv") end
364 function go.diabetes() return _classify("../data/diabetes.csv") end
365 -----
366 -- ## Start
367 if pcall(debug.getlocal, 4, 1)
368 then return {ROW=ROW, ROWS=ROWS, NUM=NUM, SYM=SYM, THE=THE, lib=lib}
369 else THE = cli(THE,help)
370   demos(THE,go) end
371

```