

```

1 local help= []
2
3 SHORTR: semi-supervised multi-objective optimization
4 (c) 2022 Tim Menzies tim@ieee.org BSD2 license
5
6 Explore N points via  $O(\log^2(N))$  evaluations. Generate a
7 human-readable summary of that space. In pass1, find and
8 eval two distant points using multi-objective criteria.
9 Everything nearest the worst is pruned and we recurse on
10 the rest. This algorithm is only approximate so, in pass2,
11 we do it all again, starting with the better items seen in
12 pass1. Explain the final results with a decision tree that
13 recursively discretizes numerics (on their ability to
14 distinguish best/worst things found in pass2).
15
16 USAGE:
17 lua shortr.lua [OPTIONS]
18
19 OPTIONS:
20
21 OPTIONS:
22 -M --Min min size of space = .5
23 -b --bins max number of bins = 16
24 -k --k Bayes hack: low attribute frequency = 2
25 -m --m Bayes hack: low class frequency = 1
26 -s --some max number of nums to keep = 256
27 -w --wait wait this number before testing = 10
28
29 OPTIONS (other):
30 -f --file file = data/auto93.csv
31 -g --go start-up goal = nothing
32 -h --help show help = false
33 -s --seed seed = 10015]]
34
35 -----
36 -- ## Names
37 --
38 -- 'the' stores settings for this code. <p>As to the classes used by this system:
39 --
40 -- 'Row' hold the 'cells' or record and a pointer ('of') back to the
41 -- container that made them.
42 -- 'Col' summarizes columns. One 'Col' can be for
43 -- numerics or symbolic columns (denoted with 'aCol.nums').
44 -- 'Data' holds many 'Row's, summarized in a table 'aData.cols'.
45 -- (where 'aData.cols.x' holds independent columns and
46 -- 'aData.cols.y' holds dependent columns).
47 -- 'Bin' is a helper class that summarizes what dependent 'ys' values are
48 -- found between 'lo' and 'hi' of an independent column.
49 -- 'NB' is an application class that implements a Naive Bayes classifier.
50 local b4={}; for x, _ in pairs(_ENV) do b4[x]=x end
51 local _ = require"lib"
52 local Abcd = require"abcd"
53
54 local argmax,atom,big,cli,csv,demos = _argmax,_atom,_big,_cli,_csv,_demos
55 local fmt,it,map,o,oo,per,push = _fmt,_it,_map,_o,_oo,_per,_push
56 local R,sort,splice,sum = _R,_sort,_splice,_sum
57
58 local the={}
59 help:sub("[^-][^%s+)](%u)%s%([%s+)]",function(key,x) the[key] = atom(x) end)
60
61 local Col,Data,Row,Bin,NB = {}, {}, {}, {}

```

```

62 -----
63 -- ## class Col
64 -- Summaries a column of data. Uses different types for numeric or other data.
65 --
66 --> .NEW(at?:int, txt?:?str) :Col -> constructor of columns.
67 -- '.ok' is set to false after every update then set back
68 -- to true if ever we update the columns (see 'Col.ok').
69 function Col.NEW(at,txt)
70     return {n = 0, at=at or 0, txt=txt or "",
71             ok = false, kept={},
72             div=0, mid=0} end
73
74 --> .NUM(at?:int, txt?:?str) :Col -> constructor, specialized for numerics.
75 -- Numbers have a weight (-1,1) as well as the mandate to keep
76 -- no more than 'aNum.nums' samples.
77 function Col.NUM(at,txt,some, i)
78     i = Col.NEW(at,txt) -- numerics are an extension to general columns.
79     i.w = Col.WEIGHT(txt)
80     i.nums = some or the.some -- if non-nil the i.nums is a numeric
81     return i end
82 -- ## Factory to make Cols
83
84 --> .GOAL(x:[str]) :bool ->
85 --> .NUMP(x:[str]) :bool ->
86 --> .KLASS(x:[str]) :bool ->
87 --> .SKIP(x:[str]) :bool -> recognize different column types
88 function Col.GOAL(x) return (x or ""):find("[+-]$" end
89 function Col.NUMP(x) return (x or ""):find("[A-Z]" end
90 function Col.KLASS(x) return (x or ""):find"$" end
91 function Col.SKIP(x) return (x or ""):find"$" end
92
93 --> .WEIGHT(x:[str]) :(-1|1) -> assign column weight.e.g. "-1" means "minimize",
94 function Col.WEIGHT(x) return (x or ""):find"$" and -1 or 1 end
95
96 --> .COLS(names:[str]) :tab -> constructor (builds 'Col's from list of 'names').
97 -- Returns a table that stores dependents in '.y', independents in '.x',
98 -- the klass (if it exists) in '.klass'. Caveat:
99 -- only if we are not '.SKIP()'ping them.
100 function Col.COLS(names)
101     local i={x={}, y={}, names=names, klass=nil}
102     for at,txt in pairs(names) do
103         local new = Col.NUMP(txt) and Col.NUM(at,txt) or Col.NEW(at,txt)
104         if not Col.SKIP(txt) then
105             push(Col.GOAL(txt) and i.y or i.x, new)
106             if Col.KLASS(txt) then i.klass=new end end end
107         return i end
108     -- ## Update
109
110 --> .add(i:Col, v:any, inc?:int) :Col -> update 'i' with 'v' ( inc times)
111 -- Numeric columns keep a sample of the numbers while other columns track the
112 -- frequency of symbols seen so far. The larger the sample, the less often
113 -- we update the numerics.
114 function Col.add(i,v,inc)
115     inc = inc or 1
116     if v == ""
117     then i.n = i.n + inc
118         if i.nums
119         then for _=1,inc do
120             if #i.kept < i.nums then i.ok=false;push(i.kept,v)
121             elseif R() < i.nums/i.n then i.ok=false;i.kept[R(#i.kept)]=v end end
122         else i.ok = false
123             i.kept[v] = inc + (i.kept[v] or 0) end end
124         return i end
125     -- ## Computing derived properties
126
127 --> .ok(i:Col) -> ensure that the current contents are up to date. Returns 'kept'.
128 -- E.g. update 'mid'dle' and 'div'ersity (median and standard deviation for
129 -- numerics; and 'mode' and 'entropy' for others).<p>To understand the idiom
130 -- "(per(.9) - per(.1))/2.56", recall that &pm;1 and &pm;2 standard deviations
131 -- select 66.66% and 95% of the mass. In between (at &pm;1.28), we get to 90%. So to
132 -- find 1 standard deviation, divide 90th - 10th percentile by twice 1.28 (2.56).
133 function Col.ok(i)
134     if not i.ok
135     then i.div, i.mid = 0, 0
136         if i.nums
137         then i.kept = sort(i.kept) -- very fast since "kept" is small
138             i.mid = per(i.kept, .5) -- median
139             i.div = (per(i.kept, .9) - per(i.kept, .1)) / 2.56 -- stdev
140         else local most = -1 -- find the mode and ent
141             for x,n in pairs(i.kept) do
142                 if n > most then most, i.mid = n, x end
143                 if n > 0 then i.div=i.div - n/i.n*math.log(n/i.n,2) end end end end
144             i.ok = true
145         return i.kept end
146     -- ## Querying
147 -- Most of these need to call 'Col.ok()' first (to ensure column is up to date).
148
149 --> .lo(i:Col) :num ->
150 --> .hi(i:Col) :num ->
151 --> .div(i:Col) :num ->
152 --> .mid(i:Col) :any -> 'lo'west number, 'hi'ghest number, 'div'ersity, 'mid'dle numb
153     er.
154 function Col.lo(i) Col.ok(i); return i.kept[1] end
155 function Col.hi(i) Col.ok(i); return i.kept[#i.kept] end
156 function Col.div(i) Col.ok(i); return i.div end
157 function Col.mid(i) Col.ok(i); return i.mid end
158
159 --> .norm(i:Col,x:num) :0..1 -> normalize 'x' 0..1 for lo..hi.
160 function Col.norm(i,x)
161     local a=Col.ok(i); return a[a]-a[1] < 1E-9 and 0 or (x-a[1])/(a[a]-a[1]) end
162 -- ## For Discretization
163
164 --> .bin(i:Col,x:any) :any -> round numeric 'x' to nearest '(hi-lo)/the.bins'
165 -- (and for non-numerics, just return 'x').
166 function Col.bin(i,x)
167     if i.nums then
168         local lo,hi = Col.lo(i), Col.hi(i)
169         local b=(hi - lo)/the.bins
170         x = lo==hi and 1 or math.floor(x/b+.5)*b end
171     return x end
172
173 --> .bin(i:Col,j:Col) :Col -> returns a combination of two columns.
174 function Col.merge(i,j, k)
175     k = (i.nums and Col.NUM or Col.NEW)(i.at, i.txt)
176     for _ ,kept in pairs(i.kept, j.kept) do
177         for v,inc in pairs(kept) do Col.add(k,v,inc) end end
178     return k end
179
180 --> .simpler(i:col,this:col,that:col):bool->am 'i' simpler than 'this' and 'that'?
181 function Col.simpler(i,this,that)
182     print(o(i.ys),o(this.ys),o(that.ys))
183     return Col.div(i) <= (this.ys.n*Col.div(this)+that.ys.n*Col.div(that))/i.ys.n end
184 -- ## For Naive Bayes
185 function Col.like(i,x,prior)

```

```

166 if i.nums
167 then local sd,mu=Col.div(i), Col.mid(i)
168     return sd==0 and (x==mu and 1 or 0) or
169         math.exp(-1*(x - mu)^2/(2*sd^2)) / (sd*(2*math.pi)^0.5))
170 else return ((i.kept[x] or 0)+the.m*prior)/(i.n+the.m) end end

```

```

191 -- ## Row
192 function Row.NEW(of,cells) return {of=of,cells=cells,evaluated=false} end
193
194 function Row.better(i,j)
195   local s1, s2, ys = 0, 0, i.of.cols.y
196   for _,c in pairs(ys) do
197     local x,y = i.cells[c.at], j.cells[c.at]
198     x,y = Col.norm(c, x), Col.norm(c, y)
199     s1 = s1 - 2.7183*(c.w * (x-y)/#ys)
200     s2 = s2 - 2.7183*(c.w * (y-x)/#ys) end
201   return s1/#ys < s2/#ys end
202
203 function Row.klass(i) return i.cells[i.of.cols.klass.at] end
204
205 -----
206 -- ## Data
207 function Data.NEW(t) return {rows={}, cols=Col.COLS(t)} end
208
209 function Data.ROWS(src,fun)
210   if type(src)=="table" then for _,t in pairs(src) do fun(t) end
211   else for t in csv(src) do fun(t) end end end
212
213 function Data.LOAD(src, i)
214   Data.ROWS(src,function(t)
215     if i then Data.add(i,t) else i=Data.NEW(t) end end); return i end
216
217 function Data.clone(i,init)
218   local j=Data.NEW(i.cols.names)
219   for _,t in pairs(init or {}) do Data.add(j,t) end; return j end
220
221 function Data.add(i,t)
222   t = t.cells and t or Row.NEW(i,t)
223   push(i.rows, t)
224   for _,cols in pairs(i.cols.x, i.cols.y) do
225     for _,c in pairs(cols) do Col.add(c, t.cells[c.at]) end end
226   return t end
227
228 function Data.mids(i,cols)
229   local t={}
230   for _,c in pairs(cols or i.cols.y) do t[c.txt] = Col.mid(c) end;return t end
231
232 function Data.like(i,row, nklasses, nrows)
233   local prior,like,inc,x
234   prior = (#i.rows + the.k) / (nrows + the.k * nklasses)
235   like = math.log(prior)
236   for _,col in pairs(i.cols.x) do
237     x = row.cells[col.at]
238     if x and x ~= "" then
239       inc = Col.like(col,x,prior)
240       like = like + math.log(inc) end end
241   return like end
242
243 -----
244 -- ## NB
245 function NB.NEW(src,report)
246   local i = {overall=nil, dict={}, list={}}
247   report = report or print
248   Data.ROWS(src, function(row)
249     if not i.overall then i.overall = Data.NEW(row) else -- (0) eat row1
250     row = Data.add(i.overall, row) over XX add to overall
251     if #i.overall.rows > the.wait then report(Row.klass(row), NB.guess(i,row)) end
252     NB.train(i,row) end end) -- add tp rows's klass
253   return i end
254
255 function NB.train(i,row)
256   local k1 = Row.klass(row)
257   i.dict[k1] = i.dict[k1] or push(i.list, Data.clone(i.overall)) -- klass is known
258   i.dict[k1].txt = k1 -- each klass knows its name
259   Data.add(i.dict[k1],row) end -- update klass with row
260
261 function NB.guess(i,row)
262   return argmax(i.dict,
263     function(klass) return Data.like(klass,row,#i.list,#i.overall.rows) end) end
264
265 -----
266 -- ## Bin
267 function Bin.NEW(xlo, xhi, ys) return {lo=xlo, hi=xhi, ys=ys} end
268 function Bin.add(i,x,y)
269   i.lo = math.min(i.lo, x)
270   i.hi = math.max(i.hi, x)
271   Col.add(i.ys, y) end
272
273 function Bin.merge(i,j, min)
274   local k = Col.merge(i,j)
275   if i.ys.n < min or j.ys.n < min or Col.simpler(k,i,j) then return k end end
276
277 function Bin.BINS(listOfRows,col)
278   local n,list, dict = 0, {}, {}
279   for label,rows in pairs(listOfRows) do
280     for _,row in pairs(rows) do
281       local v = row.cells[col.at]
282       if v ~= "" then
283         n = n + 1
284         local pos = Col.bin(col,v)
285         print(pos, v)
286         dict[pos] = dict[pos] or push(list, Bin.NEW(v,v,Col.NEW(col.at,col.txt)))
287       end
288     end
289   end
290   list = sort(list, function(a,b) return a.txt < b.txt end)
291   list = col.names and Bin.MERGES(list, n*the.Min) or list
292   return {bins= list,
293     div = sum(list,function(z) return Col.div(z.ys)*z.ys.n/n end)} end
294
295 function Bin.MERGES(b4, min)
296   local n,now = 1, {}
297   while n <= #b4 do
298     local merged = n<#b4 and Bin.merge(b4[n], b4[n+1], min)
299     now[#now+1] = merged or b4[n]
300     n = n + (merged and 2 or 1) end
301   if #now < #b4
302   then return Bin.MERGES(now,min) -- loop to look for other merges
303   else -- stretch the bins to cover any gaps from minus infinity to plus infinity
304     for n=2,#now do now[n].lo = now[n-1].hi end
305     now[1].lo, now[#now].hi = -big, big
306     return now end end

```

```

303 -----
304 -- To disable a test, relabel it from 'Go' to 'No'.
305 local Go,No = {},{}
306
307 function Go.BINS( i,t,m,left,right)
308   i= Data.LOAD(the.file)
309   t= sort(i.rows,Row.better)
310   m= (#t)^.5
311   left = splice(t,1,m)
312   right= splice(i.rows,#t - m)
313   for n,col in pairs(i.cols.x) do
314     if n==1 then
315       Bin.BINS((left,right),col)
316     end
317     print(n,col.txt, col.nums and true or false) end
318   return true end
319
320 function Go.THE() oo(the); return true end
321
322 function Go.ROWS( d)
323   Data.ROWS(the.file,function(row)
324     if not d then d=Data.NEW(row) else
325       Data.add(d,row) end end)
326   oo(Data.mids(d))
327   return true end
328
329 function Go.STATS()
330   oo(Data.mids(Data.LOAD(the.file) )); return true end
331
332 function Go.ORDER( i,t,m,left,right)
333   i= Data.LOAD(the.file)
334   t= sort(i.rows,Row.better)
335   m= (#t)^.5
336   left = Data.clone(i,splice(t,1,m))
337   right= Data.clone(i,splice(i.rows,#t - m))
338   print("all", o(Data.mids(i)))
339   print("hes", o(Data.mids(left)))
340   print("res", o(Data.mids(right)))
341   return true end
342
343 function Go.DIABETES(f, i,t,a)
344   a = Abcd.NEW()
345   NB.NEW(f or "data/diabetes.csv",function(x,y) Abcd.add(a,x,y) end)
346   Abcd.pretty(f,Abcd.report(a))
347   return true end
348
349 function Go.SOYBEAN() return Go.DIABETES("data/soybean.csv") end
350
351 -----
352 if pcall(debug.getlocal, 4, 1)
353 then return {DATA=DATA,ROW=ROW, COL=COL, the=the,lib=lib}
354 else the = cli(the,help)
355   demos(the,Go) end

```