

```

1  --- vim: ts=2 sw=2 et:
2  local b4,help = {},{}
3  SAW2: best or rest multi-objective optimization.
4  (c) 2022 Tim Menzies, tim@ieee.org
5  "I think the highest and lowest points are the important ones.
6  Anything else is just...in between." ~ Jim Morrison
7
8  USAGE: lua saw2.lua [OPTIONS]
9
10 OPTIONS:
11 -b --bins      max bins          = 16
12 -s --seed      random number seed = 10019
13 -S --some      number of nums to keep = 256
14 -p --p         distance coefficient = 2
15
16 OPTIONS (other):
17 -f --file      where to find data  = ../etc/data/auto93.csv
18 -h --help      show help          = false
19 -r --rnd        rounding rules     = %5.2f
20 -g --go         start up action    = nothing
21
22 Usage of the works is permitted provided that this instrument is
23 retained with the works, so that any entity that uses the works is
24 notified of this instrument.  DISCLAIMER:THE WORKS ARE WITHOUT WARRANTY. ]] --[[
25
26 ## Coding Conventions
27
28 - _Separate policy from mechanism:
29 - All "magic parameters" that control code behavior should
30   be part of that help text. Allow for '-h' on the command line
31   to print help.  Parse that string to set the options.
32 - _Encapsulation: Use polymorphic but not inheritance (simpler debugging).
33   Use UPPERCASE for class names. All classes need a 'new' constructor.
34 - _Dialogue independence: Isolate and separate operating system interaction.
35 - _Test-driven development: The 'go' functions store tests.
36   Tests should be silent unless they fail. -tests can be disabled by
37   renaming from 'go.fun' to 'no.fun'. Tests should return 'true' if the
38   test passes. On exit, return number of failed tests.
39 - Less is more: Code 80 chars wide, or less. Functions in 1 line,
40   if you can. Indent with two spaces. Divide code into 120 line (or less) pages.
41   Use '\.' instead of 'self'. Use '\.' to denote the last created class/
42   Use '\.' for anonymous variable.s
43   Minimize use of local (exception: define all functions as local
44   at top of file).
45
46 ## About the Learning
47
48 - Beware missing values (marked in "?" ) and avoid them
49 - Where possible all learning should be incremental.
50 - XXX tables very useful
51 - XXX table have cols. cols are num, syms. ranges --]]
52
53 --- ## Namespace
54 local the={}
55 local _,big,clone, csv,demos,discretize,dist, eg,entropy,fmt,gap,like,lt
56 local map,merged,mid,mode,mu,norm, o,obj, oo,pdf,per,push
57 local rand, range, rangeB4, rnd, rnds, rowB4, slice, sort, some, same, sd, string2thing, sym, t
58 hese
59 local NUM, SYM, RANGE, EGS, COLS, ROW
60 for k, _ in pairs ( _ENV) do b4[k]=k end -- At end, use 'b4' to find rogue vars.

```

```

61 -----
62 -- ## Utils
63 -- Misc
64 big=math.huge
65 rand=math.random
66 fmt=string.format
67 same = function(x) return x end
68
69 -- Sorting
70 function sort(t,f)      table.sort(#t>0 and t or map(t,same), f); return t end
71 function lt(x)          return function(a,b) return a[x] < b[x] end end
72
73 -- Query and update
74 function map(t,f, u)    u={};for k,v in pairs(t) do u[1+#u]=f(v) end; return u end
75 function push(t,x)      t[1+#t]=x; return x end
76 function slice(t,i,j,k, u)
77   i,j = i or 1, j or #t
78   k = (k or 1)
79   k = (j - i) / n
80   u={}; for n=1,j,k do u[1+#u] = t[n] end return u end
81
82 -- "Strings 2 things" coercion.
83 function string2thing(x)
84   x = x:match("%s*(-)%s*")
85   if x=="true" then return true elseif x=="false" then return false end
86   return math.tointeger(x) or tonumber(x) or x end
87
88 function csv(csvfile)
89   file = io.input(csvfile)
90   return function(line, row)
91     line=io.read()
92     if not line then io.close(csvfile) else
93       row={} for x in line:match("([^\r\n]*)") do push(row,string2thing(x)) end
94       return row end end
95
96 -- "Things 2 strings" coercion.
97 function oo(t) print(o(t)) end
98 function o(t, u)
99   if #t>0 then return ("%.table.concat(map(t,tostring),"").")..t else
100     u={}; for k,v in pairs(t) do u[1+#u] = fmt("%.5s",k,v) end
101     return (t.is or "")..t["%.table.concat(sort(u),"")."] end
102
103 function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
104 function rnd(x,f)
105   return fmt(type(x)=="number" and (x-x//1 and f or the.rnd) or"%s",x) end
106
107 -- Polymorphic objects.
108 function obj(name, t,new)
109   function new(kl,...)
110     local x=metatable({},kl); kl.new(x,...); return x end
111     t = {__tostring=o, is=name or ""}; t.__index=t
112     _ = t
113     return setmetatable(t, (__call=new)) end
114
115 -----
116 -- ## Objects
117
118 -- ### NUM
119 -- For a stream of 'add'itions, incrementally maintain 'mu,sd'.
120 -- 'Norm'alize data for distance and discretization calcs (see 'dist' and 'bin')
121
122 -- Comment on 'like'lihood that something belongs to this distribution.
123 NUM=obj"NUM"
124 function _new(i,at,txt)
125   i.at=at or 0; i.txt=txt or ""; i.lo,i.hi=big, -big
126   i.n,i.mu,i.m2,i.sd = 0,0,0,0; i.w=(txt or ""):find"-$" and -1 or 1 end
127
128 function _add(i,x, d)
129   if x=="?" then return x end
130   i.n = i.n + 1
131   d = x - i.mu
132   i.mu = i.mu + d/i.n
133   i.m2 = i.m2 + d*(x - i.mu)
134   i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n - 1))^0.5)
135   i.lo = math.min(i.lo,x)
136   i.hi = math.max(i.hi,x) end
137
138 function _bin(i,x,n, b) b=(i.hi-i.lo)/n; return math.floor(x/b+0.5)*b end
139 function _mid(i) return i.mu end
140
141 function _norm(i,x) return i.hi-i.lo<1E-9 and 0 or (x-i.lo)/(i.hi-i.lo+1/big) end
142
143 function _dist(i, x,y)
144   if x=="?" and y=="?" then return 1 end
145   if x=="?" then y = i:norm(y); x = y<.5 and 1 or 0
146   elseif y=="?" then x = i:norm(x); y = x<.5 and 1 or 0
147   else x,y = i:norm(x), i:norm(y) end
148   return math.abs(x - y) end
149
150 function _like(i,x,_, e)
151   return (x < i.mu - 4*i.sd and 0 or x > i.mu + 4*i.sd and 0 or
152     2.7183^(-(x - i.mu)^2 / (z + 2*i.sd^2))/(z + (math.pi*2*i.sd^2)^.5)) end

```

```

151 -----
152 -- ### SYM
153 -- For a stream of 'add'itions, incrementally maintain count of 'all' symbols.
154 -- Using that info, report 'dist', mode ('mid') symbol, and entropy
155 -- ('div') of this distribution.
156 -- Comment on 'like'lihood that something belongs to this distribution.
157 SYM=obj"SYM"
158 function _new(i,at,txt) i.at=at or 0; i.txt=txt or ""; i.n,i.all = 0,{} end
159 function _add(i,x,n)
160   if x=="?" then return x end
161   i.n=i.n+1; i.all[x] = (n or 1) + (i.all[x] or 0) end
162
163 function _dist(i,x,y) return (a==b and 0 or 1) end
164
165 function _mid(i)
166   m=0; for y,n in pairs(i.all) do if n>m then m,x=n,y end end; return x end
167
168 function _div(i, n,e)
169   e=0; for k,n in pairs(i.all) do e=e-n/i.n*math.log(n/i.n,2) end; return e end
170
171 function _like(i,x,prior) return (c.all[x] or 0) + the.m*prior/(c.n+the.m) end
172
173 -----
174 -- ### RANGE
175 -- For a stream of 'add'itions, incrementally maintain counts of 'x' and 'y'.
176 -- Summarize 'x' as the 'lo,hi' seen so far and summarize 'y' in 'SYM' counts
177 -- in 'y.all' (and get counts there using 'of').
178 -- Support range sorting ('_lt') and printing ('_tostring').
179 -- Check if this range's 'x' values select's for a particular row.
180 -- 'Merge' adjacent ranges if the entropy of the whole is less than the parts.
181 RANGE=obj"RANGE"
182 function _new(i,col,lo,hi,y)
183   i.cols, i.x, i.y = col, ({lo=lo or big, hi=hi or -big}), (y or SYM()) end
184
185 function _add(i,x,y)
186   if x=="?" then return x end
187   i.x.lo = math.min(i.x.lo,x)
188   i.x.hi = math.max(i.x.hi,x)
189   i.y:add(x,y) end
190
191 function _lt(i,j) return i.col.at == j.col.at and i.x.lo < j.x.lo end
192 function _of(i,x) return i.y.all[x] or 0 end
193
194 function _selects(i,t, x)
195   t = t.cells and t.cells or t
196   x = t[i.at]
197   return x=="?" or (i.x.lo==i.x.hi and i.x.lo==x) or (i.x.lo<=x and x<i.x.hi) end
198
199 function _tostring(i)
200   local x, lo, hi = i.txt, i.x.lo, i.x.hi
201   if lo == hi then return fmt("%s==%s",x, lo)
202   elseif hi == big then return fmt("%s>=%s",x, lo)
203   elseif lo == -big then return fmt("%s<=%s", x, hi)
204   else return fmt("%s<=%s<%s",lo,x,hi) end end
205
206 function _merge(i,j,n0, k)
207   if i.at == j.at then
208     k = SYM(i.y.at, i.y.txt)
209     i,j = i.y, j.y
210     for x,n in pairs(i.all) do sym(k,x,n) end
211     for x,n in pairs(j.all) do sym(k,x,n) end
212     if i.y.n<(n0 or 0) or j.y.n<(n0 or 0) or (ent(i)*i.n+ent(j)*j.n)/k.n > ent(k)
213       then return RANGE(i.col, i.lo, j.hi, k) end end
214
215 -----
216 -- ### ROW
217 -- Using knowledge of the 'base' geometry of the data, support distance calcs
218 -- i ('_sub' and 'around') as well as multi-objective ranking ('_lt').
219 ROW=obj"ROW"
220 function _new(i, eg, cells) i.base,i.cells = eg,cells end
221 function _lt(i,j, sl,s2,e,y,a,b)
222   y = i.base.cols.y
223   sl, s2, e = 0, 0, math.exp(1)
224   for _ ,col in pairs(y) do
225     a = col:norm(i.cells[col.at])
226     b = col:norm(j.cells[col.at])
227     sl = sl - e^(col.w * (a - b) / #y)
228     s2 = s2 - e^(col.w * (b - a) / #y) end
229   return sl/#y < s2/#y end
230
231 function _sub(i,j)
232   for _ ,col in pairs(i.base.cols.x) do
233     a,b = i.cells[col.at], j.cells[col.at]
234     inc = a=="?" and b=="?" and 1 or col:dist(a,b)
235     d = d + inc*the.p end
236   return (d / (#i.base.cols.x) ^ (1/the.p) end
237
238 function _around(i,rows)
239   return sort(map(rows or i.base.rows, function(j) return (dist=i-j,row=j) end),
240     lt="dist") end

```

```

239 --- ### COLS
240 -- Factory for converting column 'names' to 'NUM's ad 'SYM's.
241 -- Store all columns in 'all', and for all columns we are not skipping,
242 -- store the independent and dependent columns distributions in 'x' and 'y'.
243 COLS=obj"COLS"
244 function _new(i, names, head, row, col)
245   i.names=names; i.all={}; i.y={}; i.x={}
246   for at,txt in pairs(names) do
247     col = push(i.all, (txt:find("[A-Z]" and NUM or SYM) (at, txt)))
248     col.goalp = txt:find("[4-5]" and true or false)
249   if not txt:find("S" then
250     if txt:find("S" then i.klass=col end
251     push(col.goalp and i.y or i.x, col) end end end
252 ---### EGS
253 -- For a stream of 'add'itions, incrementally store rows, summarized in 'cols'.
254 -- When 'add'ing, build new rows for new data. Otherwise reuse rows across
255 -- multiple sets of examples.
256 -- Supporting 'copy'ing of this structure, without or without rows of data.
257 -- Report how much this set of examples 'like' a new row.
258 -- Discretize columns as 'ranges' that distinguish two sets of rows
259 -- (merging irrelevant distinctions).
260 -- Summarize the 'mid'point of these examples.
261 EGS=obj"EGS"
262 function _new(i, names) i.rows, i.cols = {}, COLS(names) end
263 function _load(f, i)
264   for row in csv(the.file) do if i then i:add(row) else i=EGS(row) end end
265   return i end
266
267 function _add(i, row, cells)
268   cells = push(i.rows, row.cells and row or ROW(i, row)).cells
269   for n, col in pairs(i.cols.all) do col:add(cells[n]) end end
270
271 function _mid(i, cols)
272   return map(cols or i.cols.y, function(c) return c:mid() end) end
273
274 function _copy(i, rows, j)
275   j=EGS(i.cols.names); for _, r in pairs(rows or {}) do j:add(r) end; return j end
276
277 function _like(i, t, overall, nHypotheses, c)
278   prior = (#i.rows + the.k) / (overall + the.k * nHypotheses)
279   like = math.log(prior)
280   for at, x in pairs(t) do
281     c=i.cols.all[at]
282     if x=="?" and not c.goalp then
283       like = math.log(col:like(x)) + like end end
284   return like end
285
286 local _merge, _xpd, _ranges
287 function _ranges(i, one, two, t)
288   t={}; for _, c in pairs(i.cols.x) do t[c.at]=_ranges(c, one, two) end; return t end
289
290 function _ranges(col, yes, no, out, x, d)
291   out = {}
292   for _, what in pairs({rows=yes, klass=true}, {rows=no, klass=false}) do
293     for _, row in pairs(what.rows) do x = row.cells[col.at]; if x=="?" then
294       d = col:discretize(x, the.bins)
295       out[d] = out[d] or RANGE{col, x, x}
296       out[d]:add(x, what.klass) end end end
297   return _xpd(_merge(sort(out))) end
298
299 function _merge(b4, a, b, c, j, n, tmp)
300   j, n, tmp = 1, #b4, {}
301   while j<=n do
302     a, b = b4[j], b4[j+1]
303     if b then c = a:merged(b); if c then a, j = c, j+1 end end
304     tmp[#tmp+1] = a
305     j = j+1 end
306   return #tmp==#b4 and tmp or _merge(tmp) end
307
308 function _xpd(t)
309   for j=2, #t do t[j].lo=t[j-1].hi end; t[1].lo, t[#t].hi = -big, big; return t end
310
311

```

```

312 --- ### DEMOS
313 local go, no = {}, {}
314
315 -- Convert help string to a table. Check command line for any updates.
316 function these(f1, f2, k, x)
317   for n, flag in ipairs(arg) do if flag==f1 or flag==f2 then
318     x = x=="false" and true" or x=="true" and false" or arg[n+1] end end
319   the[k] = string2thing(x) end
320
321 -- Run the demos, resetting settings and random number see before each.
322 -- Return number of failures.
323 function demos( fails, names, defaults, status)
324   fails=0 -- this code will return number of failures
325   names, defaults = {}, {}
326   for k, f in pairs(go) do if type(f)=="function" then push(names, k) end end
327   for k, v in pairs(the) do defaults[k]=v end
328   if go[the.go] then names=(the.go) end
329   for _, one in pairs(sort(names)) do -- for all we want to do
330     for k, v in pairs(defaults) do the[k]=v end -- set settings to defaults
331     math.randomseed(the.seed or 10019) -- reset random number seed
332     io.stderr:write("\n")
333     status = go[one]() -- run demo
334     if status == true then
335       print("Error", one, status)
336       fails = fails + 1 end end
337   return fails end -- return total failure count
338
339 -- Simple stuff
340 function go.the() return type(the.bins)=="number" end
341 function go.sort( t) return 0==sort((100, 3, 4, 2, 10, 0))[1] end
342 function go.num( n, mu, sd)
343   n, mu, sd = NUM(), 10, 1
344   for i=1, 10^4 do
345     n:add(mu+sd*math.sqrt(-2*math.log(rand()))*math.cos(2*math.pi*rand())) end
346   return math.abs(n.mu - mu) < 0.05 and math.abs(n.sd - sd) < 0.5 end
347
348 -- Can we read rows off the disk?
349 function go.rows( n, m)
350   m, n=0, 0; for row in csv(the.file) do m=m+1; n=n+#row; end; return n/m==8 end
351
352 -- Can we turn a list of names into columns?
353 function go.cols( i)
354   i=COLS("name", "Age", "ShoeSize-")
355   return i.y[1].w == -1 end
356
357 -- Can we read data, summarized as columns?
358 function go.egs( it)
359   it = EGS.load(the.file); return math.abs(2970 - it.cols.y[1].mu) < 1 end
360
361 -- Can we discretize
362 function go.ranges( it, n, a, b)
363   it = EGS.load(the.file)
364   print(o(rnds(it:mid()))))
365   it.rows = sort(it.rows)
366   n = (#it.rows)^.5
367   a, b = slice(it.rows, 1, n), slice(it.rows, n+1, #it.rows, 3*n)
368   print(o(rnds(it:copy(a):mid()), o(rnds(it:copy(b):mid()))))
369   --oo(a:mid())
370   --oo(b:mid())
371   return math.abs(2970 - it.cols.y[1].mu) < 1 end
372

```

```

373 --- ## Main
374 -- Parse help text for flags and defaults, check CLI for updates. Maybe print
375 -- the help (with some pretty colors). Run the demos. Check for rogue vars.
376 -- Exit, reporting number of failures.
377 help:gsb("\n ([-][^%s+)]%s)+([[-][^%s+)]|^n)%s([%^s+)", these)
378 if the.help then
379   print(help:gsb("%y%u+", "%27[31m%127[0m")
380     :gsb("(%s)([-][-]?[^%s+)%s)", "%127[33m%227[0m%3*)", ""))
381 else
382   local fails = demos()
383   for k, v in pairs(_ENV) do if not b4[k] then print("?", k, type(v)) end end
384   os.exit(fails) end
385 -- function SOME() return (all={}, ok=false, n=0) end
386 -- function some(i, x)
387 --   if x=="?" then return x end
388 --   i.n = 1 + i.n
389 --   if #i.all < the.some then i.ok=false; push(i.all, x)
390 --   elseif rand() < the.some/i.n then i.ok=false; i.all[#i.all]=x end end
391 -- function per(i, p)
392 --   i.all = i.ok and i.all or sort(i.all); i.ok=true
393 --   return i.all[math.max(1, math.min(#i.all, (p or .5)*#i.all//1))] end
394

```