

```

1  --- vim: ts=2 sw=2 et :
2  local b4,help = {},{}
3  CHOP: best or rest multi-objective optimization.
4  (c) 2022 Tim Menzies, timm@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 chop.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 local _ = {}
27 local _big, clone, csv, demos, discretize, dist, eg, entropy, fmt, gap, like, lt
28 local map, merge, mid, mode, mu, norm, num, o, obj, oo, pdf, per, push, rand, range
29 local rnd, rnds, rowB4, slice, sort, some, same, sd, string2thing, sym, these
30 local NUM, SYM, RANGE, EOS, COLS, ROW
31 for k, _ in pairs(_ENV) do b4[k]=k end -- At end, use 'b4' to find rogue vars.
32
33 --- ## Coding Conventions
34 -- Separate policy from mechanism.
35 -- All "magic parameters" that control code behavior should be part
36 -- of that help text. Allow for '-h' on the command line to print
37 -- help. Parse that string to set the options.
38 -- Dialogue independence: Isolate and separate operating system interaction.
39 -- Test-driven development: The 'go' functions store tests.
40 -- Tests should be silent unless they -- fail. -tests can be
41 -- disabled by renaming from 'go.fun' to 'no.fun'. Tests should
42 -- return 'true' if the test passes. On exit, return number of
43 -- failed tests.
44 -- Less is more: Code 80 chars wide, or less. Functions in 1 line,
45 -- if you can. Indent with two spaces. Divide code into 120 line (or
46 -- less) pages. Use 'i' instead of 'self'. Use '.' to denote the
47 -- last created class/ Use '_' for anonymous variable.s Minimize
48 -- use of local (exception: define all functions as local at top of
49 -- file).
50 -- Encapsulation: Use polymorphism but no inheritance (simpler
51 -- debugging). All classes get a 'new' constructor.
52 -- Use UPPERCASE for class names.
53
54 --- ## About the Learning
55 -- Data is stored in ROWs.
56 -- Beware missing values (marked in "?") and avoid them
57 -- Where possible all learning should be incremental.
58 -- Standard deviation and entropy generalized to 'div' (diversity);
59 -- Mean and mode generalized to 'mid' (middle);
60 -- Rows are created once and shared between different sets of
61 -- examples (so we can accumulate statistics on how we are progressing
62 -- inside each row).
63 -- When a row is first created, it is assigned to a 'base'; i.e.
64 -- a place to store the 'lo,hi' values for all numerics.
65 -- XXX tables very useful
66 -- XXX table have cols. cols are num, syms. ranges

```

```

70 --- ## Utils
71 --- Misc
72 big=math.huge
73 rand=math.random
74 fmt=string.format
75 same = function(x) return x end
76
77 --- Sorting
78 function sort(t,f) table.sort(#t>0 and t or map(t,same), f); return t end
79 function lt(x) return function(a,b) return a[x] < b[x] end end
80
81 --- Query and update
82 function map(t,f, u) u={};for k,v in pairs(t) do u[1+#u]=f(v) end; return u end
83 function push(t,x) t[1+#t]=x; return x end
84 function slice(t,i,j,k, u)
85 i,j = (i or 1)/1, (j or #t)/1
86 k = (k and (j-i)/k or 1)/1
87 u={}; for n=1,j,k do u[1+#u] = t[n] end return u end
88
89 --- "Strings 2 things" coercion.
90 function string2thing(x)
91 x = x:match("%s*(-)%s*$")
92 if x=="true" then return true elseif x=="false" then return false end
93 return math.tointeger(x) or tonumber(x) or x end
94
95 function csv(csvfile)
96 csvfile = io.input(csvfile)
97 return function(line, row)
98 line=io.read()
99 if not line then io.close(csvfile) else
100 row={}; for x in line:gmatch("[^,]+") do push(row, string2thing(x)) end
101 return row end end
102
103 --- "Things 2 strings" coercion.
104 function oo(t) print(o(t)) end
105 function o(t)
106 if #t>0 then return "["..table.concat(map(t,tostring),"").."]" else
107 u={}; for k,v in pairs(t) do u[1+#u] = fmt("%s%s",k,v) end
108 return (t.is or "").."{"..table.concat(sort(u,"")..")}" end end
109
110 function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
111 function rnd(x,f)
112 return fmt(type(x)=="number" and (x-=x//1 and f or the.rnd) or "%s",x) end
113
114 --- Polymorphic objects.
115 function obj(name, t,new)
116 function new(kl,...)
117 local x=setmetatable({},kl); kl.new(x,...); return x end
118 t = {__tostring=o, is=name or ""; t.__index=t
119 _ = t
120 return setmetatable(t, {__call=new}) end
121
122 -----
123 --- ## Objects
124
125 --- ## NUM
126 -- For a stream of 'add'itions, incrementally maintain 'mu,sd'.
127 -- Norm'alize data for distance and discretization calcs
128 -- (see 'dist' and 'range').
129 -- Comment on 'like'lihood that something belongs to this distribution.
130 NUM=obj"NUM"
131 function _new(i,at,txt)
132 i.at=at or 0; i.txt=txt or ""; i.lo,i.hi=big, -big
133 i.n,i.mu,i.m2,i.sd = 0,0,0,0; i.w=(txt or ""):find("-$") and -1 or 1 end
134
135 function _add(i,x, d)
136 if x=="?" then return x end
137 i.n = i.n + 1
138 d = x - i.mu
139 i.mu = i.mu + d/i.n
140 i.m2 = i.m2 + d*(x - i.mu)
141 i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n - 1))^0.5)
142 i.lo = math.min(i.lo,x)
143 i.hi = math.max(i.hi,x) end
144
145 function _range(i,x,n, b) b=(i.hi-i.lo)/n; return math.floor(x/b+0.5)*b end
146 function _mid(i) return i.mu end
147
148 function _norm(i,x) return i.hi-i.lo<1E-9 and 0 or (x-i.lo)/(i.hi-i.lo+1/big) end
149
150 function _dist(i, x,y)
151 if x=="?" and y=="?" then return 1 end
152 if x=="?" then y = i:norm(y); x = y<.5 and 1 or 0
153 elseif y=="?" then x = i:norm(x); y = x<.5 and 1 or 0
154 else x,y = i:norm(x), i:norm(y) end
155 return math.abs(x - y) end
156
157 function _like(i,x,_, e)
158 return (x < i.mu - 4*i.sd and 0 or x > i.mu + 4*i.sd and 0 or
159 2.7183*(-(x - i.mu)^2 / (z + 2*i.sd^2)))/(z + (math.pi*2*i.sd^2)^.5)) end

```

```

160 -----
161 --- ## SYM
162 -- For a stream of 'add'itions, incrementally maintain count of 'all' symbols.
163 -- Using that info, report 'dist', mode ('mid') symbol, and entropy
164 -- ('div') of this distribution.
165 -- Comment on 'like'lihood that something belongs to this distribution.
166 -- Discretization of a symbol just returns that sym ('range').
167 SYM=obj"SYM"
168 function _new(i,at,txt) i.at=at or 0; i.txt=txt or ""; i.n,i.all = 0,{} end
169 function _add(i,x,n)
170 if x=="?" then return x end
171 n = n or 1
172 i.n=i.n+n; i.all[x] = n + (i.all[x] or 0) end
173
174 function _range(i,x,_) return x end
175 function _dist(i,x,y) return (a==b and 0 or 1) end
176
177 function _mid(i)
178 m=0; for y,n in pairs(i.all) do if n>m then m,x=n,y end end; return x end
179
180 function _div(i, n,e)
181 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
182
183 function _like(i,x,prior) return ((c.all[x] or 0) + the.m*prior)/(c.n+the.m) end
184
185 --- ## RANGE
186 -- For a stream of 'add'itions, incrementally maintain counts of 'x' and 'y'.
187 -- Summarize 'x' as the 'lo,hi' seen so far and summarize 'y' in 'SYM' counts
188 -- in 'y.all' (and get counts there using 'o.f').
189 -- Support range sorting ('_lt') and printing ('_tostring').
190 -- Check if this range's 'x' values 'select's for a particular row.
191 -- 'Merge' adjacent ranges if the entropy of the whole is less than the parts.
192 RANGE=obj"RANGE"
193 function _new(i,col,lo,hi,y)
194 i.col, i.x, i.y = col, {lo=lo or big, hi=hi or -big}, (y or SYM()) end
195
196 function _add(i,x,y)
197 if x=="?" then return x end
198 i.x.lo = math.min(i.x.lo,x)
199 i.x.hi = math.max(i.x.hi,x)
200 i.y:add(y) end
201
202 function _lt(i,j) return i.x.lo < j.x.lo end
203 function _of(i,x) return i.y.all[x] or 0 end
204
205 function _selects(i,t, x)
206 t = t.cells and t.cells or t
207 x = t[i.at]
208 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
209
210 function _tostring(i)
211 local x, lo, hi = i.col.txt, i.x.lo, i.x.hi
212 if lo == hi then return fmt("%s==%",x, lo)
213 elseif hi == big then return fmt("%s>=%s",x, lo)
214 elseif lo == -big then return fmt("%s<=%s", x, hi)
215 else return fmt("%s<=%s<=%s",lo,x,hi) end end
216
217 function _merge(i,j,n0, k)
218 k = SYM(i.col.at, i.col.txt)
219 for x,n in pairs(i.y.all) do k:add(x,n) end
220 for x,n in pairs(j.y.all) do k:add(x,n) end
221 if i.y.nc(n0 or 0) or j.y.nc(n0 or 0) or (
222 (i.y:div(i)*i.y.n + j.y:div(j)*j.y.n)/k.n >= .99*k:div())
223 then return RANGE(i.col, i.x.lo, j.x.hi, k) end end
224
225 --- ## ROW
226 -- Using knowledge of the 'base' geometry of the data, support distance calcs
227 -- i ('_sub' and 'around') as well as multi-objective ranking ('_lt').
228 ROW=obj"ROW"
229 function _new(i,eg, cells) i.base,i.cells = eg,cells end
230 function _lt(i,j, s1,s2,e,y,a,b)
231 y = i.base.cols.y
232 s1, s2, e = 0, 0, math.exp(1)
233 for _col in pairs(y) do
234 a = col:norm(i.cells[col.at])
235 b = col:norm(j.cells[col.at])
236 s1 = s1 - e^(col.w * (a - b) / #y)
237 s2 = s2 - e^(col.w * (b - a) / #y) end
238 return s1/#y < s2/#y end
239
240 function _sub(i,j)
241 for _col in pairs(i.base.cols.x) do
242 a,b = i.cells[col.at], j.cells[col.at]
243 inc = a=="?" and b=="?" and 1 or col:dist(a,b)
244 d = d + inc*the.p end
245 return (d / (#i.base.cols.x) ^ (1/the.p) end
246
247 function _around(i,rows)
248 return sort(map(rows or i.base.rows, function(j) return (dist=i-j,row=j) end),
249 lt"dist") end

```

```

250 -----
251 -- ### COLS
252 -- - Factory for converting column 'names' to 'NUM's ad 'SYM's.
253 -- - Store all columns in - 'all', and for all columns we are not skipping,
254 -- - store the independent and dependent columns distributions in 'x' and 'y'.
255 COLS=obj{"COLS"
256 function _new(i, names, head, row, col)
257   i.names=names; i.all={}; i.y={}; i.x={}
258   for at,txt in pairs(names) do
259     col = push(i.all, (txt::find("[A-Z]" and NUM or SYM) (at, txt))
260     col.goalp = txt::find("[4-5]" and true or false
261     if not txt::find("S" then
262       if txt::find("I" then i.klass=col end
263       push(col.goalp and i.y or i.x, col) end end end
264 -----
265 -- ## EGS
266 -- - For a stream of 'add'itions, incrementally store rows, summarized in 'cols'.
267 -- - When 'add'ing, build new rows for new data. Otherwise reuse rows across
268 -- - multiple sets of examples.
269 -- - Supporting 'copy'ing of this structure, without or without rows of data.
270 -- - Report how much this set of examples 'like' a new row.
271 -- - Discretize columns as 'ranges' that distinguish two sets of rows
272 -- - (merging irrelevant distinctions).
273 -- Summarize the 'mid'point of these examples.
274 EGS=obj{"EGS"
275 function _new(i, names) i.rows, i.cols = {}, COLS(names) end
276 function _load(f, i)
277   for row in csv(the.file) do if i then i:add(row) else i=EGS(row) end end
278   return i end
279
280 function _add(i, row, cells)
281   cells = push(i.rows, row.cells and row or ROW(i, row)).cells
282   for n, col in pairs(i.cols.all) do col:add(cells[n]) end end
283
284 function _mid(i, cols)
285   return map(cols or i.cols.y, function(c) return c:mid() end) end
286
287 function _copy(i, rows, j)
288   j=EGS(i.cols.names); for __, r in pairs(rows or {}) do j:add(r) end; return j end
289
290 function _like(i, t, overall, nHypotheses, c)
291   prior = (#i.rows + the.k) / (overall + the.k * nHypotheses)
292   like = math.log(prior)
293   for at, x in pairs(t) do
294     c=i.cols.all.at[at]
295     if x=="*" and not c.goalp then
296       like = math.log(col:like(x)) + like end end
297   return like end
298
299 local _merge, _xpad, _ranges
300 function _ranges(i, one, two, t)
301   t={}; for __, c in pairs(i.cols.x) do t[c.at]=_ranges(c, one, two) end; return t end
302
303 function _ranges(col, yes, no, out, x, d)
304   out = {}
305   for __, what in pairs(({rows=yes, klass=true}, {rows=no, klass=false}) do
306     for __, row in pairs(what.rows) do x = row.cells[col.at]; if x=="*" then
307       d = col.range(x, the.bins)
308       out[d] = out[d] or RANGE(col, x, x)
309       out[d]:add(x, what.klass) end end end
310   return _xpad(_merge(sort(map(out, same))), (#yes+#no)^.5) end
311   --return _xpad(_merge(sort(out))) end
312
313 function _merge(b4, min, a, b, c, j, n, tmp)
314   j, n, tmp = 1, #b4, {}
315   while j<=n do
316     a, b = b4[j], b4[j+1]
317     if b then c = a:merge(b, min); if c then a, j = c, j+1 end end
318     tmp[#tmp+1] = a
319     j = j+1 end
320   return #tmp==#b4 and tmp or _merge(tmp, min) end
321
322 function _xpad(t)
323   for j=2, #t do t[j].lo=t[j-1].hi end
324   t[1].x.lo, t[#t].x.hi = -big, big
325   return t end

```

```

326 -----
327 -- ## DEMOS
328 local go, no = {}, {}
329
330 -- Convert help string to a table. Check command line for any updates.
331 function these(f1, f2, k, x)
332   for n, flag in ipairs(arg) do if flag==f1 or flag==f2 then
333     x = x.."false" and true" or x=="true" and "false" or arg[n+1] end end
334   the[k] = string2thing(x) end
335
336 -- Run the demos, resetting settings and random number see before each.
337 -- Return number of failures.
338 function demos( fails, names, defaults, status)
339   fails=0 -- this code will return number of failures
340   names, defaults = {}, {}
341   for k, f in pairs(go) do if type(f)=="function" then push(names, k) end end
342   for k, v in pairs(the) do defaults[k]=v end
343   if go(the.go) then names=(the.go) end
344   for __, one in pairs(sort(names)) do -- for all we want to do
345     for k, v in pairs(defaults) do the[k]=v end -- set settings to defaults
346     math.randomseed(the.seed or 10019) -- reset random number seed
347     io.stderr:write(" ")
348     status = go(one)() -- run demo
349     if status == true then
350       print("Error, one, status)
351       fails = fails + 1 end end
352     return fails end
353
354 -- Simple stuff
355 function go.the() return type(the.bins)=="number" end
356 function go.sort( t) return 0==sort((100,3,4,2,10,0))[1] end
357 function go.slice( t, u)
358   t = {10,20,30,40,50,60,70,80,90,100,110,120,130,140}
359   u = slice(t, 3, #t, 3)
360   t = slice(t, 3, 5)
361   return #t==3 and #u==4 end
362
363 function go.num( n, mu, sd)
364   n, mu, sd = NUM(), 10, 1
365   for i=1, 10^4 do
366     n:add(mu+sd*math.sqrt(-2*math.log(rand()))*math.cos(2*math.pi*rand())) end
367   return math.abs(n.mu - mu) < 0.05 and math.abs(n.sd - sd) < 0.5 end
368
369 -- Can we read rows off the disk?
370 function go.rows( n, m)
371   m, n=0, 0; for row in csv(the.file) do m=m+1; n=n+#row; end; return n/m==8 end
372
373 -- Can we turn a list of names into columns?
374 function go.cols( i)
375   i=COLS{"Name", "Age", "ShoeSize"}
376   return i.y[1].w == -1 end
377
378 -- Can we read data, summarized as columns?
379 function go.egs( i, t)
380   it = EGS.load(the.file); return math.abs(2970 - it.cols.y[1].mu) < 1 end
381
382 -- Can we discretize
383 function go.ranges( it, n, best, rest, min)
384   it = EGS.load(the.file)
385   print("all", o(rnds(it:mid())) )
386   it.rows = sort(it.rows)
387   for j, row in pairs(sort(it.rows)) do row.klass = 1+j//(#it.rows*.35/6) end
388   n = (#it.rows)^.5
389   best, rest = slice(it.rows, 1, n), slice(it.rows, n+1, #it.rows, 3*n)
390   print("best", o(rnds(it:copy(best):mid())))
391   print("rest", o(rnds(it:copy(rest):mid())))
392   for __, ranges in pairs(it:ranges(best, rest)) do
393     print("
394     for at, range in pairs(ranges) do
395       print(range) end end
396     --oo(a:mid())
397     --oo(b:mid())
398   return math.abs(2970 - it.cols.y[1].mu) < 1 end

```

```

399 -----
400 -- ## Main
401
402 -- Parse help text for flags and defaults, check CLI for updates.
403 -- Maybe print the help (with some pretty colors).
404 -- Run the demos.
405 -- Check for rogue vars.
406 -- Exit, reporting number of failures.
407 help:gsub("u ([-]|[%s+])%s)+([-]|[-]|[%s+])^u)%s(%s)+", these)
408 if the.help then
409   print(help:gsub("%u%u+", "%27[31m%127[0m"
410     :gsub("(%s)([-]|[-]|[%s+])%s)", "%127[33m%27[0m%3]", ""))
411 else
412   local fails = demos()
413   for k, v in pairs(_ENV) do if not b4[k] then print("?", k, type(v)) end end
414   os.exit(fails) end
415
416 -- function SOME() return (all={}, ok=false, n=0) end
417 -- function some(i, x)
418 --   if x=="*" then return x end
419 --   i.n = 1 + i.n
420 --   if #i.all < the.some then i.ok=false; push(i.all, x)
421 --   elseif rand() < the.some/i.n then i.ok=false; i.all[#i.all]=x end end
422
423 -- function per(i, p)
424 --   i.all = i.ok and i.all or sort(i.all); i.ok=true
425 --   return i.all[math.max(1, math.min(#i.all, (p or .5)*#i.all//1))] end

```