

```

1 --- vim: ts=2 sw=2 et:
2 local b4,help = {},{[[
3 CHOP: 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 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 = 45.2f
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 be part
30 -- of that help text. Allow for '-h' on the command line to print
31 -- help. Parse that string to set the options.
32 -- - Dialogue independence.: Isolate and separate operating system interaction.
33 -- - Test-driven development.: The 'go' functions store tests.
34 -- Tests should be silent unless they -- fail. -tests can be
35 -- disabled by renaming from 'go.fun' to 'no.fun'. Tests should
36 -- return 'true' if the test passes. On exit, return number of
37 -- failed tests.
38 -- - Less is more: Code 80 chars wide, or less. Functions in 1 line,
39 -- if you can. Indent with two spaces. Divide code into 120 line (or
40 -- less) pages. Use '_' instead of 'self'. Use '_' to denote the
41 -- last created class/ Use '_' for anonymous variable.s Minimize
42 -- use of local (exception: define all functions as local at top of
43 -- file).
44 -- - Encapsulation.: Use polymorphism but no inheritance (simpler
45 -- debugging). All classes get a 'new' constructor.
46 -- Use UPPERCASE for class names.
47
48 ---
49 -- ## About the Learning
50
51 -- - Data is stored in ROWS
52 -- - Beware missing values (marked in "?") and avoid them
53 -- - Where possible all learning should be incremental.
54 -- - Standard deviation and entropy generalized to 'div' (diversity);
55 -- - Mean and mode generalized to 'mid' (middle);
56 -- - Rows are created once and shared between different sets of
57 -- examples (so we can accumulate statistics on how we are progressing
58 -- inside each row).
59 -- - When a row is first created, it is assigned to a 'base'; i.e.
60 -- a place to store the 'lo,hi' values for all numerics.
61 -- - XXX tables very usefual
62 -- - XXX table have cols. cols are num, syms. ranges
63
64 ---
65 -- ## Namespace
66 local thes={}
67 local _big,clone,csv,demos,discretize,dist,eg,entropy,fmt,gap,like,lt
68 local map,merged,mid,mode,mu,norm,num,o,obj,oo,pdf,per,push
69 local rand,range,rangeB4,rnd,rnds,rowB4,slice,sort,some,same,sd,string2thing,sym,t
70 local NUM,SYM,RANGE,EGS,COLS,ROW
71 for k, _ in pairs(_ENV) do b4[k]=k end -- At end, use 'b4' to find rogue vars.

```

```

72 -- ## Utils
73 -- Misc
74 big=math.huge
75 rand=math.random
76 fmt=string.format
77 same = function(x) return x end
78
79 -- Sorting
80 function sort(t,f) table.sort(#t>0 and t or map(t,same), f); return t end
81 function lt(x) return function(a,b) return a[x] < b[x] end end
82
83 -- Query and update
84 function map(t,f, u) u={};for k,v in pairs(t) do u[1+#u]=f(v) end; return u end
85 function push(t,x) t[1+#t]=x; return x end
86 function slice(t,i,j,k, u)
87 i,j = i or 1,j or #t
88 k = (k or 1)
89 k = (j - i)/k
90 u={}; for n=i,j,k//1 do u[1+#u] = t[n] end return u end
91
92 -- "Strings 2 things" coercion.
93 function string2thing(x)
94 x = x:match("%s*(.-)%s*")
95 if x=="mu" then return true elseif x=="false" then return false end
96 return math.tointeger(x) or tonumber(x) or x end
97
98 function csv(csvfile)
99 csvfile = io.input(csvfile)
100 return function(line, row)
101 line=io.read()
102 if not line then io.close(csvfile) else
103 row={}; for x in line:match("([^\,]*)") do push(row,string2thing(x)) end
104 return row end end
105
106 -- "Things 2 strings" coercion.
107 function oo(t) print(o(t)) end
108 function o(t, u)
109 if #t>0 then return "["..table.concat(map(t,tostring)," " ..").."]" else
110 u={}; for k,v in pairs(t) do u[1+#u] = fmt("%s%s",k,v) end
111 return (t.is or "").."["..table.concat(sort(u)," " ..").."]" end
112
113 function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
114 function rnd(x,f)
115 return fmt(type(x)=="number" and (x-x//1 and f or the.rnd) or"%s",x) end
116
117 -- Polymorphic objects.
118 function obj(name, t,new)
119 function new(kl,...)
120 local x=setmetatable({},kl); kl.new(x,...); return x end
121 t = {__tostring=0, is=name or ""}; t.__index=t
122 _ = t
123 return setmetatable(t, {__call=new}) end
124
125 --- ## Objects
126
127 -- ### NUM
128 -- - For a stream of 'add'itions, incrementally maintain 'mu,sd'.
129 -- - 'Norm'alize data for distance and discretization calcs (see 'dist' and 'bin')
130
131 -- Comment on 'like'lihood that something belongs to this distribution.
132 NUM=obj("NUM")
133 function _new(i,at,txt)
134 i,at=at or 0; i,txt=txt or ""; i.lo,i.hi=big, -big
135 i.n,i.mu,i.m2,i.sd = 0,0,0,0; i.w=(txt or ""):find("-$") and -1 or 1 end
136
137 function _add(i,x, d)
138 if x=="?" then return x end
139 i.n = i.n + 1
140 d = x - i.mu
141 i.mu = i.mu + d/i.n
142 i.m2 = i.m2 + d*(x - i.mu)
143 i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n - 1))^0.5)
144 i.lo = math.min(i.lo,x)
145 i.hi = math.max(i.hi,x) end
146
147 function _bin(i,x,n, b) b=(i.hi-i.lo)/n; return math.floor(x/b+0.5)*b end
148 function _mid(i) return i.mu end
149
150 function _norm(i,x) return i.hi-i.lo<1E-9 and 0 or (x-i.lo)/(i.hi-i.lo+1/big) end
151
152 function _dist(i, x,y, e)
153 if x=="?" and y=="?" then return 1 end
154 if x=="?" then y = i:norm(y); x = y<.5 and 1 or 0
155 elseif y=="?" then x = i:norm(x); y = x<.5 and 1 or 0
156 else x,y = i:norm(x), i:norm(y) end
157 return math.abs(x - y) end
158
159 function _like(i,x,_, e)
160 return (x < i.mu - 4*i.sd and 0 or x > i.mu + 4*i.sd and 0 or
161 2.7183*(-(x - i.mu)^2 / (z + 2*i.sd^2))/(z + (math.pi*2*i.sd^2)^.5)) end

```

```

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

```

```

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

```

```

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

```

```

383 ---### Main
384 -- Parse help text for flags and defaults, check CLI for updates.
385 -- Maybe print the help (with some pretty colors).
386 -- Run the demos.
387 -- Check for rogue vars.
388 -- Exit, reporting number of failures.
389 -- Help: gsub("m [0-9%]+)%s)+([0-9%]+)%s)+([0-9%]+)%s)+)", these)
390 if the.help then
391   print(help:gsb("u%u+", "%27[31m%127[0m")
392     :gsb("%s)[0-9%]+)%s)+", "%127[33m%27[0m%3*"), "")
393 else
394   local fails = demos()
395   for k, v in pairs(_ENV) do if not b4[k] then print("?", k, type(v)) end end
396   os.exit(fails) end
397 -- function SOME() return {all={}, ok=false, n=0} end
398 -- function some(i, x)
399   if x=="?" then return x end
400   i.n = 1 + i.n
401   if #i.all < the.some then i.ok=false; push(i.all, x)
402   elseif rand() < the.some/i.n then i.ok=false; i.all[rand(#i.all)] = x end end
403   return i.all[math.max(1, math.min(#i.all, (p or .5)*#i.all//1))] end

```