

```

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

```

```

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

```

```

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

```

```

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

```

```

324 -----
325 -- ## DEMOS
326 local go,no={},{}
327
328 -- Convert help string to a table. Check command line for any updates.
329 function these(f1,f2,k,x)
330   for n,flag in ipairs(arg) do if flag==f1 or flag==f2 then
331     x = x=="false" and"true" or x=="true" and"false" or arg[n+1] end end
332   the[k] = string2thing(x) end
333
334 -- Run the demos, resetting settings and random number see before each.
335 -- Return number of failures.
336 function demos( fails,names,defaults,status)
337   fails=0 -- this code will return number of failures
338   names, defaults = {},{}
339   for k,f in pairs(go) do if type(f)=="function" then push(names,k) end end
340   for k,v in pairs(the) do defaults[k]=v end
341   if go[the.go] then names=(the.go) end
342   for __,one in pairs(sort(names)) do -- for all we want to do
343     for k,v in pairs(defaults) do the[k]=v end -- set settings to defaults
344     math.randomseed(the.seed or 10019) -- reset random number seed
345     io.stderr:write("(")
346     status = go[one]() -- run demo
347     if status == true then
348       print("("..Error,one,status)
349       fails = fails + 1 end end -- update fails
350   return fails end -- return total failure count
351
352 -- Simple stuff
353 function go.the() return type(the.bins)=="number" end
354 function go.sort( t) return 0==sort((100,3,4,2,10,0))[1] end
355 function go.slice( t,u)
356   t = {10,20,30,40,50,60,70,80,90,100,110,120,130,140}
357   u = slice(t,3,#t,3)
358   t = slice(t,3,5)
359   return #t==3 and #u==4 end
360
361 function go.num( n,mu,sd)
362   n, mu, sd = NUM(), 10, 1
363   for i=1,10^4 do
364     n:add(mu+sd*math.sqrt(-2*math.log(rand()))*math.cos(2*math.pi*rand())) end
365   return math.abs(n.mu - mu) < 0.05 and math.abs(n.sd - sd) < 0.5 end
366
367 -- Can we read rows off the disk?
368 function go.rows(n,m)
369   m,n=0,0; for row in csv(the.file) do m=m+1; n=n+#row; end; return n/m==8 end
370
371 -- Can we turn a list of names into columns?
372 function go.cols( l)
373   i=COLS{"Name","Age","ShoeSize-"}
374   return i.y[1].w == -1 end
375
376 -- Can we read data, summarized as columns?
377 function go.egs( it)
378   it = EGS.load(the.file); return math.abs(2970 - it.cols.y[1].mu) < 1 end
379
380 -- Can we discretize
381 function go.ranges( it,n,best,rest)
382   it = EGS.load(the.file)
383   print("all",o(rnds(it:mid()))))
384   it.rows = sort(it.rows)
385   for j,row in pairs(sort(it.rows)) do row.klass = 1+j//(#it.rows*.35/6) end
386   n = (#it.rows)*.5
387   best,rest = slice(it.rows,1,n), slice(it.rows, n+1, #it.rows, 3*n)
388   print("best",o(rnds(it:copy(best):mid()))))
389   print("rest",o(rnds(it:copy(rest):mid()))))
390   for __,ranges in pairs(it:ranges(best,rest)) do
391     print("
392     for at,range in pairs(ranges) do
393       print(range) end end
394   --oo(a:mid())
395   --oo(b:mid())
396   return math.abs(2970 - it.cols.y[1].mu) < 1 end

```

```

397 -----
398 -- ## Main
399
400 -- Parse help text for flags and defaults, check CLI for updates.
401 -- Maybe print the help (with some pretty colors).
402 -- Run the demos.
403 -- Check for rogue vars.
404 -- Exit, reporting number of failures.
405 help:gsub("u ([-]|[%s+])[%s]+([-]|[%s+])^u)[^%s](^%s+)",these)
406 if the.help then
407   print(help:gsub("%u%u+", "%27[31m%127[0m]",
408     :gsub("(%s)([-]|[%s+])[%s]+", "%127[33m%27[0m%3*]", ""))
409 else
410   local fails = demos()
411   for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
412   os.exit(fails) end
413 -- function SOME() return {all={}, ok=false, n=0} end
414 -- function some(i,x)
415 --   if x=="?" then return x end
416 --   i.n = 1 + i.n
417 --   if #i.all < the.some then i.ok=false; push(i.all, x)
418 --   elseif rand() < the.some/i.n then i.ok=false; i.all[#i.all]=x end end
419
420 -- function per(i,p)
421 --   i.all = i.ok and i.all or sort(i.all); i.ok=true
422 --   return i.all[math.max(1, math.min(#i.all, (p or .5)*#i.all//1))] end

```