

```

1  #!/usr/bin/env lua
2  --
3  --
4  --
5  --
6  --
7  --
8  --
9
10 local your, our={}, {b4={}, help=[[
11 duo.lua [OPTIONS]
12 (c)2022 Tim Menzies, MIT license (2 clause)
13 Data miners using/used by optimizers.
14 Understand N items after log(N) probes, or less.
15
16 -file    ../..data/au093.csv
17 -ample   512
18 -far     .9
19 -best    .5
20 -help    false
21 -dull    .5
22 -rest    3
23 -seed    10019
24 -Small   .35
25 -rnd     %.2f
26 -task    -
27 -p       2]]}
28
29 for k, _ in pairs(_ENV) do our.b4[k] = k end
30 local any, asserts, cells, copy, first, fmt, go, id, main, many, map
31 local merge, new, o, push, rand, randi, ranges, rnd, rogues, rows, same
32 local second, seconds, settings, slots, sort, super, thing, things, xpect
33 local COLS, EG, EGS, NUM, RANGE, SAMPLE, SYM
34 local class= function(t, new)
35     function new(_, ...) return t.new(...) end
36     t.__index=t
37     return setmetatable(t, {__call=new}) end
38
39 -- Copyright (c) 2022, Tim Menzies
40 --
41 -- Redistribution and use in source and binary forms, with or without
42 -- modification, are permitted provided that the following conditions are met.
43 -- (1) Redistributions of source code must retain the above copyright notice,
44 -- this list of conditions and the following disclaimer. (2) Redistributions
45 -- in binary form must reproduce the above copyright notice, this list of
46 -- conditions and the following disclaimer in the documentation and/or other
47 -- materials provided with the distribution.
48 --
49 -- THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS
50 -- IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO,
51 -- THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR
52 -- PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR
53 -- CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL,
54 -- EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO,
55 -- PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR
56 -- PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF
57 -- LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING
58 -- NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS
59 -- SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE
60

```

```

60 --
61 --
62 --
63 --
64 --
65 --
66 COLS=class{}
67 function COLS.new(t, i, where, now)
68     i = new({all={}, x={}, y={}}, COLS)
69     for at, s in pairs(t) do
70         now = push(i.all, (s:find"^[A-Z]" and NUM or SYM) (at, s))
71         if not s:find"." then
72             push((s:find"-" or s:find"+") and i.y or i.x, now) end end
73     return i end
74
75 function COLS.__tostring(i, txt)
76     function txt(c) return c.txt end
77     return fmt("COLS{all %s\n\t: x %s\n\t: y %s", o(i.all, txt), o(i.x, txt), o(i.y, txt)) end
78
79 function COLS.add(i, t, add)
80     function add(col, x) x=t[col.at]; col:add(x); return x end
81     return map(i.all, add) end
82
83 -----
84 EG=class{}
85 function EG.new(t) return new({has=t, id=id()}, EG) end
86
87 function EG.__tostring(i) return fmt("EG%s%s", i.id, o(i.has), #i.has) end
88
89 function EG.better(i, j, cols)
90     local s1, s2, e, n, a, b = 0, 0, 10, #cols
91     for _, col in pairs(cols) do
92         a = col:norm(i.has[col.at])
93         b = col:norm(j.has[col.at])
94         s1 = s1 - e^(col.w * (a-b)/n)
95         s2 = s2 - e^(col.w * (b-a)/n) end
96     return s1/n < s2/n end
97
98 function EG.col(i, cols)
99     return map(cols, function(col) return i.has[col.at] end) end
100
101 function EG.dist(i, j, egs, a, b, d, n)
102     d, n = 0, #egs.cols.x + 1E-31
103     for _, col in pairs(egs.cols.x) do
104         a, b = i.has[col.at], j.has[col.at]
105         d = d + col:dist(a, b) ^ your.p end
106     return (d/n) ^ (1/your.p) end
107
108 -----
109 EGS=class{}
110 function EGS.new() return new({rows={}, cols=nil}, EGS) end
111
112 function EGS.__tostring(i) return fmt("EGS{#rows %s: cols %s", #i.rows, i.cols) end
113
114 function EGS.add(i, row)
115     row = row.has and row.has or row
116     if i.cols then push(i.rows, EG(i.cols:add(row))) else i.cols=COLS(row) end end
117
118 function EGS.clone(i, inits, j)
119     j = EGS()
120     j:add(map(i.cols.all, function(col) return col.txt end))
121     for _, x in pairs(inits or {}) do j:add(x) end
122     return j end
123
124 function EGS.far(i, eg1, rows, fun, tmp)
125     fun = function(eg2) return {eg2, eg1:dist(eg2, i)} end
126     tmp = sort(map(rows, fun), seconds)
127     return table.unpack(tmp[#tmp*your.far//1]) end
128
129 function EGS.file(i, file) for row in rows(file) do i:add(row) end; return i end
130
131 function EGS.mid(i, cols, mid)
132     function mid(col) return col:mid() end
133     return map(cols or i.cols.y, mid) end
134
135 function EGS.halve(i, rows)
136     local c, l, r, ls, rs, cosine, some
137     function cosine(row, a, b)
138         a, b = row:dist(l, i), row:dist(r, i); return {(a^2+c^2-b^2)/(2*c), row} end
139

```

```

137 rows = rows or i.rows
138 some = #rows > your.ample and many(rows, your.ample) or rows
139 l = i:far(any(rows), some)
140 r,c = i:far(l, some)
141 ls,rs = i:clone(), i:clone()
142 for n,pair in pairs(sort(map(rows,cosine), firsts)) do
143   (n <= #rows//2 and ls or rs):add(pair[2]) end
144 return ls,rs,l,r,c end
145
146 function EGS.ranges(i,j, all,there, ranges)
147   all = {}
148   for n,here in pairs(i.cols.x) do
149     there = j.cols.x[n]
150     ranges = here:ranges(there)
151     if #ranges>1 then push(all, {xpect(ranges),ranges}) end end
152 return map(sort(all,firsts),second) end
153
154 function EGS.xcluster(i,top,lvl)
155   local split, left, right,kid1, kid2
156   top, lvl = top or i, lvl or 0
157   ls,rs = (top or i):halve(i.rows)
158   if #i.rows >= 2*(#top.rows)^your.small then
159     split, kid1, kid2 = i:splitter(top), i:clone(), i:clone()
160     for _,row in pairs(i.rows) do
161       (split:selects(row) and kid1 or kid2):add(row) end
162     if #kid1.rows ~= #i.rows then left = kid1:xcluster(top,lvl+1) end
163     if #kid2.rows ~= #i.rows then right = kid2:xcluster(top,lvl+1) end
164   end
165   return {here=i, split=split, left=left, right=right} end
166
167 -----
168 NUM=class{}
169 function NUM.new(at,s, big)
170   big = math.huge
171   return new({lo=big, hi=-big, at=at or 0, txt=s or "",
172     n=0, mu=0, m2=0, sd=0, _all=SAMPLE(),
173     w=(s or ""):find("-" and -1 or 1),NUM) end
174
175 function NUM.__tostring(i)
176   return fmt("NUM{at %s:txt %s:n %s:lo %s:hi %s:mu %s:sd %s}",
177     i.at, i.txt, i.n, i.lo, i.hi, rnd(i.mu), rnd(i.div())) end
178
179 function NUM.add(i,x, d,pos)
180   if x~="?" then
181     i.n = i.n+1
182     d = x - i.mu
183     i.mu = i.mu + d/i.n
184     i.m2 = i.m2 + d*(x-i.mu)
185     i.lo = math.min(x,i.lo); i.hi = math.max(x,i.hi)
186     i._all:add(x) end
187   return x end
188
189 function NUM.dist(i,a,b)
190   if a=="?" and b=="?" then a,b =1,0
191   elseif a=="?" then b = i:norm(b); a=b>.5 and 0 or 1
192   elseif b=="?" then a = i:norm(a); b=a>.5 and 0 or 1
193   else a,b = i:norm(a), i:norm(b) end
194   return math.abs(a-b) end
195
196 function NUM.div(i) return i.n <2 and 0 or (i.m2/(i.n-1))^0.5 end
197
198 function NUM.merge(i,j, k)
199   k= NUM(i.at, i.txt)
200   for _,x in pairs(i._all,it) do k:add(x) end
201   for _,x in pairs(j._all,it) do k:add(x) end
202   return k end
203
204 function NUM.mid(i) return i.mu end
205
206 function NUM.norm(i,x) return i.hi-i.lo < 1E-9 and 0 or (x-i.lo)/(i.hi-i.lo) end
207
208 function NUM.ranges(i,j,ykind, tmp,xys)
209   xys={}
210   for _,x in pairs(i._all,it) do push(xys,{x=x,y="best"}) end
211   for _,x in pairs(j._all,it) do push(xys,{x=x,y="rest"}) end
212   return merge( ranges(xys,i, ykind or SYM,
213     (#xys)^your.dull,
214     xpect{i,j}*your.Small)) end

```

```

214 -----
215 RANGE=class{}
216 function RANGE.new(col,lo,hi,ys)
217   return new({n=0, col=col, lo=lo, hi=hi or lo, ys=ys or SYM{}},RANGE) end
218
219 function RANGE.__lt(i,j) return i:div() < j:div() end
220
221 function RANGE.__tostring(i)
222   if i.lo == i.hi then return fmt("%s==%s", i.col.txt, i.lo) end
223   if i.lo == -math.huge then return fmt("%s<%s", i.col.txt, i.hi) end
224   if i.hi == math.huge then return fmt("%s>=%s", i.col.txt, i.lo) end
225   return fmt("%s<=%s<%s", i.lo, i.col.txt, i.hi) end
226
227 function RANGE.add(i,x,y,inc)
228   inc = inc or 1
229   i.n = i.n + inc
230   i.hi = math.max(x,i.hi)
231   i.ys:add(y, inc) end
232
233 function RANGE.div(i) return i.ys:div() end
234
235 function RANGE.selects(i,row, x)
236   x=row.has[col.at]; return x=="?" or i.lo<=x and x<i.hi end
237
238 -----
239 SAMPLE=class{}
240 function SAMPLE.new() return new({n=0,it={},ok=false,max=your.ample},SAMPLE) end
241
242 function SAMPLE.add(i,x, pos)
243   i.n = i.n + 1
244   if #i.it < i.max then pos= #i.it + 1
245   elseif rand() < #i.it/i.n then pos= #i.it * rand() end
246   if pos then i.ok = false; i.it[pos//1]= x end end
247
248 function SAMPLE.all(i) if not i.ok then i.ok=true;sort(i.it)end; return i.it end
249
250 -----
251 SYM=class{}
252 function SYM.new(at,s)
253   return new({at=at or 0,txt=s or "",has={},n=0,most=0,mode=nil},SYM) end
254
255 function SYM.__tostring(i)
256   return fmt("SYM{at %s:txt %s:mode %s:has %s}",
257     i.at, i.txt, i.mode, o(i.has)) end
258
259 function SYM.add(i,x, inc)
260   if x ~= "?" then
261     inc = inc or 1
262     i.n = i.n+inc
263     i.has[x] = inc + (i.has[x] or 0)
264     if i.has[x] > i.most then i.most, i.mode = i.has[x], x end end
265   return x end
266
267 function SYM.dist(i,a,b) return a=="?" and b=="?" and 1 or a==b and 0 or 1 end
268
269 function SYM.div(i, e)
270   e=0;for _,v in pairs(i.has) do e=e - v/i.n*math.log(v/i.n,2) end; return e end
271
272 function SYM.merge(i,j, k)
273   k= SYM(i.at, i.txt)
274   for x,count in pairs(i.has) do k:add(x,count) end
275   for x,count in pairs(j.has) do k:add(x,count) end
276   return k end
277
278 function SYM.mid(i) return i.mode end
279
280 function SYM.ranges(i,j, t)
281   t = {}
282   for _,pair in pairs({i.has,"bests"}, {j.has,"rests"}) do
283     for x,inc in pairs(pair[1]) do
284       t[x] = t[x] or RANGE(i,x)
285       t[x]:add(x, pair[2], inc) end end
286   return map(t) end

```

```

285 --
286 --
287 --
288 --
289 --
290
291 fmt = string.format
292 new = setmetatable
293 same = function(x,...) return x end
294
295 function any(t) return t[randi(1,#t)] end
296
297 function asserts(test,msg)
298   msg=msg or ""
299   if test then return print("PASS: "..msg) end
300   our.failures = our.failures + 1
301   print("FAIL: "..msg)
302   if your.Debug then assert(test,msg) end end
303
304 function copy(t, u)
305   if type(t)~="table" then return t end
306   u={};for k,v in pairs(t) do u[k]=copy(v) end;return new(u,getmetatable(t)) end
307
308 function first(a,b) return a[1] end
309
310 function firsts(a,b) return a[1] < b[1] end
311
312 function id() our.id = 1+(our.id or 0); return our.id end
313
314 function many(t,n, u) u={};for j=1,n do push(u,any(t)) end; return u end
315
316 function map(t,f, u)
317   u={};for _,v in pairs(t) do u[1+#u]=(f or same)(v) end; return u end
318
319 function o(t,f, u,key)
320   key= function(k)
321     if t[k] then return fmt(":%s%s", k, rnd((f or same)(t[k]))) end end
322   u = #t>0 and map(map(t,f),rnd) or map(slots(t),key)
323   return "{..table.concat(u, " ")..}" end
324
325 function rand(lo,hi)
326   your.seed = (16807 * your.seed) % 2147483647
327   return (lo or 0) + ((hi or 1) - (lo or 0)) * your.seed / 2147483647 end
328
329 function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
330
331 function push(t,x) table.insert(t,x); return x end
332
333 function rnd(x)
334   return fmt(type(x)=="number" and x~=x//1 and your.rnd or"%s",x) end
335
336 function rows(file, x)
337   file = io.input(file)
338   return function()
339     x=io.read(); if x then return things(x) else io.close(file) end end end
340
341 function main( defaults,tasks)
342   tasks = your.task=="all" and slots(go) or {your.task}
343   defaults=copy(your)
344   our.failures=0
345   for _,x in pairs(tasks) do
346     if type(our.go[x]) == "function" then our.go[x]() end
347     your = copy(defaults) end
348   rogues()
349   return our.failures end
350
351 function merge(b4, j,tmp,merged,one,two)
352   j, tmp = 0, {}
353   while j < #b4 do
354     j = j + 1
355     one, two = b4[j], b4[j+1]
356     if two then
357       merged = one.ys:merge(two.ys)
358       local after=merged:div()
359       local b4=xpect{one.ys,two.ys}
360       --print(o{before=b4, one=one.ys.n, two=two.ys.n,after=after,frac=math.abs(
after-b4)/b4})
361
362       if after+b4> 0.01 and after<= b4 or math.abs(after-b4)/b4 < .1 then
363         j = j+1
364         one = RANGE(one.col, one.lo, two.hi, merged) end end
365       push(tmp,one) end
366       return #tmp==#b4 and b4 or merge(tmp) end
367
368 function ranges(xys,col,ykind, small, dull, one,out)
369   out = {}
370   xys = sort(xys, function(a,b) return a.x < b.x end)
371   one = push(out, RANGE(col, xys[1].x, xys[1].x, ykind()))
372   for j,xy in pairs(xys) do
373     if j < #xys - small and -- enough items remaining after split
374       xy.x ~= xys[j+1].x and -- next item is different (so can split here)
375       one.n > small and -- one has enough items
376       one.hi - one.lo > dull -- one is not trivially small
377     then one = push(out, RANGE(col, one.hi, xy.x, ykind())) end
378     one:add(xy.x, xy.y) end
379   out[1].lo = -math.huge
380   out[#out].hi = math.huge
381   return out end
382
383 function rogues()
384   for k,v in pairs(_ENV) do
385     if not our.b4[k] then print("??",k,type(v)) end end end
386
387 function second(t) return t[2] end
388
389 function seconds(a,b) return a[2] < b[2] end
390
391 function settings(help, t)
392   t={}
393   help:gsub("\n[-][^%s+][^\\n]*%s([%s]+)", function(slot, x)
394     for n,flag in ipairs(arg) do
395       if flag:sub(1,1)=="-" and slot:match("^"..flag:sub(2)..".*")
396       then x=x=="false" and "true" or x=="true" and "false" or arg[n+1] end end
397     t[slot] = thing(x) end
398   if t.help then print(t.help) end
399   return t end
400
401 function slots(t,u) u={};for x,_ in pairs(t) do u[1+#u]=x end;return sort(u) end
402
403 function sort(t,f) table.sort(t,f); return t end
404
405 function thing(x)
406   x = x:match("^%s*(-)%s*$")
407   if x=="true" then return true elseif x=="false" then return false end
408   return tonumber(x) or x end
409
410 function things(x,sep, t)
411   t={};for y in x:gmatch(sep or"([^\,]+)") do t[1+#t]=thing(y) end; return t end
412
413 function xpect(t)
414   local m,d = 0,0
415   for _,z in pairs(t) do m=m+z.n; d=d+z.n*z:div() end; return d/m end

```

```

416  ---
417  ---
418  ---
419  ---
420
421  our.go, our.no = {},{}; go=our.go
422  function go.settings() print("your",o(your)) end
423
424  function go.sample() print(EGS():file(your.file)) end
425
426  function go.clone( a,b)
427    a = EGS():file(your.file)
428    b = a:clone(a.rows)
429    asserts(#a.rows == #b.rows,"cloning rows")
430    asserts(tostring(a.cols.all[1])==tostring(b.cols.all[1]),"cloning cols")
431  end
432
433  function go.dist( t,a,egl,eg2)
434    a = EGS():file(your.file)
435    egl = any(a.rows)
436    print(o(egl:col(a.cols.x)))
437    t={}
438    for j=1,20 do
439      eg2 = any(a.rows)
440      push(t, {egl:dist(eg2,a),eg2}) end
441    for _,pair in pairs(sort(t,firsts)) do
442      print(o(pair[2]:col(a.cols.x)),rnd(pair[1])) end end
443
444  function go.halve( a,b)
445    a,b = EGS():file(your.file):halve()
446    print(o(a:mid()))
447    print(o(b:mid())) end
448
449  function go.ranges( a,b,x,col2)
450    a,b = EGS():file(your.file):halve()
451    for n,col1 in pairs(a.cols.x) do
452      col2 = b.cols.x[n]
453      col1:ranges(col2) end end
454  -- x = a:delta(b)
455  -- print(x,type(x))
456  -- print(">>", x.lo, x.hi)
457  -- end
458
459  your = settings(our.help)
460  os.exit( main() )

```