

```

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/auto93.csv
17 -ample   512
18 -far     .9
19 -best    .5
20 -help    false
21 -chill   .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, firsts, 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
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\ntx %s\nty %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 %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, eggs, 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 %s %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, egl, rows, fun, tmp)
125   fun = function(eg2) return {eg2, egl: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   rows = rows or i.rows
140   some = #rows > your.ample and many(rows, your.ample) or rows
141   l = i:far(any(rows), some)
142   r, c = i:far(l, some)
143   ls, rs = i:clone(), i:clone()
144   for n, pair in pairs(sort(map(rows, cosine), firsts)) do
145     (n <= #rows/2 and ls or rs):add(pair[2]) end
146   return ls, rs, l, r, c end
147
148 -- XXX ranges2 suspicious. d=0 and morerangesis 0
149 function EGS.ranges(i, j, all, there, ranges)
150   all = {}
151   for n, here in pairs(i.cols.x) do
152     there = j.cols.x[n]
153     ranges = here:ranges(there)
154     if #ranges > 1 then push(all, {xpect(ranges, here.txt .. "ranges"), ranges}) end
155   end
156   --for k, v in pairs(sort(all, firsts)) do
157   --  print(v[1], #v[2], v[2][1].col.txt) end
158   return map(sort(all, firsts), second) end
159
160 function EGS.xcluster(i, top, lvl)
161   local split, left, right, kid1, kid2
162   top, lvl = top or i, lvl or 0
163   ls, rs = (top or i):halve(i.rows)
164   if #i.rows >= 2*(#top.rows)^your.small then
165     split, kid1, kid2 = i:splitter(top, i:clone(), i:clone())
166     for _, row in pairs(i.rows) do
167       (split:selects(row) and kid1 or kid2):add(row) end
168     if #kid1.rows ~= #i.rows then left = kid1:xcluster(top, lvl+1) end
169     if #kid2.rows ~= #i.rows then right = kid2:xcluster(top, lvl+1) end
170     end
171     return (here=i, split=split, left=left, right=right) end
172
173 -----
174 NUM=class{}
175 function NUM.new(at, s, big)
176   big = math.huge
177   return new({lo=big, hi=-big, at=at or 0, txt=s or "",
178     n=0, mu=0, m2=0, sd=0, all=SAMPLE(),
179     w=(s or ""):find"-" and -1 or 1}, NUM) end
180
181 function NUM.__tostring(i)
182   return fmt("NUM{at %s txt %s n %s do %s hi %s mu %s sd %s}",
183     i.at, i.txt, i.n, i.lo, i.hi, rnd(i.mu), rnd(i:div())) end
184
185 function NUM.add(i, x, d, pos)
186   if x=="?" then
187     i.n = i.n+1
188     d = x - i.mu
189     i.mu = i.mu + d/i.n
190     i.m2 = i.m2 + d*(x-i.mu)
191     i.lo = math.min(x, i.lo); i.hi = math.max(x, i.hi)
192     i.all:add(x) end
193   return x end
194
195 function NUM.dist(i, a, b)
196   if a=="?" and b=="?" then a, b = 1, 0
197   elseif a=="?" then b = i:norm(b); a=b>.5 and 0 or 1
198

```

```

195 elseif b=="?" then a = i:norm(a); b=a>.5 and 0 or 1
196 else a,b = i:norm(a), i:norm(b) end
197 return math.abs(a-b) end
198
199 function NUM.div(i) return i.n < 2 and 0 or (i.m2/(i.n-1))^0.5 end
200
201 function NUM.merge(i,j, k)
202   k= NUM(i.at, i.txt)
203   for _,x in pairs(i._all.it) do k:add(x) end
204   for _,x in pairs(j._all.it) do k:add(x) end
205   return k end
206
207
208 function NUM.mid(i) return i.mu end
209
210
211 function NUM.norm(i,x) return i.hi-i.lo < 1E-9 and 0 or (x-i.lo)/(i.hi-i.lo) end
212
213 function NUM.ranges(i,j,ykind, tmp,xys)
214   xys={}
215   for _,x in pairs(i._all.it) do push(xys,{x=x,y="best"}) end
216   for _,x in pairs(j._all.it) do push(xys,{x=x,y="rest"}) end
217   return merge( ranges(xys,i, ykind or SYM,
218     (#xys)*your.dull,
219     xpect{i,j}*your.Small)) end
220
221 -----
222 RANGE=class{}
223 function RANGE.new(col,lo,hi,ys)
224   return new({n=0, col=col, lo=lo, hi=hi or lo, ys=ys or SYM()},RANGE) end
225
226 function RANGE.__lt(i,j) return i:div() < j:div() end
227
228 function RANGE.__tostring(i)
229   if i.lo == i.hi then return fmt("%s==%s", i.col.txt, i.lo) end
230   if i.lo == -math.huge then return fmt("%s<%s", i.col.txt, i.hi) end
231   if i.hi == math.huge then return fmt("%s>=%s", i.col.txt, i.lo) end
232   return fmt("%s<=%s<%s", i.lo, i.col.txt, i.hi) end
233
234 function RANGE.add(i,x,y,inc)
235   inc = inc or 1
236   i.n = i.n + inc
237   i.hi = math.max(x,i.hi)
238   i.ys:add(y, inc) end
239
240 function RANGE.div(i) return i.ys:div() end
241
242 function RANGE.selects(i,row, x)
243   x=row.has[col.at]; return x=="?" or i.lo<=x and x<i.hi end
244
245 -----
246 SAMPLE=class{}
247 function SAMPLE.new() return new({n=0,it={},ok=false,max=your.ample},SAMPLE) end
248
249 function SAMPLE.add(i,x, pos)
250   i.n = i.n + 1
251   if #i.it < i.max then pos= #i.it + 1
252   elseif rand() < #i.it/i.n then pos= #i.it * rand() end
253   if pos then i.ok = false; i.it[pos//1]= x end end
254
255 function SAMPLE.all(i) if not i.ok then i.ok=true;sort(i.it)end; return i.it end
256
257 -----
258 SYM=class{}
259 function SYM.new(at,s)
260   return new({at=at or 0,txt=s or "",has={},n=0,most=0,mode=nil},SYM) end
261
262 function SYM.__tostring(i)
263   return fmt("SYM{:at %s.txt %s :mode %s :has %s}",
264     i.at, i.txt, i.mode, o(i.has)) end
265
266 function SYM.add(i,x, inc)
267   if x ~= "?" then
268     inc = inc or 1
269     i.n = i.n+inc
270     i.has[x] = inc + (i.has[x] or 0)
271     if i.has[x] > i.most then i.most, i.mode = i.has[x], x end end
272   return x end
273
274 function SYM.dist(i,a,b) return a=="?" and b=="?" and 1 or a==b and 0 or 1 end
275
276 function SYM.div(i, e)
277   e=0;for _,v in pairs(i.has) do e=e - v/i.n*math.log(v/i.n,2) end; return e end
278
279 function SYM.merge(i,j, k)
280   k= SYM(i.at, i.txt)
281   for x,count in pairs(i.has) do k:add(x,count) end
282   for x,count in pairs(j.has) do k:add(x,count) end
283   return k end
284
285 function SYM.mid(i) return i.mode end
286
287 function SYM.ranges(i,j, t)
288   t = {}
289   for _,pair in pairs({i.has,"bests"}, {j.has,"rests"}) do
290     for x,inc in pairs(pair[1]) do
291       t[x] = t[x] or RANGE(i,x)
292       print("inc",i.txt,inc)
293       t[x]:add(x, pair[2], inc) end end
294   return map(t) end
295
296

```

```

290 --
291 --
292 --
293 --
294 --
295 --
296 fmt = string.format
297 new = setmetatable
298 same = function(x,...) return x end
299
300 function any(t) return t[randi(1,#t)] end
301
302 function asserts(test,msg)
303   msg=msg or ""
304   if test then return print("PASS:".msg) end
305   our.failures = our.failures + 1
306   print("FAIL:".msg)
307   if your.Debug then assert(test,msg) end end
308
309 function copy(t, u)
310   if type(t)~="table" then return t end
311   u={};for k,v in pairs(t) do u[k]=copy(v) end;return new(u,getmetatable(t)) end
312
313 function first(a,b) return a[1] end
314
315 function firsts(a,b) return a[1] < b[1] end
316
317 function id() our.id = 1+(our.id or 0); return our.id end
318
319 function many(t,n, u) u={};for j=1,n do push(u,any(t)) end; return u end
320
321 function map(t,f, u)
322   u={};for _,v in pairs(t) do u[1+#u]=(f or same)(v) end; return u end
323
324 function main( defaults,tasks)
325   defaults=copy(your)
326   our.failures=0
327   tasks = your.task=="all" and slots(go) or {your.task}
328   for _,x in pairs(tasks) do
329     if type(our.go[x]) == "function" then our.go[x]() end
330     your = copy(defaults) end
331   rogues()
332   return our.failures end
333
334 function merge(b4, j,tmp,merged,one,two)
335   j, tmp = 0, {}
336   while j < #b4 do
337     j = j + 1
338     one, two = b4[j], b4[j+1]
339     if two then
340       merged = one.ys:merge(two.ys)
341       local after=merged:div()
342       local b4=xpect(one.ys,two.ys)
343       if after+b4< 0.01 or after<= b4 or math.abs(after-b4)/b4 < .1 then
344         j = j+1
345         one = RANGE(one.col, one.lo, two.hi, merged) end end
346     push(tmp,one) end
347   return #tmp==#b4 and b4 or merge(tmp) end
348
349 function o(t,f, u,key)
350   key= function(k)
351     if t[k] then return fmt(":%s %s", k, rnd((f or same)(t[k]))) end end
352   u = #t>0 and map(map(t,f),rnd) or map(slots(t),key)
353   return "{..table.concat(u, " ")}.." end
354
355 function push(t,x) table.insert(t,x); return x end
356
357 function rand(lo,hi)
358   your.seed = (16807 * your.seed) % 2147483647
359   return (lo or 0) + ((hi or 1) - (lo or 0)) * your.seed / 2147483647 end
360
361 function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
362
363 function ranges(xys,col,ykind, small, dull, one,out)
364   out = {}
365   xys = sort(xys, function(a,b) return a.x < b.x end)
366   one = push(out, RANGE(col, xys[1].x, xys[1].x, ykind()))
367   for j,xy in pairs(xys) do
368     if j < #xys - small and -- enough items remaining after split
369     xy.x ~- xys[j+1].x and -- next item is different (so can split here)
370     one.n > small and -- one has enough items
371     one.hi - one.lo > dull and -- one is not trivially small
372     then one = push(out, RANGE(col, one.hi, xy.x, ykind())) end
373     one:add(xy.x, xy.y) end
374   out[1].lo = -math.huge
375   out[#out].hi = math.huge
376   return out end
377
378 function rnd(x)
379   return fmt(type(x)=="number" and x~x//1 and your.rnd or"%s",x) end
380
381 function rogues()
382   for k,v in pairs(_ENV) do
383     if not our.b4[k] then print("?",k,type(v)) end end end
384
385 function rows(file, x)
386   file = io.input(file)
387   return function()
388     x=io.read(); if x then return things(x) else io.close(file) end end end
389
390 function second(t) return t[2] end
391
392 function seconds(a,b) return a[2] < b[2] end
393
394 function settings(help, t)
395   t={}
396   help:gsub("\n [-](^%s+)[^%n]*%s([%^s]+)", function(slot, x)
397     for n,flag in ipairs(arg) do
398       if flag:sub(1,1)=="-" and slot:match("^[^].flag:sub(2).."*)
399       then x=x.."false" and true" or x=="true" and "false" or arg[n+1] end end
400       t[slot] = thing(x) end
401     if t.help then print(t.help) end
402     return t end
403
404 function slots(t,u) u={};for x,_ in pairs(t) do u[1+#u]=x end;return sort(u) end
405
406 function sort(t,f) table.sort(t,f); return t end
407
408 function thing(x)
409   x = x:match"%s*(-)%s*$"
410   if x=="true" then return true elseif x=="false" then return false end
411   return tonumber(x) or x end
412
413 function things(x,sep, t)
414   t={};for y in x:gmatch(sep or"([^\s]+)") do t[1+#t]=thing(y) end; return t end
415
416 function xpect(t,s)
417   local m,d = 0,0
418   for _,z in pairs(t) do m=m+z.n; d=d+z.n*z:div() end; print(o{d=d,m=m},s or "");
419   return d/m end
420
421

```

```

419 --
420 --
421 --
422 --
423 --
424
425 our.go, our.no = {},{}; go=our.go
426 function go.settings() print("your",o(your)) end
427
428 function go.sample() print(EGS():file(your.file)) end
429
430 function go.clone( a,b)
431 a= EGS():file(your.file)
432 b= a:clone(a.rows)
433 asserts(#a.rows == #b.rows,"cloning rows")
434 asserts(tostring(a.cols.all[1])==tostring(b.cols.all[1]),"cloning cols")
435 end
436
437 function go.dist( t,a,eg1,eg2)
438 a= EGS():file(your.file)
439 eg1 = any(a.rows)
440 print(o(eg1:col(a.cols.x)))
441 t={}
442 for j=1,20 do
443 eg2 = any(a.rows)
444 push(t, {eg1:dist(eg2,a),eg2}) end
445 for _,pair in pairs(sort(t,firsts)) do
446 print(o(pair[2]:col(a.cols.x)),rnd(pair[1])) end end
447
448 function go.halve( a,b)
449 a,b = EGS():file(your.file):halve()
450 print(o(a:mid()))
451 print(o(b:mid())) end
452
453 function go.ranges( a,b,x,col2)
454 a,b = EGS():file(your.file):halve()
455 for n,coll in pairs(a.cols.x) do
456 col2 = b.cols.x[n]
457 print("")
458 for _, range in pairs(coll:ranges(col2)) do
459 print(coll.txt, range.lo, range.hi) end end end
460
461 function go.ranges2( a,b,x,col2)
462 a,b = EGS():file(your.file):halve()
463 a:ranges(b) end
464 -- x = a:delta(b)
465 -- print(x,type(x))
466 -- print(">>", x.lo, x.hi)
467 -- end
468
469 your = settings(our.help)
470 os.exit( main() )

```