```lua
1   #!/usr/bin/env lua
2   --       __
3   --      /\ \
4   --      \ \ \         ___         ___
5   --      /'_` \       /\_\/\_\     /_/`\
6   --    /\ \L\ \     \ \ \\_\ \     /\  \L\ \
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    -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,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 MERCHNTABILITY 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  EG=class{}
84  function EG.new(t) return new({has=t, id=id()},EG) end
85
86  function EG.__tostring(i) return fmt("EG%s%s %s", i.id,o(i.has),#i.has) end
87
88  function EG.better(i,j,cols)
89    local s1,s2,e,n,a,b = 0,0,10,#cols
90    for _,col in pairs(cols) do
91      a  = col:norm(i.has[col.at])
92      b  = col:norm(j.has[col.at])
93      s1 = s1 - e^(col.w * (a-b)/n)
94      s2 = s2 - e^(col.w * (b-a)/n) end
95    return s1/n < s2/n end
96
97  function EG.col(i,cols)
98    return map(cols, function(col) return i.has[col.at] end) end
99
100 function EG.dist(i,j,egs,     a,b,d,n)
101   d,n = 0, #egs.cols.x + 1E-31
102   for _,col in pairs(egs.cols.x) do
103     a,b = i.has[col.at], j.has[col.at]
104     d   = d + col:dist(a,b) ^ your.p end
105   return (d/n) ^ (1/your.p) end
106 -- ----------------------------------------------------------------------------
107 EGS=class{}
108 function EGS.new() return new({rows={}, cols=nil}, EGS) end
109
110 function EGS.__tostring(i) return fmt("EGS{#rows %s:cols %s", #i.rows,i.cols) end
111
112 function EGS.add(i,row)
113   row = row.has and row.has or row
114   if i.cols then push(i.rows,EG(i.cols:add(row))) else i.cols=COLS(row) end end
115
116 function EGS.clone(i,inits,     j)
117   j = EGS()
118   j:add(map(i.cols.all, function(col) return col.txt end))
119   for _,x in pairs(inits or {}) do  j:add(x) end
120   return j end
121
122 function EGS.far(i,eg1,rows,     fun,tmp)
123   fun = function(eg2) return {eg2, eg1:dist(eg2,i)} end
124   tmp = sort(map(rows, fun), seconds)
125   return table.unpack(tmp[#tmp*your.far//1] ) end
126
127 function EGS.file(i,file) for row in rows(file) do i:add(row) end; return i end
128
129 function EGS.mid(i,cols,      mid)
130   function mid(col)  return col:mid() end
131   return map(cols or i.cols.y, mid) end
132
133 function EGS.halve(i,rows)
134   local c,l,r,ls,rs,cosine,some
135   function cosine(row,      a,b)
136     a,b = row:dist(l,i), row:dist(r,i); return {(a^2+c^2-b^2)/(2*c),row} end
```

page 2

```lua
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 -- XXX ranges2 suspicious. d=0 and morerangesis 0
147 function EGS.ranges(i,j,     all,there, ranges)
148    all = {}
149    for n,here in pairs(i.cols.x) do
150      there = j.cols.x[n]
151      ranges = here:ranges(there)
152      if #ranges> 1 then push(all, {xpect(ranges,here.txt .. "ranges"),ranges}) end
153      end
154    --for k,v  in pairs(sort(all,firsts)) do
155      -- print(v[1], #v[2], v[2][1].col.txt) end
156    return map(sort(all,firsts),second) end
157
158 function EGS.xcluster(i,top,lvl)
159    local split, left, right,kid1, kid2
160    top, lvl = top or i, lvl or 0
161    ls,rs = (top or i):halve(i.rows)
162    if #i.rows >= 2*(#top.rows)^your.small then
163      split, kid1, kid2 = i:splitter(top), i:clone(), i:clone()
164      for _,row in pairs(i.rows) do
165        (split:selects(row) and kid1 or kid2):add(row) end
166      if #kid1.rows ~= #i.rows then left  = kid1:xcluster(top,lvl+1) end
167      if #kid2.rows ~= #i.rows then right = kid2:xcluster(top,lvl+1) end
168    end
169    return {here=i, split=split, left=left, right=right} end
170 -- -------------------------------------------------------------------------
171 NUM=class{}
172 function NUM.new(at,s, big)
173    big = math.huge
174    return new({lo=big, hi=-big, at=at or 0, txt=s or "",
175             n=0, mu=0, m2=0, sd=0,_all=SAMPLE(),
176             w=(s or ""):find"-" and -1 or 1},NUM) end
177
178 function NUM.__tostring(i)
179    return fmt("NUM{:at %s :txt %s :n %s :lo %s :hi %s :mu %s :sd %s}",
180             i.at, i.txt,  i.n, i.lo, i.hi, rnd(i.mu), rnd(i:div())) end
181
182 function NUM.add(i,x,      d,pos)
183    if x~="?" then
184      i.n  = i.n+1
185      d    = x - i.mu
186      i.mu = i.mu + d/i.n
187      i.m2 = i.m2 + d*(x-i.mu)
188      i.lo = math.min(x,i.lo); i.hi = math.max(x,i.hi)
189      i._all:add(x) end
190    return x end
191
192 function NUM.dist(i,a,b)
193    if     a=="?" and b=="?" then a,b =1,0
194    elseif a=="?"           then b    = i:norm(b); a=b>.5 and 0 or 1
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 function NUM.mid(i) return i.mu end
208
209 function NUM.norm(i,x) return i.hi-i.lo < 1E-9 and 0 or (x-i.lo)/(i.hi-i.lo) end
210
211 function NUM.ranges(i,j,ykind,        tmp,xys)
212    xys={}
213    for _,x in pairs(i._all.it) do push(xys,{x=x,y="best"}) end
214    for _,x in pairs(j._all.it) do push(xys,{x=x,y="rest"}) end
215    return merge( ranges(xys,i,  ykind or SYM,
216                         (#xys)^your.dull,
217                         xpect{i,j}*your.Small)) end
218 -- -------------------------------------------------------------------------
219 RANGE=class{}
220 function RANGE.new(col,lo,hi,ys)
221    return new({n=0, col=col, lo=lo, hi=hi or lo, ys=ys or SYM()},RANGE) end
222
223 function RANGE.__lt(i,j) return i:div() < j:div() end
224
225 function RANGE.__tostring(i)
226    if i.lo == i.hi        then return fmt("%s == %s", i.col.txt, i.lo) end
227    if i.lo == -math.huge then return fmt("%s < %s",  i.col.txt, i.hi) end
228    if i.hi ==  math.huge then return fmt("%s >= %s", i.col.txt, i.lo) end
229    return fmt("%s <= %s < %s", i.lo, i.col.txt, i.hi) end
230
231 function RANGE.add(i,x,y,inc)
232    inc  = inc or 1
233    i.n  = i.n + inc
234    i.hi = math.max(x,i.hi)
235    i.ys:add(y, inc) end
236
237 function RANGE.div(i) return i.ys:div() end
238
239 function RANGE.selects(i,row,     x)
240    x=row.has[col.at]; return x=="?" or i.lo<=x and x<i.hi end
241 -- -------------------------------------------------------------------------
242 SAMPLE=class{}
243 function SAMPLE.new() return new({n=0,it={},ok=false,max=your.ample},SAMPLE) end
244
245 function SAMPLE.add(i,x,     pos)
246    i.n = i.n +1
247    if      #i.it < i.max       then pos= #i.it + 1
248    elseif rand() < #i.it/i.n then pos= #i.it * rand() end
249    if pos then i.ok = false; i.it[pos//1]= x end end
250
251 function SAMPLE.all(i) if not i.ok then i.ok=true;sort(i.it)end; return i.it end
252 -- -------------------------------------------------------------------------
253 SYM=class{}
254 function SYM.new(at,s)
255    return new({at=at or 0,txt=s or "",has={},n=0,most=0,mode=nil},SYM) end
256
257 function SYM.__tostring(i)
258    return fmt("SYM{:at %s :txt %s :mode %s :has %s}",
259             i.at, i.txt, i.mode, o(i.has)) end
260
261 function SYM.add(i,x, inc)
262    if x ~= "?" then
263      inc = inc or 1
264      i.n = i.n+inc
265      i.has[x] = inc  + (i.has[x] or 0)
266      if i.has[x] > i.most then i.most, i.mode = i.has[x], x end end
267    return x end
268
269 function SYM.dist(i,a,b) return a=="?" and b=="?" and 1 or a==b and 0 or 1 end
270
271 function SYM.div(i,     e)
272    e=0;for _,v in pairs(i.has) do e=e - v/i.n*math.log(v/i.n,2) end; return e end
273
274 function SYM.merge(i,j,     k)
275    k= SYM(i.at, i.txt)
276    for x,count in pairs(i.has) do k:add(x,count) end
277    for x,count in pairs(j.has) do k:add(x,count) end
278    return k end
279
280 function SYM.mid(i) return i.mode end
281
282 function SYM.ranges(i,j,     t)
283    t = {}
284    for _,pair in pairs{{i.has,"bests"}, {j.has,"rests"}} do
285      for x,inc in pairs(pair[1]) do
286        t[x] = t[x] or RANGE(i,x)
287        print("inc",i.txt,inc)
288        t[x]:add(x, pair[2], inc) end end
289    return map(t) end
```

page 4

```lua
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    tasks = your.task=="all" and slots(go) or {your.task}
326    defaults=copy(your)
327    our.failures=0
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()))
```

```lua
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 -- 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"([^,]+)") 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
```

```lua
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,col1 in pairs(a.cols.x) do
456     col2 = b.cols.x[n]
457     print("")
458     for _, range in pairs(col1:ranges(col2)) do
459       print(col1.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() )
```