

```

1  #!/usr/bin/env lua
2  ---
3  ---
4  ---
5  ---
6  ---
7  ---
8  local F=require"fun"
9  local the=F.options[[
10
11  ./duo.lua [OPTIONS]
12  (c)2022 Tim Menzies, MIT license
13
14  Data miners using/used by optimizers.
15  Understand N items after log(N) probes, or less.
16
17  OPTIONS:
18  -ample      when enough is enough      = 512
19  -Debug      on error, dump stack and halt = false
20  -enough     use (#t)^enough             = .5
21  -far        how far to go               = .9
22  -file       read data from file         = ../etc/data/auto93.csv
23  -help       show help                   = false
24  -p          distance coefficient         = 2
25  -seed       random number seed          = 10019
26  -task       start up actions            = donothing]]
27
28  local EGS, NUM, RANGE, SYM = {}, {}, {}, {}
29  local map,fmt,new,sort,push,o,oo = F.map,F.fmt,F.new,F.sort,F.push,F.oo,F.oo
30  local any = F.any
31
32  function RANGE.new(k,col,lo,hi,b,B,r,R)
33  return new(k,{col=col,lo=lo,hi=hi or lo,b=b,B=B,r=r,R=R}) end
34
35  function RANGE.__lt(i,j) return i:val() < j:val() end
36  function RANGE.merge(i,j,k, lo,hi)
37  lo = math.min(i.lo, j.lo)
38  hi = math.max(i.hi, j.hi)
39  k = RANGE:new(i.col,lo,hi,i.b+j.b,i.B,i.r+j.r, j.R)
40  if k:val() > i:val() and j:val() then return k end end
41
42  function RANGE.__tostring(i)
43  if i.lo == i.hi then return fmt("%s==%s", i.col.txt, i.lo) end
44  if i.lo == -math.huge then return fmt("%s<%s", i.col.txt, i.hi) end
45  if i.hi == math.huge then return fmt("%s>=%s", i.col.txt, i.lo) end
46  return fmt("%s<=%s<%s", i.lo, i.col.txt, i.hi) end
47
48  function RANGE.val(i, z,B,R)
49  z=1E-31; B,R = i.B+z, i.R+z; return (i.b/B)^2/(i.b/B + i.r/R) end
50
51  function RANGE.selects(i,row, x)
52  x=row.has[col.at]; return x=="?" or i.lo<=x and x<i.hi end
53
54  function NUM.new(k,at,s)
55  return new(k,{at=at,txt=s,w=s:find"- " and -1 or 1,_has={},
56  ok=false, lo=math.huge, hi=math.huge}) end
57
58  function NUM.add(i,x)
59  if x == "?" then
60  i.ok = false
61  push(i._has, x)
62  if x < i.lo then i.lo = x end
63  if x > i.hi then i.hi = x end end
64  return x end
65
66  function NUM.dist(i,a,b)
67  if a=="?" and b=="?" then a,b=1,0
68  elseif a=="?" then b = i:norm(b); a=b>.5 and 0 or 1
69  elseif b=="?" then a = i:norm(a); b=a>.5 and 0 or 1
70  else
71  a, b = i:norm(a), i:norm(b) end
72  return math.abs(a-b) end
73
74  function NUM.has(i)
75  if not i.ok then sort(i._has); i.ok=true end; return i._has end
76
77  function NUM.norm(i,x)
78  return i.hi - i.lo<1E-9 and 0 or (x - i.lo)/(i.hi - i.lo) end
79
80  -- compare to old above
81  function NUM.ranges(i,j,lo,hi)
82  local z,is,js,lo,hi,m0,m1,m2,n0,n1,n2,step,most,best,r1,r2
83  is,js = i:has(), j:has()
84  lo = math.min(is[1], js[1])
85  hi = math.max(is[#is], js[#js])
86  gap, max = (hi - lo)/16, -1
87  for x=lo,hi,gap do
88  -- col, lo hi, b B r R
89  local b =
90  RANGE:new(i,lo,hi,
91  if hi-lo < 2*gap then
92  z = 1E-32
93  m0, m2 = fun.search(is, lo),fun.bsearch(is, hi+z)
94  n0, n2 = fun.bsearch(js, lo),fun.bsearch(js, hi+z)
95  -- col,lo hi,b B r R
96  best = nil
97  for mid in lo,hi,gap do
98  if mid > lo and k < hi then
99  m1 = bsearch(is, mid+z)
100 n1 = bsearch(js, mid+z)
101 r1 = RANGE:new(i, lo,mid,m1-m0,i.n,m2-(m1+1),j.n)
102 r2 = RANGE:new(i, mid+z,hi, n1-n0,i.n,n2-(n1+1),j.n)
103 if r1:val() > max then best, max = r1, r1:val() end
104 if r2:val() > max then best, max = r2, r2:val() end end end end
105 if best
106 then return i:ranges(j, best.lo, best.hi)
107 else return RANGE:new(i, lo,hi,m2-m0,i.n,n2-n0,j.n) end end
108
109  function SYM.new(k,at,s) return new(k,{at=at,txt=s,_has={}}) end
110  function SYM.add(i,x)
111  if x=="?" then i._has[x]=1+(i._has[x] or 0)end;return x end
112
113  function SYM.dist(i,a,b) return a=="?" and b=="?" and 1 or a==b and 0 or 1 end
114  function SYM.has(i) return i.has end
115  function SYM.ranges(i,j)
116  return lib.mapp(i._has,
117  function(x,n) return RANGE:new(i,x,x,n,i.n,(j._has[x] or 0),j.n) end) end
118
119  function EGS.new(k,file, i)
120  i = new(k,{rows={}, cols=nil, x={}, y={}})
121  if file then for row in F.rows(file) do i:add(row) end end
122  return i end
123
124  function EGS.add(i,t)
125  local add,now,where = function(col) return col:add(t[col.at]) end
126  if i.cols then
127  push(i._rows, map(i.cols, add))
128  else
129  i.cols = {}
130  for n,x in pairs(t) do
131  now = push(i.cols, (x:find"[A-Z]" and NUM or SYM):new(n,x))
132  if not x:find"." then
133  push((x:find"+" or x:find"-") and i.y or i.x, now) end end end end
134
135  function EGS.clone(i,init, j)
136  j = EGS:new()

```

```

137  j:add(map(i.cols, function(col) return col.txt end))
138  for _,row in pairs(inits or {}) do j = j:add(row) end
139  return j end
140
141  function EGS.cluster(i,top,lvl, tmp1,tmp2,left,right)
142  top = top or i
143  lvl = lvl or 0
144  print(fmt("%s%s", string.rep(".",lvl),#i._rows))
145  if #i._rows >= 2*(#top._rows)^the.enough then
146  tmp1, tmp2 = top:half(i._rows)
147  if #tmp1._rows < #i._rows then left = tmp1:cluster(top,lvl+1) end
148  if #tmp2._rows < #i._rows then right = tmp2:cluster(top,lvl+1) end
149  end
150  return (here=i, left=left, right=right) end
151
152  function EGS.dist(i,r1,r2)
153  local d,n,inc = 0, (#i.x)+1E-31
154  for _,col in pairs(i.x) do
155  inc = col:dist(r1[col.at], r2[col.at])
156  d = d + inc^the.p end
157  return (d/n)^(1/the.p) end
158
159  function EGS.far(i,r1,rows, act,tmp)
160  act = function(r2) return {r2, i:dist(r1,r2)} end
161  tmp = sort(map(rows,act), F.seconds)
162  return table.unpack(tmp[#tmp^the.far//1] ) end
163
164  function EGS.half(i,rows)
165  print(11)
166  local some,left,right,c,cosine,lefts,rights
167  rows = rows or i._rows
168  some = #rows > the.ample and F.many(rows, the.ample) or rows
169  left = i:far(any(rows), some)
170  right,c = i:far(left, some)
171  function cosine(r, a,b)
172  a, b = i:dist(r,left), i:dist(r,right); return {(a^2+c^2-b^2)/(2*c),r} end
173  lefts,rights = i:clone(), i:clone()
174  for n,pair in pairs(sort(map(rows,cosine), F.seconds)) do
175  (n <= (#rows)/2 and lefts or rights):add(pair[2] ) end
176  return lefts,rights,left,right,c end
177
178  ---
179  local no,go={},{}
180  local asserts=F.asserts
181
182  function go.half( a,b)
183  a,b=EGS:new(the.file):half()
184
185  end
186
187  function go.any( t,x,n)
188  t={}; for i=1,10 do t[i+#t] = i end
189  n=0; for i=1,5000 do x=F.any(t); n = 1 <= x and x <=10 and n+1 or 0 end
190  asserts(n==5000,"any") end
191
192  function go.bsearch( t,x,a,b)
193  t={}
194  for j =1,10^6 do push(t,100*math.random()//1) end
195  table.sort(t);
196  for j =1,1000 do
197  x=F.any(t)
198  a,b = F.brange(t,x)
199  assert(t[a-1] ~= x)
200  assert(t[b+1] ~= x)
201  for k=a,b do assert(t[k] == x) end end end
202
203  function no.fail() asserts(fail,"checking crashes"); print(no.thi.ng) end
204  function go.oo( u) oo{10,20,30} end
205
206  function go.rows( t)
207  for row in F.rows(the.file) do t=row end
208  asserts(type(t[1])=="number","is number")
209  asserts(t[1]==4, "is four")
210  asserts(#t==8,"is eight") end
211
212  function go.egs( i,t)
213  i=EGS:new(the.file); map(i.y,oo); asserts(i.y[1].lo==1613,"lo")
214  t=i.y[1]:has(); asserts(1613==t[1],"lo2") asserts(5140== t[#t],"hi");
215  asserts(i.y[1].ok,"ok") end
216
217  function go.dist( i, t,a,b,d)
218  i=EGS:new(the.file)
219  t= i._rows
220  for j=1,100 do
221  a,b= any(t), any(t)
222  d= i:dist(a,b)
223  assert(0<= d and d <= 1) end end
224
225  the(go)

```