

```

1  #!/usr/bin/env lua
2  ---
3  
4  ---
5  --- .lua
6  ---
7  ---
8  local F=require"fun"; local the=F.options[[
9
10 ./duo.lua [OPTIONS]
11 (c)2022 Tim Menzies, MIT license
12
13 Data miners using/used by optimizers.
14 Understand N items after log(N) probes, or less.
15
16 OPTIONS:
17 -ample      when enough is enough      = 512
18 -Debug      on error, dump stack and halt = false
19 -enough      use (#t)^enough            = .5
20 -far         how far to go              = .9
21 -file        read data from file         = data/duo93.csv
22 -help        show help                  = false
23 -p           distance coefficient        = 2
24 -seed        random number seed         = 10019
25 -task        start up actions           = donothing]]
26
27 local EGS, NUM, RANGE, SYM = {}, {}, {}, {}
28 local map,fmt,new,sort,push,o,oo = F.map,F.fmt,F.new,F.sort,F.push,F.oo,F.oo
29 local any = F.any
30
31 function RANGE.new(k,col,lo,hi,b,B,r,R)
32   return new(k,{col=col,lo=lo,hi=hi or lo,b=b,B=B,r=r,R=R}) end
33
34 function RANGE.__lt(i,j) return i:val() < j:val() end
35 function RANGE.merge(i,j,k, lo,hi)
36   lo = math.min(i.lo, j.lo)
37   hi = math.max(i.hi, j.hi)
38   k = RANGE:new(i.col,lo,hi,i.b+j.b,i.B,i.r+j.r, j.R)
39   if k:val() > i:val() and j:val() then return k end end
40
41 function RANGE.__tostring(i)
42   if i.lo == i.hi then return fmt("%s==%s", i.col.txt, i.lo) end
43   if i.lo == -math.huge then return fmt("%s<%s", i.col.txt, i.hi) end
44   if i.hi == math.huge then return fmt("%s>%s", i.col.txt, i.lo) end
45   return fmt("%s<=%s<%s", i.lo, i.col.txt, i.hi) end
46
47 function RANGE.val(i, z,B,R)
48   z=E-31; B,R = i.B+z, i.R+z; return (i.b/B)^2/( i.b/B + i.r/R) end
49
50 function RANGE.selects(i,row, x)
51   x=row.has[col.at]; return x=="?" or i.lo<=x and x<=i.hi end
52
53 function NUM.new(k,at,s)
54   return new(k,{at=at,txt=s,w=s:find("-" and -1 or 1,_has={},
55     ok=false, lo=math.huge, hi=-math.huge)) end
56
57 function NUM.add(i,x)
58   if x ~= "?" then
59     i.ok = false
60     push(i._has, x)
61     if x < i.lo then i.lo = x end
62     if x > i.hi then i.hi = x end end
63   return x end
64
65 function NUM.dist(i,a,b)
66   if a=="?" and b=="?" then a,b=1,0
67   elseif a=="?" then b = i:norm(b); a=b^.5 and 0 or 1
68   elseif b=="?" then a = i:norm(a); b=a^.5 and 0 or 1
69   else
70     a, b = i:norm(a), i:norm(b) end
71   return math.abs(a-b) end
72
73 function NUM.has(i)
74   if not i.ok then sort(i._has); i.ok=true end; return i._has end
75
76 function NUM.norm(i,x)
77   return i.hi - i.lo<E-9 and 0 or (x - i.lo)/(i.hi - i.lo) end
78
79 -- compare to old above
80 function NUM.ranges(i,j,lo,hi)
81   local z,is,js,lo,hi,m0,m1,m2,n0,n1,n2,step,most,best,r1,r2
82   is,js = i:has(), j:has()
83   lo = math.min(is[1], js[1])
84   hi = math.max(is[#is], js[#js])
85   gap, max = (hi - lo)/16, -1
86   for x=lo,hi,gap do
87     -- col, lo hi, b B r R
88     local b =
89     RANGE:new(i,lo,hi,
90       if hi-lo < 2*gap then
91         z = 1E-32
92         m0, m2 = fun.search(is, lo),fun.bsearch(is, hi+z)
93         n0, n2 = fun.bsearch(js, lo),fun.bsearch(js, hi+z)
94         -- col,lo hi,b B r R
95         best = nil
96         for mid in lo,hi,gap do
97           if mid > lo and k < hi then
98             m1 = bsearch(is, mid+z)
99             n1 = bsearch(js, mid+z)
100            r1 = RANGE:new(i, lo,mid,m1-m0,i.n,m2-(m1+1),j.n)
101            r2 = RANGE:new(i, mid+z,hi, n1-n0,i.n,n2-(n1+1),j.n)
102            if r1:val() > max then best, max = r1, r1:val() end
103            if r2:val() > max then best, max = r2, r2:val() end end end end
104          if best
105            then return i:ranges(j, best.lo, best.hi)
106            else return RANGE:new(i, lo,hi,m2-m0,i.n,n2-n0,j.n) end end
107
108 function SYM.new(k,at,s) return new(k,{at=at,txt=s,_has={}}) end
109 function SYM.add(i,x)
110   if x=="?" then i._has[x]=1+(i._has[x] or 0)end;return x end
111
112 function SYM.dist(i,a,b) return a=="?" and b=="?" and 1 or a==b and 0 or 1 end
113 function SYM.has(i) return i.has end
114 function SYM.ranges(i, j)
115   return lib.mapp(i._has, -- col lohi b B r R
116     function(x,n) return RANGE:new(i,x,x,n,i.n,(j._has[x] or 0),j.n) end) end
117
118 function EGS.new(k,file, i)
119   i = new(k,{_rows={}, cols=nil, x={}, y={}})
120   if file then for row in F.rows(file) do i:add(row) end end
121   return i end
122
123 function EGS.add(i,t)
124   local add,now,where = function(col) return col:add(t[col.at]) end
125   if i.cols
126     then push(i._rows, map(i.cols, add))
127     else i.cols = {}
128     for n,x in pairs(t) do
129       now = (x:find("[A-Z]" and NUM or SYM):new(n,x)
130       push(i.cols, now)
131       if not x:find"." then
132         where = (x:find"+" or x:find"-") and i.y or i.x
133         push(where, 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(init 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", 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
180 local no,go={},{}
181 local asserts=F.asserts
182
183 function go.half( a,b)
184   a,b=EGS:new(the.file):half()
185
186   end
187
188 function go.any( t,x,n)
189   t={}; for i=1,10 do t[i+#t] = i end
190   n=0; for i=1,5000 do x=F.any(t); n = 1 <= x and x <=10 and n+1 or 0 end
191   asserts(n==5000,"any") end
192
193 function go.bsearch( t,x,a,b)
194   t={}
195   for j =1,10^6 do push(t,100*math.random()/1) end
196   table.sort(t);
197   for j =1,1000 do
198     x=F.any(t)
199     a,b = F.brange(t,x)
200     assert(t[a-1] ~= x)
201     assert(t[b+1] ~= x)
202     for k=a,b do assert(t[k] == x) end end end
203
204 function no.fail() asserts(fail,"checking crashes"); print(no.thi.ng) end
205 function go.oo( u) oo{10,20,30} end
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)

```