

```

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  -rnd        default round               = %5.2f
26  -seed       random number seed         = 10019
27  -task       start up actions            = donothing]]
28
29  local EGS, NUM, RANGE, SYM = {}, {}, {}, {}
30  local any, asserts, brange, firsts, fmt, many, map =
31    F.any,F.asserts,F.brange,F.firsts,F.fmt,F.many,F.map
32  local new, o, oo, push, rows, seconds, sort =
33    F.new,F.o,F.oo,F.push,F.rows,F.seconds,F.sort
34
35  --- ## RANGE
36  function RANGE.new(k,col,lo,hi,b,B,r,R)
37    return new(k,{col=col,lo=lo,hi=hi or lo,b=b,B=B,r=r,R=R}) end
38
39  function RANGE.__lt(i,j) return i:val() < j:val() end
40  function RANGE.merge(i,j,k, lo,hi)
41    lo = math.min(i.lo, j.lo)
42    hi = math.max(i.hi, j.hi)
43    k = RANGE:new(i.col,lo,hi,i.b+j.b,i.B,i.r+j.r, j.R)
44    if k:val() > i:val() and j:val() then return k end end
45
46  function RANGE.__tostring(i)
47    if i.lo == i.hi then return fmt("%s==%s", i.col.txt, i.lo) end
48    if i.lo == -math.huge then return fmt("%s<%s", i.col.txt, i.hi) end
49    if i.hi == math.huge then return fmt("%s>=%s", i.col.txt, i.lo) end
50    return fmt("%s<=%s<%s", i.lo, i.col.txt, i.hi) end
51
52  function RANGE.val(i, z,B,R)
53    z=1E-31; B,R = i.B+z, i.R+z; return (i.b/B)^2/( i.b/B + i.r/R) end
54
55  function RANGE.selects(i,row, x)
56    x=row.has[col.at]; return x=="?" or i.lo<=x and x<i.hi end
57
58

```

```

58 --- ## NUM
59 function NUM.new(k,at,s)
60     return new(k,{at=at,txt=s,w=s:find("-",1,1,has={},
61         ok=false, lo=math.huge, hi=-math.huge)) end
62
63 function NUM.add(i,x)
64     if x ~= "?" then
65         i.ok = false
66         push(i._has, x)
67         if x < i.lo then i.lo = x end
68         if x > i.hi then i.hi = x end end
69     return x end
70
71 function NUM.dist(i,a,b)
72     if a=="?" and b=="?" then a,b=1,0
73     elseif a=="?" then b = i:norm(b); a=b>.5 and 0 or 1
74     elseif b=="?" then a = i:norm(a); b=a>.5 and 0 or 1
75     else
76         a, b= i:norm(a), i:norm(b) end
77     return math.abs(a-b) end
78
79 function NUM.has(i)
80     if not i.ok then sort(i._has); i.ok=true end; return i._has end
81
82 function NUM.mid(i, a) a=i:has(); return a[#a//2] end
83
84 function NUM.norm(i,x)
85     return i.hi - i.lo<1E-9 and 0 or (x - i.lo)/(i.hi - i.lo) end
86
87 -- compare to old above
88 -- function NUM.ranges(i,j,lo,hi)
89 --     local z,is,js,lo,hi,m0,m1,m2,n0,n1,n2,step,most,best,r1,r2
90 --     is,js = i:has(), j:has()
91 --     lo = math.min(is[1], js[1])
92 --     hi = math.max(is[#is], js[#js])
93 --     gap, max = (hi - lo)/16, -1
94 --     for x=lo,hi,gap do
95 --         -- col, lo hi, b B r R
96 --         local b =
97 --             RANGE:new(i,lo,hi,
98 --                 if hi-lo < 2*gap then
99 --                     z = 1E-32
100 --                     m0, m2 = fun.bsearch(is, lo),fun.bsearch(is, hi+z)
101 --                     n0, n2 = fun.bsearch(js, lo),fun.bsearch(js, hi+z)
102 --                     -- col,lo hi,b B r R
103 --                     best = nil
104 --                     for mid in lo,hi,gap do
105 --                         if mid > lo and k < hi then
106 --                             n1 = bsearch(is, mid+z)
107 --                             n1 = bsearch(js, mid+z)
108 --                             r1 = RANGE:new(i, lo,mid,m1-m0,i.n,m2-(m1+1),j.n)
109 --                             r2 = RANGE:new(i, mid+z,hi, n1-n0,i.n,n2-(n1+1),j.n)
110 --                             if r1:val() > max then best, max = r1, r1:val() end
111 --                             if r2:val() > max then best, max = r2, r2:val() end end end end
112 --                     if best
113 --                         then return i:ranges(j, best.lo, best.hi)
114 --                         else return RANGE:new(i, lo,hi,m2-m0,i.n,n2-n0,j.n) end end
115 --
116
117 --- ## SYM
118 function SYM.new(k,at,s) return new(k,{at=at,txt=s,_has={},mode=nil,most=0}) end
119 function SYM.add(i,x)
120     if x=="?" then
121         i._has[x]=1+(i._has[x] or 0)
122         if i._has[x] > i.most then i.most, i.mode = i._has[x],x end
123     end
124     return x end
125
126 function SYM.dist(i,a,b) return a=="?" and b=="?" and 1 or a==b and 0 or 1 end
127 function SYM.has(i) return i._has end
128 function SYM.mid(i) return i.mode end
129 function SYM.ranges(i,j)
130     return lib.mapp(i._has, -- col lohib B r R
131         function(x,n) return RANGE:new(i,x,x,n,i.n,{j._has[x] or 0},j.n) end) end
132

```

```

132 --- ## EGS
133 function EGS.new(k,file, i)
134 i= new(k,{_rows={}, cols=nil, x={}, y={}})
135 if file then for row in rows(file) do i:add(row) end end
136 return i end
137
138 function EGS.add(i,t)
139 local add,now,where = function(col) return col:add(t[col.at]) end
140 if i.cols then
141 push(i._rows, map(i.cols, add))
142 else
143 i.cols = {}
144 for n,x in pairs(t) do
145 now = push(i.cols, {x:find"^[A-Z]" and NUM or SYM}:new(n,x))
146 if not x:find"." then
147 push({x:find"+" or x:find"-"} and i.y or i.x, now) end end end end
148
149 function EGS.clone(i,init, j)
150 j = EGS:new()
151 j:add(map(i.cols, function(col) return col.txt end))
152 for _,row in pairs(init or {}) do j = j:add(row) end
153 return j end
154
155 function EGS.mid(i,cols)
156 return map(cols or i.y, function(col) return col:mid() end) end
157
158 function EGS.dist(i,r1,r2)
159 local d,n,inc = 0, (#i.x)+1E-31
160 for _,col in pairs(i.x) do
161 inc = col:dist(r1[col.at], r2[col.at])
162 d = d + inc^the.p end
163 return (d/n)^(1/the.p) end
164
165 function EGS.far(i,r1,rows, act,tmp)
166 act = function(r2) return {r2, i:dist(r1,r2)} end
167 tmp = sort(map(rows,act), seconds)
168 return table.unpack(tmp[#tmp*the.far//1]) end
169
170 function EGS.half(i,rows)
171 local some,left,right,c,cosine,lefts,rights
172 rows = rows or i._rows
173 some = #rows > the.ample and many(rows, the.ample) or rows
174 left = i:far(any(rows), some)
175 right,c = i:far(left, some)
176 function cosine(r, a,b)
177 a, b = i:dist(r,left), i:dist(r,right); return {(a^2+c^2-b^2)/(2*c),r} end
178 lefts,rights = i:clone(), i:clone()
179 for n,pair in pairs(sort(map(rows,cosine), firsts)) do
180 (n <= (#rows)/2 and lefts or rights):add(pair[2]) end
181 return lefts,rights,left,right,c end
182
183 local rnd,show
184 function EGS.cluster(i, top)
185 local c,lefts0, rights0, lefts, rights, left, right=0
186 top = top or i
187 if #i._rows >= 2*(#top._rows)^the.enough then
188 lefts0, rights0, left, right, c = top:half(i._rows)
189 lefts = lefts0:cluster(top)
190 rights = rights0:cluster(top)
191 end
192 return {here=i, lefts=lefts, rights=rights, left=left, right=right, c=c} end
193
194 function rnd(x)
195 return fmt(type(x)=="number" and x==x//1 and the.rnd or"%s",x) end
196
197 function show(t,lvl)
198 lvl = lvl or ""
199 if t then
200 --if t.lefts
201 print(fmt("%s%s",lvl,#t.here._rows))
202 --else print(fmt("%s%s\t%s", lvl,#t.here._rows, t.here:mid())) end
203 show(t.lefts, lvl.."..")
204 show(t.rights,lvl.."..") end end
205
206

```

```

206 --- ## Tests and Demo
207 local no,go={},{}
208
209 function go.cluster( a)
210 a=EGS:new(the.file):cluster()
211 asserts(49==#a.lefts.lefts.here._rows) end
212
213 function go.half( a,b)
214 local top =EGS:new(the.file)
215 local lefts,rights,left,right,c=top:half()
216 asserts(top:dist(left,right) > .75)
217 end
218
219 function go.any( t,x,n)
220 t={}; for i=1,10 do t[1+#t] = i end
221 n=0; for i=1,5000 do x=any(t); n= 1 <= x and x <=10 and n+1 or 0 end
222 asserts(n==5000,"any") end
223
224 function go.bsearch( t,x,a,b)
225 t={}
226 for j =1,10^6 do push(t,100*math.random())//1 end
227 table.sort(t);
228 for j =1,1000 do
229 x=any(t)
230 a,b = brange(t,x)
231 assert(t[a-1] ~= x)
232 assert(t[b+1] ~= x) ----
233 for k=a,b do assert(t[k] == x) end end end
234
235 function no.fail() asserts(fail,"checking crashes"); print(no.thi.ng) end
236 function go.oo( u) asserts("10 20 30)" == fmt("%s",o{10,20,30}), "table") end
237 function go.rows( t)
238 for row in rows(the.file) do t=row end
239 asserts(type(t[1])=="number", "is number")
240 asserts(t[1]==4, "is four")
241 asserts(#t==8, "is eight") end
242
243 function go.egs( i,t)
244 i=EGS:new(the.file)
245 asserts(i.y[1].lo==1613,"lo")
246 t=i.y[1]:has(); asserts(1613==t[1],"lo2") asserts(5140== t[#t],"hi");
247 asserts(i.y[1].ok,"ok") end
248
249 function go.dist( i, t,a,b,d)
250 i=EGS:new(the.file)
251 t= i._rows
252 for j=1,100 do
253 a,b= any(t), any(t)
254 d= i:dist(a,b)
255 assert(0<= d and d <= 1) end end
256
257 the(go)

```