

```

1  #!/usr/bin/env lua
2  -- vim : ft=lua et sts=2 sw=2 ts=2 :
3  local help = {}
4  sl (OPTIONS)
5  Sublime's unsupervised bifurcation: let's infer minimal explanations.
6  (c) 2022, Tim Menzies
7
8  OPTIONS:
9  -D      stack dump on assert fails
10 -d f     data file                = etc/data/auto93.csv
11 -e f     recurse until rows^enough = .5
12 -f F     far                      = .9
13 -k P     max kept items           = 512
14 -p P     distance coefficient      = 2
15 -S P     set seed                  = 10019
16 -t S     start up action (or 'all') = nothing
17 -h       show help
18
19 KEY: f=filename F=float P=posint S=string ]]
20
21
22
23
24
25
26
27
28
29
30
31
32
33 -- Redistribution and use in source and binary forms, with or without
34 -- modification, are permitted provided that the following conditions are met:
35
36 -- 1. Redistributions of source code must retain the above copyright notice,
37 --    this list of conditions and the following disclaimer.
38
39 -- 2. Redistributions in binary form must reproduce the above copyright notice,
40 --    this list of conditions and the following disclaimer in the documentation
41 --    and/or other materials provided with the distribution.
42
43 -- THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS"
44 -- AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
45 -- IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR PARTICULAR PURPOSE ARE
46 -- DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR CONTRIBUTORS BE LIABLE
47 -- FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL
48 -- DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR
49 -- SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER
50 -- CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY,
51 -- OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE
52 -- OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
53
54 local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end --used later (to find rogues)
55 local any, asserts, big, cli, fails, firsts, fmt, goalp, help, ignorep, klassp
56 local lessp, map, main, many, max, min, morep, new, nump, o, oo, per, push
57 local r, rows, slots, sort, sum, the, thing, things, unpack
58 local CLUSTER, COLS, EGS, NUM, ROWS, SKIP, SOME, SYM = {}, {}, {}, {}, {}, {}, {}
59
60 function cli(want,x)
61   for n,got in ipairs(arg) do if got==want then
62     x = x==false and true or x==true and "false" or arg[n+1] end end
63   if x=="false" then return false else return tonumber(x) or x end end
64
65 the = {dump      = cli("-D", false),
66        data      = cli("-d", "etc/data/auto93.csv"),
67        enough    = cli("-e", .5),
68        help      = cli("-h", false),
69        far       = cli("-f", .9),
70        keep      = cli("-k", 512),
71        p         = cli("-p", 2),
72        seed      = cli("-S", 10019),
73        todo      = cli("-t", "nothing")}
74
75 -- git rid of SOME for rows
76 -- col = NUM | SYM
77 -- COLS = all:[col]+, x:[col]*, y:[col]*, klass:col*
78 -- ROWS = cols:COLS, rows:SOME
79
80

```

```

81 --
82 -- FUNCTIONS
83 --
84 --
85 -- strings
86 fmt = string.format
87
88 -- maths
89 big = math.huge
90 max = math.max
91 min = math.min
92 r   = math.random
93
94 -- column headers
95 function goalp(x) return morep(x) or lessp(x) or klassp(x) end
96 function ignorep(x) return x:find"$" end
97 function klassp(x) return x:find"$" end
98 function lessp(x) return x:find"$" end
99 function morep(x) return x:find"$" end
100 function nump(x) return x:find"[A-Z]" end
101
102 -- tables
103 unpack = table.unpack
104 function any(t) return t[#t] end
105 function firsts(a,b) return a[1] < b[1] end
106 function many(t,n, u) u={}; for i=1,n do push(u,any(t)) end; return u end
107 function per(t,p) return t[ (#t*(p or .5))/1 ] end
108 function push(t,x) table.insert(t,x); return x end
109 function sort(t,f) table.sort(t,f); return t end
110
111 -- meta
112 function map(t,f, u) u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
113 function sum(t,f, n) n=0; for _,v in pairs(t) do n=n+f(v) end; return n end
114 function slots(t, u)
115   u={}
116   for k,v in pairs(t) do k=tostring(k);if k:sub(1,1)~="_" then push(u,k) end end
117   return sort(u) end
118
119 -- print tables, recursively
120 function oo(t) print(o(t)) end
121 function o(t)
122   if type(t)~="table" then return tostring(t) end
123   local key=function(k) return fmt("%.5s",k,o(t[k])) end
124   local u = #t>0 and map(t,o) or map(slots(t),key)
125   return '{ '..table.concat(u, " " ..")" end
126
127 -- strings to things
128 function rows(file, x)
129   file = io.input(file)
130   return function()
131     x=io.read(); if x then return things(x) else io.close(file) end end end
132
133 function thing(x)
134   x = x:match"%s*(.-)%s*"
135   if x=="true" then return true elseif x=="false" then return false end
136   return tonumber(x) or x end
137
138 function things(x,sep, t)
139   t={}
140   for y in x:gmatch(sep or "([^\+])") do push(t,thing(y)) end
141   return t end
142
143 -- errors
144 fails=0
145 function asserts(test, msg)
146   print(test and "PASS: " or "FAIL: ",msg or "")
147   if not test then
148     fails=fails+1
149     if the.dump then assert(test,msg) end end end
150
151 -- objects
152 function new(k,t) k.__index=k; k.__tostring=o; return setmetatable(t,k) end

```

```

153 -- CLASSES
154
155
156 -- COLS
157
158 function COLS.new(k,row, i)
159 i= new(k,{all={},x={},y={},names=row})
160 for at,txt in ipairs(row) do push(i.all, i:col(at,txt)) end
161 return i end
162
163 function COLS.add(i,t)
164 for _,col in pairs(i.all) do col:add( t[col.at] ) end
165 return t end
166
167 function COLS.col(i,at,txt, col)
168 if ignorep(txt) then return SKIP:new(at,txt) end
169 col = (nump(txt) and NUM or SYM):new(at,txt)
170 push(goalp(txt) and i.y or i.x, col)
171 if klassp(txt) then i.klass = col end
172 return col end
173
174 -- NUM
175 function NUM.new(k,n,s)
176 return new(k,{n=0,at=n or 0,txt=s or "",has=SOME:new(),ok=false,
177 w=lessp(s or "") and -1 or 1, lo=big, hi=-big}) end
178
179 function NUM.add(i,x)
180 if x ~= "?" then
181 i.n = i.n + 1
182 if i.has:add(x) then i.ok=false end
183 i.lo,i.hi = min(x,i.lo), max(x,i.hi); end end
184
185 function NUM.dist(i,x,y)
186 if x=="?" and y=="?" then return 1
187 elseif x=="?" then y=i:norm(y); x=y<0.5 and 1 or 0
188 elseif y=="?" then x=i:norm(x); y=x<0.5 and 1 or 0
189 else x,y = i:norm(x), i:norm(y) end
190 return math.abs(x-y) end
191
192 function NUM.mid(i) return per(i:sorted(), .5) end
193
194 function NUM.norm(i,x)
195 return math.abs(i.hi-i.lo)<1E-9 and 0 or (x-i.lo)/(i.hi - i.lo) end
196
197 function NUM.sorted(i)
198 if i.ok==false then table.sort(i.has.all); i.ok=true end
199 return i.has.all end
200
201 -- ROWS
202 function ROWS.new(k,init, i)
203 i = new(k,{rows=SOME:new(), cols=nil})
204 if type(init)=="string" then for row in rows(init) do i:add(row) end end
205 if type(init)=="table" then for row in init do i:add(row) end end
206 return i end
207
208 function ROWS.add(i,row)
209 if i.cols then i.rows:add( i.cols:add(row) )
210 else i.cols = COLS:new(row) end end
211
212 function ROWS.clone(i, j) j= ROWS:new(); j:add(i.cols.names);return j end
213
214 function ROWS.dist(i,row1,row2, d,fun)
215 function fun(col) return col:dist(row1[col.at], row2[col.at])^the.p end
216 return (sum(i.cols.x, fun) / #i.cols.x)^(1/the.p) end
217
218 function ROWS.far(i,row1,rows, fun)
219 function fun(row2) return (i:dist(row1,row2), row2) end
220 return unpack(per(sort(map(rows,fun),firsts), the.far)) end
221
222 function ROWS.half(i, top)
223 local some, top,c,x,y,tmp,mid,lefts,rights,_
224 some= many(i.rows.all, the.keep)
225 top = top or i
226 _,x = top:far(any(some), some)
227 c,y = top:far(x, some)
228 tmp = sort(map(i.rows.all, function(r) return top:project(r,x,y,c) end), first
229 s)
230 mid = #i.rows.all//2
231 lefts, rights = i:clone(), i:clone()
232 for at,row in pairs(tmp) do (at <=mid and lefts or rights):add(row[2]) end
233 return lefts,rights,x,y,c, tmp[mid] end
234
235 function ROWS.mid(i,cols)
236 return map(cols or i.cols.all, function(col) return col:mid() end) end
237
238 function ROWS.project(i, r,x,y,c, a,b)
239 a,b = i:dist(r,x), i:dist(r,y); return {(a^2 + c^2 - b^2)/(2*c), r} end
240
241 -- SKIP
242 function SKIP.new(k,n,s) return new(k,{n=0,at=at or 0,txt=s or ""}) end
243 function SKIP.add(i,x) return x end
244 function SKIP.mid(i) return "?" end
245
246 -- SOME
247 function SOME.new(k,keep) return new(k,{n=0,all={}, keep=keep or the.keep}) end
248 function SOME.add(i,x)
249 i.n = i.n+1
250 if #i.all < i.keep then push(i.all,x) ; return i.all
251 elseif r() < i.keep/i.n then i.all[r(#i.all)]=x; return i.all end end
252
253 -- SYM
254 function SYM.new(k,n,s) return new(k,{n=0,at=n or 0,txt=s or "",has={},most=0})
255 end
256 function SYM.dist(i,x,y) return (x=="?" and y=="?" and 1) or (x==y and 0 or 1) end
257 function SYM.mid(i) return i.mode end
258 function SYM.div(i, fun)
259 function fun(k, p) p = -i.has[k]/i.n; return -p*math.log(p,2) end
260 return sum(i.has, fun) end
261
262 function SYM.add(i,x,inc)
263 if x == "?" then
264 inc = inc or 1
265 i.n = i.n + inc
266 i.has[x] = inc + (i.has[x] or 0)
267 if i.has[x] > i.most then i.most,i.mode=i.has[x],x end end end
268
269 function SYM.merge(i,j, k)
270 k = SYM:new(i.at,i.txt)
271 for x,n in pairs(i.has) do k:add(x,n) end
272 for x,n in pairs(j.has) do k:add(x,n) end
273 ei, ej, ejk = i:div(), j:div(), k:div()
274 if i.n==0 or j.n==0 or .99*ek <= (i.n*ei + j.n*ej)/k.n then
275 return k end end
276
277 -- CLUSTER
278 function CLUSTER.new(k,sample,top)
279 local i,enough,left,right
280 top = top or sample
281 i = new(k, {here=sample})
282 enough = top.rows.n^the.enough
283 if sample.rows.n >= 2*enough then
284 left, right, i.x, i.y, i.c, i.mid = sample:half(top)
285 if left.rows.n < sample.rows.n then
286 i.left = CLUSTER:new(left, top)
287 i.right = CLUSTER:new(right, top) end end
288 return i end
289
290 function CLUSTER.show(i,pre, here)
291 pre = pre or ""
292 here=""
293 if not i.left and not i.right then here= o(i.here:mid(i.here.cols.y)) end
294 print(fmt("%6s:%-30s%s",i.here.rows.n, pre, here))
295 for _,kid in pairs(i.left, i.right) do
296 if kid then kid:show(pre .. ".") end end end
297
298 -----
299 -- DEMOS
300
301 function EGS.nothing() return true end
302 function EGS.the() oo(the) end
303 function EGS.rand() print(r()) end
304 function EGS.clone( r,s)
305 r = ROWS:new(the.data)
306 s = r:clone()
307 for _,row in pairs(r.rows.all) do s:add(row) end
308 asserts(r.cols.x[1].lo==s.cols.x[1].lo,"clone.lo")
309 asserts(r.cols.x[1].hi==s.cols.x[1].hi,"clone.hi")
310 end
311
312 function EGS.data( r)
313 r = ROWS:new(the.data)
314 asserts(r.cols.x[1].hi == 8, "data.columns") end
315
316 function EGS.dist( r,rows,n)
317 r = ROWS:new(the.data)
318 rows = r.rows.all
319 n = NUM:new()
320 for _,row in pairs(rows) do n:add(r:dist(row, rows[1])) end
321 oo(r.cols.x[2]:sorted()) end
322
323 function EGS.many( t)
324 t={}; for j=1,100 do push(t,j) end
325 print(oo(many(t, 10))) end
326
327 function EGS.far( r,c,row1,row2)
328 r = ROWS:new(the.data)
329 row1 = r.rows.all[1]
330 c,row2 = r:far(r.rows.all[1], r.rows.all)
331 print(c,"n",o(row1),"n", o(row2)) end
332
333 function EGS.half( r,c,row1,row2)
334 local lefts,rights,x,y,x
335 r = ROWS:new(the.data)
336 oo(r:mid(r.cols.y))
337 lefts,rights,x,y,c = r:half()
338 oo(lefts:mid(lefts.cols.y ))
339 oo(rights:mid(rights.cols.y))
340 end
341
342 function EGS.cluster(r)
343 r = ROWS:new(the.data)
344 CLUSTER:new(r):show() end
345
346 -- start-up
347 if the.help then print(help) else
348 local b4={}; for k,v in pairs(the) do b4[k]=v end
349 for _,todo in pairs(the.todo=="all" and slots(EGS) or (the.todo)) do
350 for k,v in pairs(b4) do the[k]=v end
351 math.randomseed(the.seed)
352 if type(EGS[todo])=="function" then EGS[todo]() end end end
353
354 for k,v in pairs(_ENV) do if not b4[k] then print("? ",k,type(v)) end end
355 os.exit(fails)

```