

```

1 local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
2 local the,help={},{}
3
4 lua 15.lua [OPTIONS]
5 L5 == a very little LUA learning lab
6 (c)2022, Tim Menzies, BSD 2-clause license
7
8 OPTIONS (for changing the inference):
9
10 -cohen -c F cohen's small effect size      = .35
11 -far   -F F look no further than "far"     = .9
12 -keep  -k items to keep in a number       = 512
13 -leaves -l leaf size                       = 5
14 -p     -P P distance calcs coefficient     = 2
15 -seed  -S P random number seed            = 10019
16 -some  -s look only at "some" items        = 512
17
18 OPTIONS (for housekeeping):
19
20 -dump  -d exit on error, with stacktrace = false
21 -file  -f S where to get data            = ../etc/data/auto93.csv
22 -help  -h show help                      = false
23 -rnd   -r S format string                 = %5.2f
24 -todo  -t S start-up action               = nothing
25
26 KEY: S=string, P=poisint, F=float
27 {}
28
29 -----
30
31
32
33
34
35
36
37 local function Sym(at,s)
38   return { is="Sym", -- type
39            at=at or 0, -- column index
40            name=s or "", -- column name
41            n=0, -- number of items summarized in this column
42            all={}, -- all|x| = n means we've seen "n" repeats of "x"
43            most=0, -- count of the most frequently seen symbol
44            mode=nil -- the most commonly seen letter
45          } end
46
47 local function Num(at,s)
48   return { is="Num", -- type
49            at=at or 0, -- column index
50            name=s or "", -- column name
51            n=0, -- number of items summarizes in this column
52            mu=0, -- mean (updated incrementally)
53            m2=0, -- second moment (updated incrementally)
54            sd=0, -- standard deviation
55            all={}, -- a sample of items seen so far
56            lo=1E31, -- lowest number seen
57            hi=-1E31, -- highest number seen
58            w=(s or ""):find("-$") and -1 or 1 -- "-1"= minimize and "1"= maximize
59          } end
60
61 local function Egs(names)
62   return { is="egs", -- type
63            all={}, -- all the rows
64            names=names, -- list of name
65            cols={}, -- list of all columns (Nums or Syms)
66            x={}, -- independent columns (nothing marked as "skip")
67            y={}, -- dependent columns (nothing marked as "skip")
68          } end
69
70 --[[
71 ## Coding Conventions
72 - "i" not "self"
73 - if something holds a list of thing, name the holding variable "all"
74 - no inheritance
75 - when you can, write functions down on one line
76 - all config items into a global "the" variable
77 - all the test cases (or demos) are "function Demo.xxx".
78 - random seed reset so carefully, just once, at the end of the code.
79 ]]

```

data

```

80
81
82
83
84
85
86 local r = math.random
87 local fmt = string.format
88 local function push(t,x) table.insert(t,x); return x end
89
90
91
92
93
94
95 local thing,things,file2things
96 function thing(x)
97   x = x:match("^%s*(-)%s*$")
98   if x=="true" then return true elseif x=="false" then return false end
99   return tonumber(x) or x end
100
101 function things(x,sep, t)
102   t={}; for y in x:gmatch(sep or "[^,]+") do push(t,thing(y)) end
103   return t end
104
105 function file2things(file, x)
106   file = io.input(file)
107   return function()
108     x=io.read();
109     if x then return things(x) else io.close(file) end end end
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178

```

utils

COERCE

GET,SET

TABLE

PRINT

START-UP

```

153 local Demo, ok = {fails=0}
154 function ok(test,msg)
155   print(test and "PASS:" or "FAIL:",msg or "")
156   if not test then
157     Demo.fails=Demo.fails+1
158     if the.dump then assert(test,msg) end end end
159
160 function Demo.main(todo,seed)
161   for k,one in pairs(todo=="all" and slots(Demo) or {todo}) do
162     if k ~= "main" and type(Demo[one]) == "function" then
163       math.randomseed(seed)
164       Demo[one]() end end
165   for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
166   return Demo.fails end
167
168 local function settings(txt, d)
169   d={}
170   txt:gsub("(^%s+)(%s)+(-[%s%+][^%s%+])",function(long,key,short,x)
171     for n,flag in ipairs(arg) do
172       if flag==short or flag==long then
173         x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
174     if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
175       d[key] = tonumber(x) or x end end
176   if d.help then print(help) end
177   return d end
178

```

```

179 ---
180 --- UPDATE COLS
181 ---
182
183 local nump,add
184 function nump(col) return col.w end
185
186 function add(i,x,inc,      sym1,num1)
187   function sym1()
188     i.all[x] = inc + (i.all[x] or 0)
189     if i.all[x] > i.most then i.most, i.mode = i.all[x], x end
190   end
191   function num1(      d)
192     for j=1,inc do
193       d = x - i.mu
194       i.mu = i.mu + d/i.n
195       i.m2 = i.m2 + d*(x - i.mu)
196       i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n-1))^0.5)
197       i.lo = math.min(x, i.lo)
198       i.hi = math.max(x, i.hi)
199       if #i.all < the.keep then push(i.all,x)
200       elseif r(i) < the.keep/i.n then i.all[r(#i.all)]=x end end
201     end
202     inc = inc or 1
203     if x ~= "?" then
204       i.n = i.n + inc
205       if nump(i) then num1() else sym1() end end
206     return x end
207 ---
208 --- MAKE DATA
209 ---
210
211 local header,data,file2Egs
212 function header(names, i,col)
213   i = Egs(names)
214   for at,name in pairs(names) do
215     col = push(i.cols, (name:find"^[A-Z]" and Num or Sym)(at,name))
216     if not name:find"$*" then
217       push(name:find"[+}$" and i.y or i.x, col) end end
218   return i end
219
220 function data(i,row)
221   push(i.all, row)
222   for _,col in pairs(i.cols) do add(col, row[col.at]) end
223   return i end
224
225 function file2Egs(file, i)
226   for row in file2things(file) do
227     if i then data(i,row) else i = header(row) end end
228   return i end
229 ---
230 --- SUMMARIZE
231 ---
232
233 local div,mid,mids,seen
234 function mid(i)
235   return nump(i) and i.mu or i.mode end
236
237 function div(i)
238   if nump(i) then return i.sd end
239   e=0
240   map(i.all,function(n) e = e+ n/i.n * math.log(n/i.n,2) end)
241   return -e end
242
243 function mids(cols,rows, seen,out)
244   seen = function(col) return nump(col) and Num(col.at) or Sym(col.at) end
245   out = map(cols, seen)
246   for _,row in pairs(rows) do
247     for _,seen in pairs(out) do
248       add(seen, row[seen.at]) end end
249   return rnds(map(out, function(seen) return mid(seen) end)) end
250 ---
251 --- DISTANCE
252 ---
253
254 local dist,far,furthest,neighbors
255 function dist(i,row1,row2, d,n,norm,dist1,lo,hi)
256   function norm(x,lo,hi)
257     return hi-lo<1E-9 and 0 or (x-lo)/(hi-lo)
258   end
259   function dist1(col,a,b)
260     if a=="?" and b=="?" then return 1 end
261     if not nump(col) then return a==b and 0 or 1 end
262     lo,hi=col.lo, col.hi
263     if a=="?" then b=norm(b,lo,hi); a=b<.5 and 1 or 0
264     elseif b=="?" then a=norm(a,lo,hi); b=a<.5 and 1 or 0
265     else a,b = norm(a,lo,hi), norm(b,lo,hi) end
266     return math.abs(a - b)
267   end
268   d,n = 0,0
269   for _,col in pairs(i.x) do
270     d = d + dist1(col, row1[col.at], row2[col.at])^the.p
271     n = n + 1 end
272   return (d/n)^(1/the.p) end
273
274 function far(      i,r1,rows,far)
275   return per(neighbors(i,r1,rows),far or the.far) [2] end
276
277 function furthest( i,r1,rows)
278   return last(neighbors(i,r1,rows)) [2] end
279
280 function neighbors(i,r1,rows)
281   return sort(map(rows, function(r2) return {dist(i,r1,r2),r2} end),firsts) end

```

```

282 ---
283 --- CLUSTER
284 ---
285
286 local half, cluster, clusters
287 function half(i, rows,      project,row,some,east,west,easts,wests,c,mid)
288   function project(row,a,b)
289     a= dist(i,east,row)
290     b= dist(i,west,row)
291     return ((a^2 + c^2 - b^2)/(2*c), row)
292   end
293   some = many(rows,the.some)
294   east = furthest(i,any(some), some)
295   west = furthest(i,east,      some)
296   c = dist(i,east,west)
297   easts,wests = {},{}
298   for n, xrow in pairs(sort(map(rows,project),firsts)) do
299     row = xrow[2]
300     if n==#rows//2 then mid=row end
301     push(n <= #rows//2 and easts or wests, row) end
302   return easts, wests, east, west, mid end
303
304 function cluster(i,rows, here,lefts,rights)
305   rows = rows or i.all
306   here = {all=rows}
307   if #rows > 2*(#i.all)^the.leaves then
308     lefts, rights = half(i, rows)
309     if #lefts < #rows then
310       here.lefts = cluster(i,lefts)
311       here.rights= cluster(i,rights) end end
312   return here end
313
314 function clusters(i,t,pre)
315   if t then
316     pre = pre or ""
317     if not t.lefts and not t.rights then
318       print(fmt("%$s%-20s",#t.all, pre), o(mids(i.y,t.all)))
319     else
320       print(fmt("%$s%-20s",#t.all, pre))
321       clusters(i,t.lefts, ". ". pre)
322       clusters(i,t.rights, ". ". pre) end end end
323 ---
324 --- DISCRETIZE
325 ---
326
327 local sym_spans, num_spans, merge, merged
328 function sym_spans(i, j)
329   local xys,all,one,last,x,y,n = {}, {}
330   for x,n in pairs(i.all) do push(xys, {x,"this",n}) end
331   for x,n in pairs(j.all) do push(xys, {x,"that",n}) end
332   for _,tmp in ipairs(sort(xys,firsts)) do
333     x,y,n = unpack(tmp)
334     if x ~= last then
335       last = x
336       one = push(all, {lo=x, hi=x, all=Num(i.at,i.txt)}) end
337     add(one.all, y, n) end
338   return all end
339
340 function num_spans(i, j)
341   local xys,all,lo,hi,gap,one,x,y,n = {},{}
342   lo,hi = math.min(i.lo, j.lo), math.max(i.hi,j.hi)
343   gap = (hi - lo) / (6/the.cohen)
344   for _,n in pairs(i.all) do push(xys, {n,"this",1}) end
345   for _,n in pairs(j.all) do push(xys, {n,"that",1}) end
346   one = {lo=lo, hi=lo, all=Sym(i.at,i.txt)}
347   all = {one}
348   for _,tmp in ipairs(sort(xys,firsts)) do
349     x,y,n = unpack(tmp)
350     if one.hi - one.lo > gap
351     then one = push(all, {lo=one.hi, hi=x, all=Sym(i.at,i.txt)}) end
352     one.hi = x
353     add(one.all,y,n) end
354   all
355   all[1].lo = -big
356   all[#all].hi = big
357   return all end
358
359 function merge(b4,      j,n,now,a,b,both)
360   j, n, now = 0, #b4, {}
361   while j < #b4 do
362     j = j+1
363     a, b = b4[j], b4[j+1]
364     if b then
365       both = merged(a,b)
366       if both then a, j = {lo=a.lo, hi=b.hi, all=both}, j+1 end end
367     push(now,a)
368     j = j+1 end
369   return #now == #b4 and b4 or merge(now) end
370
371 function merged(i,j,      k,ei,ej,ek)
372   k = Sym(i.at,i.txt)
373   for x,n in pairs(i.all) do add(k,x,n) end
374   for x,n in pairs(j.all) do add(k,x,n) end
375   ei, ej, ek= div(i), div(j), div(k)
376   if i.n==0 or j.n==0 or 1.01*ek <= (i.n*ei + j.n*ej)/(i.n+j.n) then
377     return k end end

```

```

378 -----
379
380 function Demo.the() oo(the) end
381
382 function Demo.many(a)
383   a={1,2,3,4,5,6,7,8,9,10}; ok("{1023}" == o(many(a,3)), "manys") end
384
385 function Demo.egs()
386   ok(5140==file2Egs(the.file).y[1].hi,"reading") end
387
388 function Demo.dist(i)
389   i = file2Egs(the.file)
390   for n,row in pairs(i.all) do print(n,dist(i, i.all[1], row)) end end
391
392 function Demo.far( i,j,row1,row2,row3,d3,d9)
393   i = file2Egs(the.file)
394   for j=1,10 do
395     row1 = any(i.all)
396     row2 = far(i,row1, i.all, .9)
397     d9 = dist(i,row1,row2)
398     row3 = far(i,row1, i.all, .3)
399     d3 = dist(i,row1,row3)
400     ok(d3 < d9, "closer far") end end
401
402 function Demo.half( i,easts,wests)
403   i = file2Egs(the.file)
404   easts,wests = half(i, i.all)
405   oo(mids(i.y, easts))
406   oo(mids(i.y, wests)) end
407
408 function Demo.cluster( i)
409   i = file2Egs(the.file)
410   i = file2Egs(the.file)
411   clusters(i,cluster(i))
412 end
413
414 -----
415 the=settings(help)
416 Demo.main(the.todo, the.seed)

```