

```

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 (inference):
9 -cohen -c F cohen's small effect size = .35
10 -far -F F look no further than "far" = .9
11 -keep -k items to keep in a number = 512
12 -leaves -l leaf size = .5
13 -p -p P distance calcs coefficient = 2
14 -seed -S P random number seed = 10019
15 -some -s look only at "some" items = 512
16
17 OPTIONS (housekeeping):
18 -dump -d on error, exit+ stacktrace = false
19 -file -f S where to get data = ../etc/data/auto93.csv
20 -help -h show help = false
21 -rnd -r S format string = %5.2f
22 -todo -t S start-up action = nothing
23 ]]]-----
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

objects

```

local as = setmetatable
local function obj( t)
  t=({__tostring=o; t.__index=t
  return as(t, {__call=function( _,...) return t.new( _,...) end}) end
end

local Sym = obj() -- Where to summarize symbols
function Sym:new(at,s) return as({
  is="Sym", -- type
  at=at or 0, -- column index
  name=s or "", -- column name
  n=0, -- number of items summarized in this column
  all={}, -- all[x] = n means we've seen "n" repeats of "x"
  most=0, -- count of the most frequently seen symbol
  mode=nil -- the most commonly seen letter
}, Sym) end

local Num = obj() -- Where to summarize numbers
function Num:new(at,s) return as({
  is="Num", -- type
  at=at or 0, -- column index
  name=s or "", -- column name
  n=0, -- number of items summarizes in this column
  mu=0, -- mean (updated incrementally)
  m2=0, -- second moment (updated incrementally)
  sd=0, -- standard deviation
  all={}, -- a sample of items seen so far
  lo=1E31, -- lowest number seen; initially, big so 1st num sends it low
  hi=-1E31, -- highest number seen; initially, small so 2st num sends it hi
  w=(s or ""):find"-$" and -1 or 1 -- "1"= minimize and "1"= maximize
}, Num) end

local Egs = obj() -- Where to store examples, summarized into Syms or Nums
function Egs:new(names, i,col,here) i=as({
  is="Egs", -- type
  all={}, -- all the rows
  names=names, -- list of name
  cols={}, -- list of all columns (Nums or Syms)
  x={}, -- independent columns (nothing marked as "skip")
  y={}, -- dependent columns (nothing marked as "skip")
}, Egs)
for at,name in pairs(names) do
  col = (name:find"^[A-Z]" and Num or Sym) (at,name)
  i.cols[1+#i.cols] = col
  here = name:find"^[+]" and i.y or i.x
  if not name:find"$" then here[1 + #here] = col end end
return i end

--[[
-- Coding Conventions
- "i" not "self"
- if something holds a list of thing, name the holding variable "all"
- no inheritance
- only define a method if that is for polymorphism
- when you can, write functions down on one line
- all config items into a global "the" variable
- all the test cases (or demos) are "function Demo.xxx".
- random seed reset so carefully, just once, at the end of the code.
- usually, no line with just "end" on it
]]

```

```

101
102
103
104
105
106
107
108
109
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
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199

```

tools

COERCE

```

local r = math.random
local fmt = string.format
local unpack = table.unpack
local function push(t,x) table.insert(t,x); return x end
return tonumber(x) or x end

function things(x,sep, t)
  t={}; for y in x:gmatch(sep or "[^+]+") do t[1+#t]=thing(y) end
  return t end

function file2things(file, x)
  file = io.input(file)
  return function ()
    x=io.read();
    if x then return things(x) else io.close(file) end end end

-- GET, SET
local last,per,any,many
function last(a) return a[ #a ] end
function per(a,p) return a[ (p*#a)//1 ] end
function any(a) return a[ math.random(#a) ] end
function many(a,n, u) u={}; for j=1,n do push(u,any(a)) end; return u end

-- LIST
local firsts,sort,map,slots
function firsts(a,b) return a[1] < b[1] end
function sort(t,f) table.sort(t,f); return t end
function map(t,f, u) u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
function slots(t, u,s)
  u={}
  for k,v in pairs(t) do s=tostring(k);if s:sub(1,1)~="_" then push(u,k) end end
  return sort(u) end

-- PRINT
local oo,o, rnd, rnds
function oo(t) print(o(t)) end
function o(t,seen, key,xseen,u)
  seen = seen or {}
  if type(t)~="table" then return tostring(t) end
  if seen[t] then return "..." end
  seen[t] = t
  key = function(k) return fmt(":%s %s",k,o(t[k],seen)) end
  xseen = function(x) return o(x,seen) end
  u = #t>0 and map(t,xseen) or map(slots(t),key)
  return (t.is or "")..'{'..table.concat(u," ")..."}'" end

function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
function rnd(x,f)
  return fmt(type(x)=="number" and (x~x//1 and f or the.rnd) or "%s",x) end

-- START-UP
local Demo, ok = {fails=0}
function ok(test,msg)
  print(test and "PASS:" or "FAIL:",msg or "")
  if not test then
    Demo.fails=Demo.fails+1
    if the.dump then assert(test,msg) end end end

function Demo.main(todo,seed)
  for k,one in pairs(todo== "all" and slots(Demo) or {todo}) do
    if k ~= "main" and type(Demo[one]) == "function" then
      math.randomseed(seed)
      Demo[one]() end end
  for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
  return Demo.fails end

local function settings(txt, d)
  d={}
  txt:gsub("(^[^%s+])(%s+)(-[^%s+])[^%s+]".."%s(%s+)",
  function(long,key,short,x)
    for n,flag in ipairs(arg) do
      if flag==short or flag==long then
        x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
        if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
          d[key] = tonumber(x) or x end end
  if the.help then print(txt) end
  return d end

```

```

200 ---
201 --- UPDATE COLS
202 ---
203 ---
204 local add
205 function add(i,x, inc)
206   inc = inc or 1
207   if x ~= "?" then
208     i.n = i.n + inc
209     i:internalAdd(x,inc) end
210   return x end
211
212 function Sym.internalAdd(i,x,inc)
213   i.all[x] = inc + (i.all[x] or 0)
214   if i.all[x] > i.most then i.most = i.all[x], x end end
215
216 function Num.internalAdd(i,x,inc, d)
217   for j=1,inc do
218     d = x - i.mu
219     i.mu = i.mu + d/i.n
220     i.m2 = i.m2 + d*(x - i.mu)
221     i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n-1))^0.5)
222     i.lo = math.min(x, i.lo)
223     i.hi = math.max(x, i.hi)
224     if #i.all < the.keep then push(i.all,x)
225     elseif r() < they.keep/i.n then i.all[r(#i.all)]=x end end end
226
227 ---
228 --- MAKE DATA
229 ---
230 local file2Egs -- not "local data" (since defined above)
231 function data(i,row)
232   push(i.all, row)
233   for _,col in pairs(i.cols) do add(col, row[col.at]) end
234   return i end
235
236 function file2Egs(file, i)
237   for row in file2things(file) do
238     if i then data(i,row) else i = Egs(row) end end
239   return i end
240
241 ---
242 --- SUMMARIZE
243 ---
244 local mids
245 function mids(i,rows,cols) return i:clone(rows):mid(cols) end
246
247 function Egs.mid(i,cols)
248   return map(cols or i.y,function(col) return col:mid() end) end
249
250 function Sym.mid(i) return i.mode end
251 function Num.mid(i) return i.mu end
252
253 function Num.div(i) return i.sd end
254 function Sym.div(i, e)
255   e=0, for _,n in pairs(i.all) do e=e + n/i.n*math.log(n/i.n,2) end
256   return -e end
257
258 ---
259 --- DISTANCE
260 ---
261 local far,furthest,neighbors,dist
262 function far( i,r1,rows,far)
263   return per(neighbors(i,r1,rows),far or the.far)[2] end
264
265 function furthest( i,r1,rows)
266   return last(neighbors(i,r1,rows))[2] end
267
268 function neighbors(i,r1,rows)
269   return sort(map(rows, function(r2) return {dist(i,r1,r2),r2} end),firsts) end
270
271 function dist(i,row1,row2, d,n,a,b,inc)
272   d,n = 0,0
273   for _,col in pairs(i.x) do
274     a,b = row1[col.at], row2[col.at]
275     inc = a=="?" and b=="?" and 1 or col:dist1(a,b)
276     d = d + inc^the.p
277     n = n + 1 end
278   return (d/n)^(1/the.p) end
279
280 function Sym.dist1(i,a,b) return a==b and 0 or 1 end
281
282 function Num.dist1(i,a,b)
283   if a=="?" then b:=norm(b); a=b<.5 and 1 or 0
284   elseif b=="?" then a:=norm(a); b=a<.5 and 1 or 0
285   else a,b = i:norm(a), i:norm(b) end
286   return math.abs(a - b) end
287
288 function Num.norm(i,x)
289   return i.hi - i.lo < 1E-32 and 0 or (x - i.lo)/(i.hi - i.lo) end
290
291 ---
292 --- CLUSTER
293 ---
294 local half, cluster, clusters
295 function half(i, rows, project,row,some,left,right,lefts,rights,c,mid)
296   function project(row,a,b)
297     a = dist(i,left,row)
298     b = dist(i,right,row)
299     return {(a^2 + c^2 - b^2)/(2*c), row}
300   end
301   some = many(rows, the.some)
302   left = furthest(i,any(some), some)
303   right = furthest(i,left, some)
304   c = dist(i,left,right)
305   lefts,rights = {},{}
306   for n, projection in pairs(sort(map(rows,project),firsts)) do
307     if n==#rows//2 then mid=row end
308     push(n <= #rows//2 and lefts or rights, projection[2]) end
309   return lefts, rights, left, right, mid, c end
310
311 function cluster(i,rows, here,lefts,rights)
312   rows = rows or i.all
313   here = {all=rows}
314   if #rows >= 2* (#i.all)^the.leaves then
315     lefts, rights, here.left, here.right, here.mid = half(i, rows)
316     if #lefts < #rows then
317       here.lefts = cluster(i,lefts)
318       here.rights = cluster(i,rights) end end
319   return here end
320
321 function clusters(i,format,t,pre, front)
322   if t then
323     pre=pre or ""
324     front = fmt("%s%s",pre,#t.all)
325     if not t.lefts and not t.rights then
326       print(fmt("%%-20s%s",front, o(rnds(mids(i,t.all),format))))
327     else
328       print(front)
329       clusters(i,format,t.lefts,"|.. pre)
330       clusters(i,format,t.rights,"|.. pre) end end end

```

```

331 ---
332 --- DISCRETIZE
333 ---
334 ---
335 local merge,merged,spans,bestSpan
336 function Sym.spans(i, j)
337   local xys,all,one,last,x,y,n = {},{}
338   for x,n in pairs(i.all) do push(xys, {x,"lefts",n}) end
339   for x,n in pairs(j.all) do push(xys, {x,"rights",n}) end
340   for _,tmp in ipairs(sort(xys,firsts)) do
341     x,y,n = unpack(tmp)
342     if x ~= last then
343       last = x
344       one = push(all, {lo=x, hi=x, all=Sym(i.at,i.name)}) end
345     add(one.all, y, n) end
346   return all end
347
348 function Num.spans(i, j)
349   local xys,all,lo,hi,gap,one,x,y,n = {},{}
350   lo,hi = math.min(i.lo, j.lo), math.max(i.hi, j.hi)
351   gap = (hi - lo) / (6/the.coehen)
352   for _,n in pairs(i.all) do push(xys, {n,"lefts",1}) end
353   for _,n in pairs(j.all) do push(xys, {n,"rights",1}) end
354   one = {lo=lo, hi=lo, all=Sym(i.at,i.name)}
355   all = {one}
356   for _,tmp in ipairs(sort(xys,firsts)) do
357     x,y,n = unpack(tmp)
358     if one.hi - one.lo > gap
359     then one = push(all, {lo=one.hi, hi=x, all=one.all:clone()}) end
360     one.hi = x
361     add(one.all, y, n) end
362   all = merge(all)
363   all[1].j.lo = -math.huge
364   all[#all].hi = math.huge
365   return all end
366
367 function merge(b4, j,n,now,a,b,both)
368   j, n, now = 0, #b4, {}
369   while j < #b4 do
370     j = j+1
371     a, b = b4[j], b4[j+1]
372     if b then
373       both = a.all:merge(b.all)
374       if both
375       then a = {lo=a.lo, hi=b.hi, all=both}
376       j = j + 1 end end
377     push(now,a) end
378   return #now == #b4 and b4 or merge(now) end
379
380 function Sym.merge(i,j, k,ei,ej,ek)
381   k = i:clone()
382   for x,n in pairs(i.all) do add(k,x,n) end
383   for x,n in pairs(j.all) do add(k,x,n) end
384   ei, ej, ek = i:div(), j:div(), k:div()
385   if ek*.99 <= (i.n*ei + j.n*ej)/k.n then return k end end
386
387 function spans(egs1,egs2, spans,tmp,coll,col2)
388   spans = {}
389   for c,coll in pairs(egs1.x) do
390     col2 = egs2.x[c]
391     tmp = coll:spans(col2)
392     if #tmp> 1 then
393       for _,one in pairs(tmp) do push(spans,one) end end end
394   return spans end
395
396 function bestSpan(spans)
397   local divs,ns,n,div,stats,dist2heaven = Num(), Num()
398   function dist2heaven(s) return {(1 - n(s))^2 + (0 - div(s))^2^.5,s} end
399   function div(s) return divs:norm( s.all:div() ) end
400   function n(s) return ns:norm( s.all.n ) end
401   for _,s in pairs(spans) do
402     add(divs, s.all:div())
403     add(ns, s.all.n) end
404   return sort(map(spans, dist2heaven), firsts)[1][2] end
405
406 ---
407 --- EXPLAIN
408 ---
409 local xplain,xplains,selects,spanShow
410 function xplain(i,rows,used,
411   stop,here,left,right,lefts0,rights0,lefts1,rights1)
412   used=used or {}
413   rows = rows or i.all
414   here = {all=rows}
415   stop = (#i.all)^the.leaves
416   if #rows >= 2*stop then
417     lefts0, rights0, here.left, here.right, here.mid, here.c = half(i, rows)
418     if #lefts0 < #rows then
419       here.selector = bestSpan(spans(i:clone(lefts0),i:clone(rights0)))
420       push(used, {here.selector.all.name, here.selector.lo, here.selector.hi})
421       lefts1,rights1 = {},{}
422       for _,row in pairs(rows) do
423         push(selects(here.selector, row) and lefts1 or rights1, row) end
424       if #lefts1 > stop then here.lefts = xplain(i,lefts1,used) end
425       if #rights1 > stop then here.rights = xplain(i,rights1,used) end end end
426   return here end
427
428 function xplains(i,format,t,pre,how, sel,front)
429   pre, how = pre or "", how or ""
430   if t then
431     pre=pre or ""
432     front = fmt("%s%s%s",pre,how, #t.all, t.c and rnd(t.c) or "")
433     if t.lefts and t.rights then print(fmt("%%-35s",front)) else
434       print(fmt("%%-35s",front, o(rnds(mids(i,t.all),format))))
435     end
436     sel = t.selector
437     xplains(i,format,t.lefts, "|.. pre, spanShow(sel,.."")
438     xplains(i,format,t.rights, "|.. pre, spanShow(sel,true) ..") end end
439
440 function selects(span,row, lo,hi,at,x)
441   lo, hi, at = span.lo, span.hi, span.all.at
442   x = row[at]
443   if x=="?" then return true end
444   if lo==hi then return x==lo else return lo <= x and x < hi end end
445
446 function spanShow(span, negative, hi,lo,x,big)
447   if not span then return "" end
448   lo, hi, x, big = span.lo, span.hi, span.all.name, math.huge
449   if not negative
450   then if hi == big then return fmt("%s== %s",x,lo) end
451   if hi == big then return fmt("%s>= %s",x,lo) end
452   if lo == -big then return fmt("%s<= %s",x,hi) end
453   return fmt("%s<= %s< %s",lo,x,hi)
454   else if lo == hi then return fmt("%s!= %s",x,lo) end
455   if hi == big then return fmt("%s< %s",x,lo) end
456   if lo == -big then return fmt("%s>= %s",x,hi) end
457   return fmt("%s< %s and %s>= %s", x,lo,x,hi) end end

```

```

458 ----
459 ----
460 ----
461 ----
462 ----
463
464 function Demo.the() oo(the) end
465
466 function Demo.many(a)
467   a={1,2,3,4,5,6,7,8,9,10}; ok("{1023}" == o(many(a,3)), "manys") end
468
469 function Demo.egs()
470   ok(5140==file2Egs(the.file).y[1].hi,"reading") end
471
472 function Demo.dist(i)
473   i = file2Egs(the.file)
474   for n,row in pairs(i.all) do print(n,dist(i, i.all[1], row)) end end
475
476 function Demo.far( i,j,row1,row2,row3,d3,d9)
477   i = file2Egs(the.file)
478   for j=1,10 do
479     row1 = any(i.all)
480     row2 = far(i,row1, i.all, .9)
481     d9 = dist(i,row1,row2)
482     row3 = far(i,row1, i.all, .3)
483     d3 = dist(i,row1,row3)
484     ok(d3 < d9, "closer far") end end
485
486 function Demo.half( i, lefts, rights)
487   i = file2Egs(the.file)
488   lefts, rights = half(i, i.all)
489   oo(mids(i, lefts))
490   oo(mids(i, rights))
491   end
492
493 function Demo.cluster( i)
494   i = file2Egs(the.file)
495   clusters(i, "%0f", cluster(i)) end
496
497 function Demo.spans( i, lefts, rights)
498   i = file2Egs(the.file)
499   lefts, rights = half(i, i.all)
500   oo(bestSpan(spans(i:clone(lefts), i:clone(rights)))) end
501
502 function Demo.xplain( i,j,tmp,lefts,rights,used)
503   i = file2Egs(the.file)
504   used={}
505   xplains(i, "%0f", xplain(i, i.all, used))
506   map(sort(used, function(a,b)
507     return ((a[1] < b[1]) or
508            (a[1]==b[1] and a[2] < b[2]) or
509            (a[1]==b[1] and a[2]==b[2] and a[3] < b[3]))end), oo) end
510
511
512 -----
513 the = settings(help)
514 Demo.main(the.todo, the.seed)

```