```lua
1  local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
2  local the,help={},[[
3
4  lua l5.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
27  KEY: S=string, P=poisint, F=float
28  ]]
29
30  local as = setmetatable
31  local function obj(   t)
32    t={__tostring=o}; t.__index=t
33    return as(t, {__call=function(_,...) return t.new(_,...) end}) end
34  ------------------------------------------------------------------------------
35  ---
36  ---       _|      _|_|    _|_|_|_|_|    _|_|
37  ---     _|_|    _|    _|      _|      _|    _|
38  ---
39  ---
40
41  local Sym, Num = obj(), obj()
42  function Sym:new(at,s) return as({
43    is="Sym",        -- type
44    at=at or 0,      -- column index
45    name=s or "",    -- column name
46    n=0,             -- number of items summarized in this column
47    all={},          -- all[x] = n means we've seen "n" repeats of "x"
48    most=0,          -- count of the most frequently seen symbol
49    mode=nil         -- the most commonly seen letter
50    }, Sym) end
51
52  function Num:new(at,s) return as({
53    is="Num",        -- type
54    at=at or 0,      -- column index
55    name=s or "",    -- column name
56    n=0,             -- number of items summarizes in this column
57    mu=0,            -- mean (updated incrementally)
58    m2=0,            -- second moment (updated incrementally)
59    sd=0,            -- standard deviation
60    all={},          -- a sample of items seen so far
61    lo=1E31,         -- lowest number seen
62    hi=-1E31,        -- highest number seen
63    w=(s or ""):find"-$" and -1 or 1 -- "-1"= minimize and "1"= maximize
64    }, Num) end
65
66  local function Egs(names)  return {
67    is="egs",        -- type
68    all={},          -- all the rows
69    names=names,     -- list of name
70    cols={},         -- list of all columns  (Nums or Syms)
71    x={},            -- independent columns (nothing marked as "skip")
72    y={}             -- dependent columns (nothing marked as "skip")
73    } end
74
75  --[[
76  ## Coding Conventions
77  - "i" not "self"
78  - if something holds a list of thing, name the holding variable "all"
79  - no inheritance
80  - only define a method if that is for polymorphism
81  - when you can, write functions down on one line
82  - all config items into a global "the" variable
83  - all the test cases (or demos) are "function Demo.xxx".
84  - random seed reset so carefully, just once, at the end of the code.
85  ]]
```

```lua
86  ---
87  ---
88  ---       _|    _|  _|_|_|_|_|  _|  _|        _|_|_|
89  ---       _|    _|      _|      _|  _|      _|
90  ---         _|_|        _|      _|  _|_|_|    _|_|_|
91  ------------------------------------------------------------------------------
92  local r    = math.random
93  local fmt  = string.format
94  local function push(t,x) table.insert(t,x); return x end
95  ---
96  ---         _|_|_|    _|_|    _|_|_|    _|_|_|    _|_|
97  ---       _|        _|    _|  _|    _|  _|        _|
98
99  local thing,things,file2things
100 function thing(x)
101   x = x:match"^%s*(.-)%s*$"
102   if x=="true" then return true elseif x=="false" then return false end
103   return tonumber(x) or x end
104
105 function things(x,sep,   t)
106   t={}; for y in x:gmatch(sep or"([^,]+)") do push(t,thing(y)) end
107   return t end
108
109 function file2things(file,      x)
110   file = io.input(file)
111   return function()
112     x=io.read();
113     if x then return things(x) else io.close(file) end end end
114 ---
115 ---         _|_|_|  _|_|_|_|  _|_|_|_|_|        _|_|_|  _|_|_|_|  _|_|_|_|_|
116 ---       _|        _|            _|          _|        _|            _|
117 ---                                      '
118 local last,per,any,many
119 function last(a)          return a[ #a ] end
120 function per(a,p)         return a[ (p*#a)//1 ] end
121 function any(a)           return a[ math.random(#a) ] end
122 function many(a,n,   u) u={}; for j=1,n do push(u,any(a)) end; return u end
123 ---
124 ---         _|_|_|_|_|  _|_|_|    _|_|    _|_|_|    _|_|_|_|
125 ---
126
127 local firsts,sort,map,slots
128 function firsts(a,b)   return a[1] < b[1] end
129 function sort(t,f)     table.sort(t,f); return t end
130 function map(t,f,  u)  u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
131 function slots(t, u,s)
132   u={}
133   for k,v in pairs(t) do s=tostring(k);if s:sub(1,1)~="_" then push(u,k) end end
134   return sort(u) end
135 ---
136 ---         _|_|_|    _|_|_|    _|  _|      _|  _|_|_|_|_|
137 ---
138
139 local oo,o, rnd, rnds
140 function oo(t) print(o(t)) end
141 function o(t,seen,            key,xseen,u)
142   seen = seen or {}
143   if type(t)~="table" then return tostring(t) end
144   if seen[t]          then return "..." end
145   seen[t] = t
146   key  = function(k) return fmt(":%s %s",k,o(t[k],seen)) end
147   xseen = function(x) return o(x,seen) end
148   u = #t>0 and map(t,xseen) or map(slots(t),key)
149   return (t.is or "")..'{'..table.concat(u,"")..'}' end
150
151 function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
152 function rnd(x,f)
153   return fmt(type(x)=="number" and (x~=x//1 and f or the.rnd) or "%s",x) end
154 ---
155 ---         _|  _|_|_|_|_|    _|_|    _|_|_|    _|_|_|_|_|        _|    _|  _|_|_|
156 ---
157
158 local Demo, ok = {fails=0}
159 function ok(test,msg)
160   print(test and "PASS:"or "FAIL:",msg or "")
161   if not test then
162     Demo.fails=Demo.fails+1
163     if the.dump then assert(test,msg) end end end
164
165 function Demo.main(todo,seed)
166   for k,one in pairs(todo=="all" and slots(Demo) or {todo}) do
167     if k ~= "main" and type(Demo[one]) == "function" then
168       math.randomseed(seed)
169       Demo[one]() end end
170   for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
171   return Demo.fails
172
173 local function settings(txt,   d)
174   d={}
175   txt:gsub("\n ([-](%[^%s]+))[%s]+(-[-^%s]+)[^\n]*%s([^%s]+)",
176     function(long,key,short,x)
177       for n,flag in ipairs(arg) do
178         if flag==short or flag==long then
179           x = x=="false" and true or x=="true" or arg[n+1] end end
180       if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
181       d[key] = tonumber(x) or x end end)
182   if d.help then print(help) end
183   return d end
```

```lua
184  ---
185  ---    UPDATE COLS
186  ---
187
188  local add
189  function add(i,x, inc)
190    inc = inc or 1
191    if x ~= "?" then
192      i.n = i.n + inc
193      i:add1(x,inc) end
194    return x end
195
196  function Sym.add1(i,x,inc)
197    i.all[x] = inc + (i.all[x] or 0)
198    if i.all[x] > i.most then i.most, i.mode = i.all[x], x end end
199
200  function Num.add1(i,x,inc,     d)
201    for j=1,inc do
202      d     = x - i.mu
203      i.mu  = i.mu + d/i.n
204      i.m2  = i.m2 + d*(x - i.mu)
205      i.sd  = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n-1))^0.5)
206      i.lo  = math.min(x, i.lo)
207      i.hi  = math.max(x, i.hi)
208      if      #i.all < the.keep       then push(i.all,x)
209      elseif r()    < they.keep/i.n then i.all[r(#i.all)]=x end end end
```

```lua
210  ---
211  ---    MAKE DATA
212  ---
213
214  local header,data,file2Egs
215  function header(names,    i,col)
216    i = Egs(names)
217    for at,name in pairs(names) do
218      col = push(i.cols, (name:find"^[A-Z]" and Num or Sym)(at,name))
219      if not name:find":$" then
220        push(name:find"[-+]$" and i.y or i.x, col) end end
221    return i end
222
223  function data(i,row)
224    push(i.all, row)
225    for _,col in pairs(i.cols) do add(col, row[col.at]) end
226    return i end
227
228  function file2Egs(file,    i)
229    for row in file2things(file) do
230      if i then data(i,row) else i = header(row) end end
231    return i end
```

```lua
232  ---
233  ---    SUMMARIZE
234  ---
235
236  function Sym.mid(i) return i.mode end
237  function Sym.div(i,   e)
238    e=0; map(i.all,function(n) e = e + n/i.n * math.log(n/i.n,2) end)
239    return -e end
240
241  function Num.mid(i) return i.mu end
242  function Num.div(i) return i.sd end
243
244  function Num.clone(i) return Num(i.at, i.name) end
245  function Sym.clone(i) return Sym(i.at, i.name) end
246
247  local mids
248  function mids(cols,rows,     seen,tmp)
249    seen = function(col) return col:clone() end
250    tmp  = map(cols, seen)
251    for _,row in pairs(rows) do
252      for _,seen in pairs(tmp) do
253        add(seen, row[seen.at]) end end
254    return rnds(map(tmp, function(seen) return seen:mid() end)) end
```

```lua
255  ---
256  ---    DISTANCE
257  ---
258
259  local far,furthest,neighbors,dist
260  function far(      i,r1,rows,far)
261    return per(neighbors(i,r1,rows),far or the.far)[2] end
262
263  function furthest( i,r1,rows)
264    return last(neighbors(i,r1,rows))[2] end
265
266  function neighbors(i,r1,rows)
267    return sort(map(rows, function(r2) return {dist(i,r1,r2),r2} end),firsts) end
268
269  function dist(i,row1,row2,     d,n,a,b,inc)
270    d,n = 0,0
271    for _,col in pairs(i.x) do
272      a,b = row1[col.at], row2[col.at]
273      inc = a=="?" and b=="?" and 1 or col:dist1(a,b)
274      d = d + inc^the.p
275      n = n + 1 end
276    return (d/n)^(1/the.p) end
277
278  function Sym.dist1(i,a,b) return a==b and 0 or 1 end
279
280  function Num.dist1(i,a,b)
281    if      a=="?" then b=i:norm(b); a=b<.5 and 1 or 0
282    elseif b=="?" then a=i:norm(a); b=a<.5 and 1 or 0
283    else    a,b = i:norm(a), i:norm(b)   end
284    return math.abs(a - b) end
285
286  function Num.norm(i,x)
287    return i.hi - i.lo < 1E-32 and 0 or (x - i.lo)/(i.hi - i.lo) end
```

```lua
288  ---
289  ---    CLUSTER
290  ---
291
292  local half, cluster, clusters
293  function half(i, rows,     project,row,some,east,west,easts,wests,c,mid)
294    function project(row,a,b)
295      a= dist(i,east,row)
296      b= dist(i,west,row)
297      return {(a^2 + c^2 - b^2)/(2*c), row}
298    end ----------------------
299    some = many(rows, the.some)
300    east = furthest(i,any(some), some)
301    west = furthest(i,east,          some)
302    c    = dist(i,east,west)
303    easts,wests = {},{}
304    for n, xrow in pairs(sort(map(rows,project),firsts)) do
305      row = xrow[2]
306      if n==#rows//2 then mid=row end
307      push(n <= #rows//2 and easts or wests, row) end
308    return easts, wests, east, west, mid   end
309
310  function cluster(i,rows,  here,lefts,rights)
311    rows = rows or i.all
312    here = {all=rows}
313    if #rows > 2*(#i.all)^the.leaves then
314      lefts, rights = half(i, rows)
315      if #lefts < #rows then
316        here.lefts = cluster(i,lefts)
317        here.rights= cluster(i,rights) end end
318    return here end
319
```

```lua
320  function clusters(i,t,pre)
321    if t then
322      pre = pre or ""
323      if not t.lefts and not t.rights then
324        print(fmt("%5s %-20s",#t.all, pre), o(mids(i.y,t.all)))
325      else
326        print(fmt("%5s %-20s",#t.all, pre))
327        clusters(i,t.lefts,  "|.. ".. pre)
328        clusters(i,t.rights, "|.. ".. pre) end end end
```

```lua
329  ---
330  ---    DISCRETIZE
331  ---
332
333  local merge,merged
334  function Sym.spans(i, j)
335    local xys,all,one,last,x,y,n = {}, {}
336    for x,n in pairs(i.all) do push(xys, {x,"easts",n}) end
337    for x,n in pairs(j.all) do push(xys, {x,"wests",n}) end
338    for _,tmp in ipairs(sort(xys,firsts)) do
339      x,y,n = unpack(tmp)
340      if x ~= last then
341        last = x
342        one  = push(all, {lo=x, hi=x, all=Num(i.at,i.txt)}) end
343      add(one.all, y, n) end
344    return all end
345
346  function Num.spans(i, j)
347    local xys,all,lo,hi,gap,one,x,y,n = {},{}
348    lo,hi = math.min(i.lo, j.lo), math.max(i.hi,j.hi)
349    gap   = (hi - lo) / (6/the.cohen)
350    for _,n in pairs(i.all) do push(xys, {n,"easts",1}) end
351    for _,n in pairs(j.all) do push(xys, {n,"wests",1}) end
352    one = {lo=lo, hi=lo, all=Sym(i.at,i.txt)}
353    all = {one}
354    for _,tmp in ipairs(sort(xys,firsts)) do
355      x,y,n = unpack(tmp)
356      if    one.hi - one.lo > gap
357      then one = push(all, {lo=one.hi, hi=x, all=Sym(i.at,i.txt)}) end
358      one.hi = x
359      add(one.all,y,n) end
360    all          = merge(all)
361    all[1   ].lo = -math.huge
362    all[#all].hi =  math.huge
363    return all end
364
365  function merge(b4,       j,n,now,a,b,both)
366    j, n, now = 0, #b4, {}
367    while j < #b4 do
368      j    = j+1
369      a, b = b4[j], b4[j+1]
370      if b then
371        both = merged(a,b)
372        if both then a, j = {lo=a.lo, hi=b.hi, all=both}, j+1 end end
373      push(now,a)
374      j = j+1 end
375    return #now == #b4 and b4 or merge(now) end
376
377  function merged(i,j,      k,ei,ej,ek)
378    k = Sym(i.at,i.txt)
379    for x,n in pairs(i.all) do add(k,x,n) end
380    for x,n in pairs(j.all) do add(k,x,n) end
381    ei, ej, ek= div(i), div(j), div(k)
382    if i.n==0 or j.n==0 or 1.01*ek <= (i.n*ei + j.n*ej)/(i.n+j.n) then
       return k end end
```

```lua
383
384 -----------------------------------------------------------------------------------
385 function Demo.the() oo(the) end
386
387 function Demo.many(a)
388   a={1,2,3,4,5,6,7,8,9,10}; ok("{10 2 3}" == o(many(a,3)), "manys") end
389
390 function Demo.egs()
391   ok(5140==file2Egs(the.file).y[1].hi,"reading") end
392
393 function Demo.dist(i)
394   i = file2Egs(the.file)
395   for n,row in pairs(i.all) do print(n,dist(i, i.all[1], row)) end end
396
397 function Demo.far(  i,j,row1,row2,row3,d3,d9)
398   i = file2Egs(the.file)
399   for j=1,10 do
400     row1 = any(i.all)
401     row2 = far(i,row1, i.all, .9)
402     d9   = dist(i,row1,row2)
403     row3 = far(i,row1, i.all, .3)
404     d3   = dist(i,row1,row3)
405     ok(d3 < d9, "closer far") end end
406
407 function Demo.half(  i,easts,wests)
408   i = file2Egs(the.file)
409   easts,wests = half(i, i.all)
410   oo(mids(i.y, easts))
411   oo(mids(i.y, wests)) end
412
413 function Demo.cluster(   i)
414   i = file2Egs(the.file)
415   i = file2Egs(the.file)
416   clusters(i,cluster(i))
417  end
418
419 -----------------------------------------------------------------------------------
420 the=settings(help)
421 Demo.main(the.todo, the.seed)
```