```lua
1    --------------------------------------------------------------------------------
2    ---      _/\_         _/\_                          ___/:
3    ---     /\  /\       /\  /\                        ,'___/  :
4    ---    /\ \/ /\     /\ \/ /\                     __/   /
5    ---    \ \   /\     \ \   /\                    /  /  /
6    ---     \ \L\ \     \/\ \L\ \                   \/  />/
7    ---      \___/       \____/                      /  /_\
8    ---      \/___/       \/___/                      )'-. /
9    ---                                               )'-. /
10   ---     a little LUA learning library           ./  :\
11   ---     (c) Tim Menzies 2022, BSD-2            /.'  '
12   ---     Share and enjoy.                      '/'
13   ---                                           +
14   ---                                         '.
15   ---                                    .-"-.
16   ---                                  (    |
17   ---                                 . .-' .' )
18   ---                               ( (.   )8:
19   ---                               .   ( `-' (  `.
20   ---                               . : (   .a8a)
21   ---                              /_`( "a `a. )"'
22   ---                              (  (/  .  ' )==`
23   ---                              (   (    )  .8"   +
24   ---
25   ---
26   ---
27   local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
28   local the,help={},[[
29
30   lua l5.lua [OPTIONS]
31   L5 == a very little LUA learning lab
32   (c)2022, Tim Menzies, BSD 2-clause license
33
34   OPTIONS (for changing the inference):
35
36    -cohen  -c  F  cohen's small effect size     = .35
37    -far    -F  F  look no further than "far"    = .9
38    -keep   -k     items to keep in a number     = 512
39    -leaves -l     leaf size                      = .5
40    -p      -p  P  distance calcs coefficient    = 2
41    -seed   -S  P  random number seed            = 10019
42    -some   -s     look only at "some" items     = 512
43
44   OPTIONS (for housekeeping):
45
46    -dump   -d     exit on error, with stacktrace = false
47    -file   -f  S  where to get data              = ../etc/data/auto93.csv
48    -help   -h     show help                      = false
49    -rnd    -r  S  format string                 = %5.2f
50    -todo   -t  S  start-up action               = nothing
51
52   KEY: S=string, P=posint, F=float
53   ]]
54   local as,o = setmetatable
55   local function obj(   t)
56     t={__tostring=o}; t.__index=t
57     return as(t, {__call=function(_,...) return t.new(_,...) end}) end
58   ---      _____      _
59   ---     |  _  \    | |
60   ---     | | | |__ _| |_ __ _
61   ---     | | | / _` | __/ _` |
62   ---     | |/ / (_| | || (_| |
63   ---     |___/ \__,_|\__\__,_|
64
65   local Sym = obj() -- Where to summarize symbols
66   function Sym:new(at,s) return as({
67     is="Sym",      -- type
68     at=at or 0,    -- column index
69     name=s or "",  -- column name
70     n=0,           -- number of items summarized in this column
71     all={},        -- all[x] = n means we've seen "n" repeats of "x"
72     most=0,        -- count of the most frequently seen symbol
73     mode=nil       -- the most commonly seen letter
74     }, Sym) end
75
76   local Num = obj() -- Where to summarize numbers
77   function Num:new(at,s) return as({
78     is="Num",      -- type
79     at=at or 0,    -- column index
80     name=s or "",  -- column name
81     n=0,           -- number of items summarizes in this column
82     mu=0,          -- mean (updated incrementally)
83     m2=0,          -- second moment (updated incrementally)
84     sd=0,          -- standard deviation
85     all={},        -- a sample of items seen so far
86     lo=1E31,       -- lowest number seen
87     hi=-1E31,      -- highest number seen
88     w=(s or ""):find"-$" and -1 or 1 -- "-1"= minimize and "1"= maximize
89     }, Num) end
90
91   local Egs = obj() -- Where to store examples, summarized into Syms or Nums
92   function Egs:new(names,      i,col,here)   i=as({
93     is="Egs",      -- type
94     all={},        -- all the rows
95     names=names,   -- list of name
96     cols={},       -- list of all columns   (Nums or Syms)
97     x={},          -- independent columns (nothing marked as "skip")
98     y={}           -- dependent columns (nothing marked as "skip")
99     },Egs)
100  for at,name in pairs(names) do
101    col = (name:find"^[A-Z]" and Num or Sym)(at,name)
102    i.cols[1+#i.cols] = col
103    here = name:find"[-+|$" and i.y or i.x
104    if not name:find":$" then here[1 + #here] = col end end
105  return i end
106  ---       _____ __
107  ---      /  ___|  |   ___  __
108  ---      \ `--.|  |  / _ \/  \
109  ---       `--. \  |  | | | | |\ \
110
111  function Num.clone(i) return Num(i.at, i.name) end
112  function Sym.clone(i) return Sym(i.at, i.name) end
113  function Egs.clone(i) return Egs(i.names)       end
114  --[[
115  ## Coding Conventions
116  - "i" not "self"
117  - if something holds a list of thing, name the holding variable "all"
118  - no inheritance
119  - only define a method if that is for polymorphism
120  - when you can, write functions down on one line
121  - all config items into a global "the" variable
122  - all the test cases (or demos) are "function Demo.xxx".
123  - random seed reset so carefully, just once, at the end of the code.
124  - usually, no line with just "end" on it
125  ]]

126  ---     _   _ _   _ _
127  ---    | | | | |_(_) |___
128  ---    | | | | __| | / __|
129  ---    | |_| | |_| | \__ \
130  ---     \___/ \__|_|_|___/
131  --------------------------------------------------------------------------------
132  local r   = math.random
133  local fmt = string.format
134  local function push(t,x) table.insert(t,x); return x end
135
136  ---       _____ _____ _____ _____ _____
137  ---      |     |     |     |     |
138
139  local thing,things,file2things
140  function thing(x)
141    x = x:match"^%s*(.-)%s*$"
142    if x=="true" then return true elseif x=="false" then return false end
143    return tonumber(x) or x end
144
145  function things(x,sep,   t)
146    t={}; for y in x:gmatch(sep or"([^,]+)") do push(t,thing(y)) end
147    return t end
148
149  function file2things(file,       x)
150    file = io.input(file)
151    return function()
152      x=io.read();
153      if x then return things(x) else io.close(file) end end end
154  ---       _____ _____ _____   _____ _____ _____
155  ---      |   __|   __|_   _| |   __|   __|_   _|
156  ---      |  |  |   __| | |   |__   |   __| | |
157  ---      |_____|_____| |_|   |_____|_____| |_|
158  local last,per,any,many
159  function last(a)         return a[ #a ] end
160  function per(a,p)        return a[ (p*#a)//1 ] end
161  function any(a)          return a[ math.random(#a) ] end
162  function many(a,n,   u) u={}; for j=1,n do push(u,any(a)) end; return u end
163  ---       _     _____ _____ _____
164  ---      | |   |     |   __|_   _|
165  ---      | |__ |-   -|__   | | |
166
167  local firsts,sort,map,slots
168  function firsts(a,b)   return a[1] < b[1] end
169  function sort(t,f)     table.sort(t,f); return t end
170  function map(t,f,  u)   u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
171  function slots(t,  u,s)
172    u={}
173    for k,v in pairs(t) do s=tostring(k);if s:sub(1,1)~="_" then push(u,k) end end
174    return sort(u) end
175  ---       _____ _____ _____ _____ _____
176  ---      |  _  | __  |     |   | |_   _|
177  ---      |   __|    -|-   -| | | | | |
178
179  local oo,rnd, rnds -- local o was declared above (in "new")
180  function oo(t)  print(o(t)) end
181  function o(t,seen,           key,xseen,u)
182    seen = seen or {}
183    if type(t)~="table" then return tostring(t) end
184    if seen[t]           then return "..." end
185    seen[t] = t
186    key    = function(k) return fmt(":%s %s",k,o(t[k],seen)) end
187    xseen  = function(x) return o(x,seen) end
188    u = #t>0 and map(t,xseen) or map(slots(t),key)
189    return (t.is or "")..'{'..table.concat(u,"")..'}' end
190
191  function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
192  function rnd(x,f)
193    return fmt(type(x)=="number" and (x~=x//1 and f or the.rnd) or "%s",x) end
194  ---       _____ _____ _____ _____ _____ _____ _____
195  ---      |   __|_   _|  _  | __  |_   _|___| | | | _ |
196  ---      |__   | | | |     |    -| | | |___| | | |  _|
197
198  local Demo, ok = {fails=0}
199  function ok(test,msg)
200    print(test and "PASS:"or "FAIL:",msg or "")
201    if not test then
202      Demo.fails=Demo.fails+1
203      if the.dump then assert(test,msg) end end end

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

213  local function settings(txt,   d)
214    d={}
215    txt:gsub("\n  ([-]([^%s]+))[%s]+(-[^%s]+)[^\n]*%s([^%s]+)",
216      function(long,key,short,x)
217        for n,flag in ipairs(arg) do
218          if flag==short or flag==long then
219            x = x=="false" and true or x=="true" or "false" or arg[n+1] end end
220          if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
221          d[key] = tonumber(x) or x end end)
222    if d.help then print(help) end
223    return d end
```

```lua
---
---     UPDATE COLS
---
228  local add
229  function add(i,x, inc)
230    inc = inc or 1
231    if x ~= "?" then
232      i.n = i.n + inc
233      i:add1(x,inc) end
234    return x end
235
236  function Sym.add1(i,x,inc)
237    i.all[x] = inc + (i.all[x] or 0)
238    if i.all[x] > i.most then i.most, i.mode = i.all[x], x end end
239
240  function Num.add1(i,x,inc,      d)
241    for j=1,inc do
242      d      = x - i.mu
243      i.mu   = i.mu + d/i.n
244      i.m2   = i.m2 + d*(x - i.mu)
245      i.sd   = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n-1))^0.5)
246      i.lo   = math.min(x, i.lo)
247      i.hi   = math.max(x, i.hi)
248      if     #i.all < the.keep       then push(i.all,x)
249      elseif r()    < they.keep/i.n then i.all[r(#i.all)]=x end end end
---
---     MAKE DATA
---
254  local data,file2Egs
255  function data(i,row)
256    push(i.all, row)
257    for _,col in pairs(i.cols) do add(col, row[col.at]) end
258    return i end
259
260  function file2Egs(file,    i)
261    for row in file2things(file) do
262      if i then data(i,row) else i = Egs(row) end end
263    return i end
---
---     SUMMARIZE
---
268  function Sym.mid(i) return i.mode end
269  function Num.mid(i) return i.mu end
270
271  function Num.div(i) return i.sd end
272  function Sym.div(i,   e)
273    e=0; for _,n in pairs(i.all) do e=e + n/i.n*math.log(n/i.n,2) end
274    return -e end
275
276  function Egs.mid(i,cols)
277    return map(cols or i.y,function(col) return col:mid() end) end
278
279  local mids
280  function mids(i,rows,cols,      seen,tmp)
281    j = i:clone()
282    for _,row in pairs(rows) do data(j, row) end
283    return rnds(j:mid(cols)) end
---
---     DISTANCE
---
288  local far,furthest,neighbors,dist
289  function far(       i,r1,rows,far)
290    return per(neighbors(i,r1,rows),far or the.far)[2] end
291
292  function furthest( i,r1,rows)
293    return last(neighbors(i,r1,rows))[2] end
294
295  function neighbors(i,r1,rows)
296    return sort(map(rows, function(r2) return {dist(i,r1,r2),r2} end),firsts) end
297
298  function dist(i,row1,row2,     d,n,a,b,inc)
299    d,n = 0,0
300    for _,col in pairs(i.x) do
301      a,b = row1[col.at], row2[col.at]
302      inc = a=="?" and b=="?" and 1 or col:dist1(a,b)
303      d = d + inc^the.p
304      n = n + 1 end
305    return (d/n)^(1/the.p) end
306
307  function Sym.dist1(i,a,b)  return a==b and 0 or 1 end
308
309  function Num.dist1(i,a,b)
310    if     a=="?" then b=i:norm(b); a=b<.5 and 1 or 0
311    elseif b=="?" then a=i:norm(a); b=a<.5 and 1 or 0
312    else   a,b = i:norm(a), i:norm(b)   end
313    return math.abs(a - b) end
314
315  function Num.norm(i,x)
316    return i.hi - i.lo < 1E-32 and 0 or (x - i.lo)/(i.hi - i.lo) end
---
---     CLUSTER
---
321  local half, cluster, clusters
322  function half(i, rows,     project,row,some,east,west,easts,wests,c,mid)
323    function project(row,a,b)
324      a= dist(i,east,row)
325      b= dist(i,west,row)
326      return {(a^2 + c^2 - b^2)/(2*c), row}
327    end -----------------------
328    some = many(rows, the.some)
329    east = furthest(i,any(some), some)
330    west = furthest(i,east,            some)
331    c    = dist(i,east,west)
332    easts,wests = {},{}
333    for n, xrow in pairs(sort(map(rows,project),firsts)) do
334      row = xrow[2]
335      if n==#rows//2 then mid=row end
336      push(n <= #rows//2 and easts or wests, row) end
337    return easts, wests, east, west, mid  end
338
339  function cluster(i,rows,  here,lefts,rights)
340    rows = rows or i.all
341    here = {all=rows}
342    if #rows > 2*(#i.all)^the.leaves then
343      lefts, rights = half(i, rows)
344      if #lefts < #rows then
345        here.lefts = cluster(i,lefts)
346        here.rights= cluster(i,rights) end end
347    return here end
348
349  function clusters(i,t,pre)
350    if t then
351      pre = pre or ""
352      if not t.lefts and not t.rights then
353        print(fmt("%5s %-20s",#t.all, pre), o(mids(i,t.all)))
354      else
355        print(fmt("%5s %-20s",#t.all, pre))
356        clusters(i,t.lefts,  "|.."..pre)
357        clusters(i,t.rights, "|.."..pre) end end end
---
---     [...]
---
---
---     [...]
---
362  function Sym.spans(i, j)
363    local xys,all,one,last,x,y,n = {}, {}
364    for x,n in pairs(i.all) do push(xys, {x,"easts",n}) end
365    for x,n in pairs(j.all) do push(xys, {x,"wests",n}) end
366    for _,tmp in ipairs(sort(xys,firsts)) do
367      x,y,n = unpack(tmp)
368      if x ~= last then
369        last = x
370        one  = push(all, {lo=x, hi=x, all=Num(i.at,i.txt)}) end
371      add(one.all, y, n) end
372    return all end
373
374  local merge,merged
375  function Num.spans(i, j)
376    local xys,all,lo,hi,gap,one,x,y,n = {},{}
377    lo,hi = math.min(i.lo, j.lo), math.max(i.hi,j.hi)
378    gap   = (hi - lo) / (6/the.cohen)
379    for _,n in pairs(i.all) do push(xys, {n,"easts",1}) end
380    for _,n in pairs(j.all) do push(xys, {n,"wests",1}) end
381    one = {lo=lo, hi=lo, all=Sym(i.at,i.txt)}
382    all = {one}
383    for _,tmp in ipairs(sort(xys,firsts)) do
384      x,y,n = unpack(tmp)
385      if    one.hi - one.lo > gap
386      then one = push(all, {lo=one.hi, hi=x, all=Sym(i.at,i.txt)}) end
387      one.hi = x
388      add(one.all,y,n) end
389    all         = merge(all)
390    all[1   ].lo = -math.huge
391    all[#all].hi =  math.huge
392    return all end
393
394  function merge(b4,       j,n,now,a,b,both)
395    j, n, now = 0, #b4, {}
396    while j < #b4 do
397      j    = j+1
398      a, b = b4[j], b4[j+1]
399      if b then
400        both = merged(a,b)
401        if both then a, j = {lo=a.lo, hi=b.hi, all=both}, j+1 end end
402      push(now,a)
403      j = j+1 end
404    return #now == #b4 and b4 or merge(now) end
405
406  function merged(i,j,     k,ei,ej,ek)
407    k = Sym(i.at,i.txt)
408    for x,n in pairs(i.all) do add(k,x,n) end
409    for x,n in pairs(j.all) do add(k,x,n) end
410    ei, ej, ek= div(i), div(j), div(k)
411    if i.n==0 or j.n==0 or 1.01*ek <= (i.n*ei + j.n*ej)/(i.n+j.n) then
412      return k end end
```

```
413  -----------------------------------------------------------------------------
414  ---
415  ---      __   __   ___  __
416  ---     |  \ /  \ /__`  |
417  ---     |__/ \__/ .__/  |
418  ---
419
420  function Demo.the() oo(the) end
421
422  function Demo.many(a)
423    a={1,2,3,4,5,6,7,8,9,10}; ok("{10 2 3}" == o(many(a,3)), "manys") end
424
425  function Demo.egs()
426    ok(5140==file2Egs(the.file).y[1].hi,"reading") end
427
428  function Demo.dist(i)
429    i = file2Egs(the.file)
430    for n,row in pairs(i.all) do print(n,dist(i, i.all[1], row)) end end
431
432  function Demo.far(  i,j,row1,row2,row3,d3,d9)
433    i = file2Egs(the.file)
434    for j=1,10 do
435      row1 = any(i.all)
436      row2 = far(i,row1, i.all, .9)
437      d9   = dist(i,row1,row2)
438      row3 = far(i,row1, i.all, .3)
439      d3   = dist(i,row1,row3)
440      ok(d3 < d9, "closer far") end end
441
442  function Demo.half(  i,easts,wests)
443    i = file2Egs(the.file)
444    easts,wests = half(i, i.all)
445    oo(mids(i.y, easts))
446    oo(mids(i.y, wests)) end
447
448  function Demo.cluster(   i)
449    i = file2Egs(the.file)
450    i = file2Egs(the.file)
451    clusters(i,cluster(i)) end
452
453  -----------------------------------------------------------------------------
454  the = settings(help)
455  Demo.main(the.todo, the.seed)
```