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

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

```lua
--- ═╦╗ ─  ─ ╔══  ─  ╦═╗   ═╔═  ═╗  ╦    ═══
--- ║║║ ─ ─  ─  ║    ─  ║ ║ ║    ║║   ╚═╗
--- ╚╩╝ ─ ─  ─  ╚══  ─  ╚═╝ ╩    ╩═╝  ═══

local add
function add(i,x, inc)
  inc = inc or 1
  if x ~= "?" then
    i.n = i.n + inc
    i:add1(x,inc) end
  return x end

function Sym.add1(i,x,inc)
  i.all[x] = inc + (i.all[x] or 0)
  if i.all[x] > i.most then i.most, i.mode = i.all[x], x end end

function Num.add1(i,x,inc,    d)
  for j=1,inc do
    d     = x - i.mu
    i.mu  = i.mu + d/i.n
    i.m2  = i.m2 + d*(x - i.mu)
    i.sd  = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n-1))^0.5)
    i.lo  = math.min(x, i.lo)
    i.hi  = math.max(x, i.hi)
    if      #i.all < the.keep       then push(i.all,x)
    elseif r()     < they.keep/i.n then i.all[r(#i.all)]=x end end end

--- ╔╦╗ ═╗ ╦ ═ ╔═╗ ═══
--- ║║║ ═╩═ ╠═ ║ ║ ─
--- ╩ ╩ ═╝ ╩ ╚ ╩ ╩ ═══

local header,data,file2Egs
function header(names,    i,col)
  i = Egs(names)
  for at,name in pairs(names) do
    col = push(i.cols, (name:find"^[A-Z]" and Num or Sym)(at,name))
    if not name:find":$" then
      push(name:find"[-+]$" and i.y or i.x, col) end end
  return i end

function data(i,row)
  push(i.all, row)
  for _,col in pairs(i.cols) do add(col, row[col.at]) end
  return i end

function file2Egs(file,    i)
  for row in file2things(file) do
    if i then data(i,row) else i = header(row) end end
  return i end

--- ╔═╗ ╦ ╦ ╔╦╗ ╔╦╗ ╔═╗ ═╗ ╦ ═══
--- ╚═╗ ║ ║ ║║║ ║║║ ╠═╣ ╠╦╝  ─
--- ╚═╝ ╚═╝ ╩ ╩ ╩ ╩ ╩ ╩ ╩╚  ═══

function Sym.mid(i) return i.mode end
function Sym.div(i,  e)
  e=0; map(i.all,function(n) e = e + n/i.n * math.log(n/i.n,2) end)
  return -e end

function Num.mid(i) return i.mu end
function Num.div(i) return i.sd end

function Num.clone(i) return Num(i.at, i.name) end
function Sym.clone(i) return Sym(i.at, i.name) end

local mids
function mids(cols,rows,    seen,tmp)
  seen = function(col) return col:clone() end
  tmp  = map(cols, seen)
  for _,row in pairs(rows) do
    for _,seen in pairs(tmp) do
      add(seen, row[seen.at]) end end
  return rnds(map(tmp, function(seen) return seen:mid() end)) end

--- ╔╦╗ ═╦ ╔═╗ ═╦═ ╔═╗ ╔╗╔ ╔═╗ ═══
--- ║║║  ║ ╚═╗  ║  ╠═╣ ║║║ ║    ─
--- ╩ ╩ ═╩ ╚═╝  ╩  ╩ ╩ ╝╚╝ ╚═╝ ═══

local far,furthest,neighbors,dist
function far(        i,r1,rows,far)
  return per(neighbors(i,r1,rows),far or the.far)[2] end

function furthest( i,r1,rows)
  return last(neighbors(i,r1,rows))[2] end

function neighbors(i,r1,rows)
  return sort(map(rows, function(r2) return {dist(i,r1,r2),r2} end),firsts) end

function dist(i,row1,row2,    d,n,a,b,inc)
  d,n = 0,0
  for _,col in pairs(i.x) do
    a,b = row1[col.at], row2[col.at]
    inc = a=="?" and b=="?" and 1 or col:dist1(a,b)
    d = d + inc^the.p
    n = n + 1 end
  return (d/n)^(1/the.p) end

function Sym.dist1(i,a,b) return a==b and 0 or 1 end

function Num.dist1(i,a,b)
  if      a=="?" then b=i:norm(b); a=b<.5 and 1 or 0
  elseif b=="?" then a=i:norm(a); b=a<.5 and 1 or 0
  else    a,b = i:norm(a), i:norm(b)   end
  return math.abs(a - b) end

function Num.norm(i,x)
  return i.hi - i.lo < 1E-32 and 0 or (x - i.lo)/(i.hi - i.lo) end

--- ╔═╗ ╦    ╦ ╔═╗ ═╗ ╦ ═══
--- ║   ║    ║ ╚═╗ ╠═ ╠╦╝  ─
--- ╚═╝ ╩═╝ ╚═╝ ╚═╝ ╩ ╩╚  ═══

local half, cluster, clusters
function half(i, rows,      project,row,some,east,west,easts,wests,c,mid)
  function project(row,a,b)
    a= dist(i,east,row)
    b= dist(i,west,row)
    return {(a^2 + c^2 - b^2)/(2*c), row}
  end -----------------------
  some = many(rows, the.some)
  east = furthest(i,any(some), some)
  west = furthest(i,east,          some)
  c    = dist(i,east,west)
  easts,wests = {},{}
  for n, xrow in pairs(sort(map(rows,project),firsts)) do
    row = xrow[2]
    if n==#rows//2 then mid=row end
    push(n <= #rows//2 and easts or wests, row) end
  return easts, wests, east, west, mid  end

function cluster(i,rows,  here,lefts,rights)
  rows = rows or i.all
  here = {all=rows}
  if #rows > 2*(#i.all)^the.leaves then
    lefts, rights = half(i, rows)
    if #lefts < #rows then
      here.lefts = cluster(i,lefts)
      here.rights= cluster(i,rights) end end
  return here end
```

```lua
function clusters(i,t,pre)
  if t then
    pre = pre or ""
    if not t.lefts and not t.rights then
      print(fmt("%5s %-20s",#t.all, pre), o(mids(i.y,t.all)))
    else
      print(fmt("%5s %-20s",#t.all, pre))
      clusters(i,t.lefts,  "|.."..pre)
      clusters(i,t.rights, "|.."..pre) end end end

--- ╔╦╗ ═╦ ╔═╗ ╔═╗ ═╗ ╦ ╔═╗ ═╦═ ═╦ ═╔═ ╔═╗ ═══
--- ║║║  ║ ╚═╗ ║    ╠═ ╠╦╝ ╠═   ║   ║  ╔╝  ╠═  ─
--- ╩ ╩ ═╩ ╚═╝ ╚═╝ ╩ ╩╚  ╚═╝  ╩  ═╩ ╚══ ╚═╝ ═══

local merge,merged
function Sym.spans(i, j)
  local xys,all,one,last,x,y,n = {}, {}
  for x,n in pairs(i.all) do push(xys, {x,"easts",n}) end
  for x,n in pairs(j.all) do push(xys, {x,"wests",n}) end
  for _,tmp in ipairs(sort(xys,firsts)) do
    x,y,n = unpack(tmp)
    if x ~= last then
      last = x
      one  = push(all, {lo=x, hi=x, all=Num(i.at,i.txt)}) end
    add(one.all, y, n) end
  return all end

function Num.spans(i, j)
  local xys,all,lo,hi,gap,one,x,y,n = {},{}
  lo,hi = math.min(i.lo, j.lo), math.max(i.hi,j.hi)
  gap   = (hi - lo) / (6/the.cohen)
  for _,n in pairs(i.all) do push(xys, {n,"easts",1}) end
  for _,n in pairs(j.all) do push(xys, {n,"wests",1}) end
  one = {lo=lo, hi=lo, all=Sym(i.at,i.txt)}
  all = {one}
  for _,tmp in ipairs(sort(xys,firsts)) do
    x,y,n = unpack(tmp)
    if    one.hi - one.lo > gap
    then one = push(all, {lo=one.hi, hi=x, all=Sym(i.at,i.txt)}) end
    one.hi = x
    add(one.all,y,n)  end
  all        = merge(all)
  all[1   ].lo = -math.huge
  all[#all].hi =  math.huge
  return all end

function merge(b4,       j,n,now,a,b,both)
  j, n, now = 0, #b4, {}
  while j < #b4 do
    j    = j+1
    a, b = b4[j], b4[j+1]
    if b then
      both = merged(a,b)
      if both then a, j = {lo=a.lo, hi=b.hi, all=both}, j+1 end end
    push(now,a)
    j = j+1 end
  return #now == #b4 and b4 or merge(now)  end

function merged(i,j,     k,ei,ej,ek)
  k = Sym(i.at,i.txt)
  for x,n in pairs(i.all)  do add(k,x,n) end
  for x,n in pairs(j.all)  do add(k,x,n) end
  ei, ej, ek= div(i), div(j), div(k)
  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
```

```
392  ------------------------------------------------------------------------------
393  ---
394  ---    __  __   __ _  __
395  ---   |  \/  | / _` |(_) _ _
396  ---   | |\/| || (_| || || ' \
397  ---   |_|  |_| \__,_||_||_||_|
398
399  function Demo.the() oo(the) end
400
401  function Demo.many(a)
402    a={1,2,3,4,5,6,7,8,9,10}; ok("{10 2 3}" == o(many(a,3)), "manys") end
403
404  function Demo.egs()
405    ok(5140==file2Egs(the.file).y[1].hi,"reading") end
406
407  function Demo.dist(i)
408    i = file2Egs(the.file)
409    for n,row in pairs(i.all) do print(n,dist(i, i.all[1], row)) end end
410
411  function Demo.far(  i,j,row1,row2,row3,d3,d9)
412    i = file2Egs(the.file)
413    for j=1,10 do
414      row1 = any(i.all)
415      row2 = far(i,row1, i.all, .9)
416      d9   = dist(i,row1,row2)
417      row3 = far(i,row1, i.all, .3)
418      d3   = dist(i,row1,row3)
419      ok(d3 < d9, "closer far") end end
420
421  function Demo.half(  i,easts,wests)
422    i = file2Egs(the.file)
423    easts,wests = half(i, i.all)
424    oo(mids(i.y, easts))
425    oo(mids(i.y, wests)) end
426
427  function Demo.cluster(   i)
428    i = file2Egs(the.file)
429    i = file2Egs(the.file)
430    clusters(i,cluster(i))
431   end
432
433  ------------------------------------------------------------------------------
434  the=settings(help)
435  Demo.main(the.todo, the.seed)
```