```lua
   1  -- vim : ft=lua et sts=2 sw=2 ts=2 :
   2  local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end --used later (to find rogues)
   3  local help = [[
   4
   5  sl == S.U.B.L.I.M.E. == Sublime's unsupervised
   6  bifurcation: let's infer minimal explanations.
   7  (c) 2022, Tim Menzies, BSD 2-clause license.
   8
   9  USAGE:
  10    lua sl.lua [OPTIONS]
  11
  12  OPTIONS:
  13    -Dump      stack dump on assert fails = false
  14    -data    N  data file               = etc/data/auto93.csv
  15    -enough  F  recurse until rows^enough = .5
  16    -far     F  far                      = .9
  17    -keep    P  max kept items           = 512
  18    -p       P  distance coefficient     = 2
  19    -seed    P  set seed                 = 10019
  20    -todo    S  start up action (or 'all') = nothing
  21    -help       show help                = false
  22
  23  KEY: N=fileName F=float P=posint S=string
  24  ]]
  25  local any,asserts,big,cli,fails,firsts,fmt,goalp,ignorep,klassp
  26  local lessp,map,main,many,max,min,morep,new,nump,o,oo,per,pop,push
  27  local r,rows,slots,sort,sum,thing,things,unpack
  28  local CLUSTER, COLS, EGS, NUM, ROWS, SKIP, SOME, SYM = {},{},{},{},{},{},{}
  29
  30  local the={}
  31  help:gsub("\n  [-]([^%s]+)[^\n]*%s([^%s]+)",function(key,x)
  32    for n,flag in ipairs(arg) do
  33      if flag:sub(1,1)=="-" and key:find("^"..flag:sub(2)..".*") then
  34        x = x=="false" and true or arg[n+1] end end
  35    if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
  36      the[key] = tonumber(x) or x end end )
  37
  38  -- ------------------------------------------------------------------------------
  39  -- this code reads csv files where the words on line1 define column types.
  40  function ignorep(x)  return x:find":$" end     -- columns to ignore
  41  function klassp(x)   return x:find"!$" end     -- symbolic goals to achieve
  42  function lessp(x)    return x:find"-$" end     -- number goals to minimize
  43  function morep(x)    return x:find"+$" end     -- numeric goals to maximize
  44  function nump(x)     return x:find"^[A-Z]" end -- numeric columns
  45  function goalp(x)    return morep(x) or lessp(x) or klassp(x) end
  46
  47  -- strings
  48  fmt = string.format
  49
  50  -- maths
  51  big = math.huge
  52  max = math.max
  53  min = math.min
  54  r   = math.random
  55
  56  -- tables
  57  pop = table.remove
  58  unpack = table.unpack
  59  function any(t)          return t[r(#t)] end
  60  function firsts(a,b)     return a[1] < b[1] end
  61  function many(t,n, u)  u={}; for i=1,n do push(u,any(t)) end; return u end
  62  function per(t,p)        return t[ (#t*(p or .5))//1 ] end
  63  function push(t,x)       table.insert(t,x); return x end
  64  function sort(t,f)       table.sort(t,f); return t end
  65
  66  -- meta
  67  function map(t,f, u)  u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
  68  function sum(t,f, n)  n=0; for _,v in pairs(t) do n=n+f(v)    end; return n end
  69  function slots(t, u)
  70    u={}
  71    for k,v in pairs(t) do k=tostring(k);if k:sub(1,1)~="_" then push(u,k) end end
  72    return sort(u) end
  73
  74  -- print tables, recursively
  75  function oo(t)  print(o(t)) end
  76  function o(t)
  77    if type(t)~="table" then return tostring(t) end
  78    local key=function(k) return fmt(":%s %s",k,o(t[k])) end
  79    local u = #t>0 and map(t,o) or map(slots(t),key)
  80    return '{'..table.concat(u,"")..'}' end
  81
  82  -- strings to things
  83  function csv(file,    x)
  84    file = io.input(file)
  85    return function()
  86      x=io.read(); if x then return things(x) else io.close(file) end end end
  87
  88  function thing(x)
  89    x = x:match"^%s*(.-)%s*$"
  90    if x=="true" then return true elseif x=="false" then return false end
  91    return tonumber(x) or x end
  92
  93  function things(x,sep,  t)
  94    t={}
  95    for y in x:gmatch(sep or"([^,]+)") do push(t,thing(y)) end
  96    return t end
```

```lua
  97  --  __  __  ____  __  __  _   _  ___  __  __
  98  -- |  \/  || ___||  \/  || \ | ||_ _||  \/  |
  99  -- | |\/| || _|  | |\/| ||  \| | | | | |\/| |
 100  --
 101  function new(k,t) k.__index=k; k.__tostring=o; return setmetatable(t,k) end
 102
 103  -- COLS: turns list of column names into NUMs, SYMs, or SKIPs
 104  function COLS.new(k,row,   i)
 105    i= new(k,{all={},x={},y={},names=row})
 106    for at,txt in ipairs(row) do  push(i.all, i:col(at,txt)) end
 107    return i end
 108
 109  function COLS.add(i,t)
 110    for _,col in pairs(i.all) do col:add( t[col.at] ) end
 111    return t end
 112
 113  function COLS.col(i,at,txt,    col)
 114    if ignorep(txt) then return SKIP:new(at,txt) end
 115    col = (nump(txt) and NUM or SYM):new(at,txt)
 116    push(goalp(txt) and i.y or i.x, col)
 117    if klassp(txt) then i.klass = col end
 118    return col end
 119
 120  -- NUM: summarizes a stream of numbers
 121  function NUM.new(k,n,s)
 122    return new(k,{n=0,at=n or 0,txt=s or"",has=SOME:new(),ok=false,
 123               w=lessp(s or "") and -1 or 1, lo=big, hi=-big}) end
 124
 125  function NUM.add(i,x)
 126    if x ~= "?" then
 127      i.n = i.n + 1
 128      if i.has:add(x) then i.ok=false end
 129      i.lo,i.hi = min(x,i.lo), max(x,i.hi); end end
 130
 131  function NUM.dist(i,x,y)
 132    if     x=="?" and y=="?" then return 1
 133    elseif x=="?" then y=i:norm(y); x=y<0.5 and 1 or 0
 134    elseif y=="?" then x=i:norm(x); y=x<0.5 and 1 or 0
 135    else   x,y = i:norm(x), i:norm(y) end
 136    return math.abs(x-y) end
 137
 138  function NUM.mid(i) return per(i:sorted(), .5) end
 139
 140  function NUM.norm(i,x)
 141    return math.abs(i.hi-i.lo)<1E-9 and 0 or (x-i.lo)/(i.hi - i.lo) end
 142
 143  function NUM.sorted(i)
 144    if i.ok==false then table.sort(i.has.all); i.ok=true end
 145    return i.has.all end
 146
 147  -- ROWS: manages 'rows', summarized in 'cols' (columns).
 148  function ROWS.new(k,inits,    i)
 149    i = new(k,{rows={},cols=nil})
 150    if type(inits)=="string" then for t in csv(inits) do i:add(t) end end
 151    if type(inits)=="table"  then for t in inits      do i:add(t) end end
 152    return i end
 153
 154  function ROWS.add(i,t)
 155    if i.cols then push(i.rows,i.cols:add(t)) else i.cols=COLS:new(t) end end
 156
 157  function ROWS.clone(i,  j) j= ROWS:new(); j:add(i.cols.names);return j end
 158
 159  function ROWS.dist(i,row1,row2,   d,fun)
 160    function fun(col) return col:dist(row1[col.at], row2[col.at])^the.p end
 161    return (sum(i.cols.x, fun)/ #i.cols.x)^(1/the.p) end
 162
 163  function ROWS.far(i,row1,rows,     fun)
 164    function fun(row2) return {i:dist(row1,row2), row2} end
 165    return unpack(per(sort(map(rows,fun), firsts), the.far)) end
 166
 167  function ROWS.half(i, top)
 168    local some, top,c,x,y,tmp,mid,lefts,rights,_
 169    some= many(i.rows, the.keep)
 170    top = top or i
 171    _,x = top:far(any(some), some)
 172    c,y = top:far(x,          some)
 173    tmp = sort(map(i.rows,function(r) return top:project(r,x,y,c) end),firsts)
 174    mid = #i.rows//2
 175    lefts, rights = i:clone(), i:clone()
 176    for at,row in pairs(tmp) do (at <=mid and lefts or rights):add(row[2]) end
 177    return lefts,rights,x,y,c, tmp[mid] end
 178
 179  function ROWS.mid(i,cols)
 180    return map(cols or i.cols.all, function(col) return col:mid() end) end
 181
 182  function ROWS.project(i, r,x,y,c,    a,b)
 183    a,b = i:dist(r,x), i:dist(r,y); return {(a^2 + c^2 - b^2)/(2*c), r} end
 184
 185  -- SKIP: summarizes things we want to ignore (so does nothing)
 186  function SKIP.new(k,n,s)  return new(k,{n=0,at=at or 0,txt=s or""}) end
 187  function SKIP.add(i,x)    return x end
 188  function SKIP.mid(i)      return "?" end
 189
 190  -- SOME: keeps a random sample on the arriving data
 191  function SOME.new(k,keep) return new(k,{n=0,all={}, keep=keep or the.keep}) end
 192  function SOME.add(i,x)
 193    i.n = i.n+1
 194    if      #i.all < i.keep then push(i.all,x)           ; return i.all
 195    elseif r()      < i.keep/i.n then i.all[r(#i.all)]=x; return i.all end end
 196
 197  -- SYM: summarizes a stream of symbols
 198  function SYM.new(k,n,s)
 199    return new(k,{n=0,at=n or 0,txt=s or"",has={},most=0}) end
 200
 201  function SYM.dist(i,x,y) return(x=="?" and y=="?" and 1) or(x==y and 0 or 1) end
 202  function SYM.mid(i)       return i.mode end
 203  function SYM.div(i,   p)
 204    return sum(i.has,function(k) p=-i.has[k]/i.n;return -p*math.log(p,2) end) end
 205
 206  function SYM.add(i,x,inc)
 207    if x ~= "?" then
 208      inc = inc or 1
 209      i.n = i.n + inc
 210      i.has[x] = inc + (i.has[x] or 0)
 211      if i.has[x] > i.most then i.most,i.mode=i.has[x],x end end end
 212
 213  function SYM.merge(i,j,    k)
 214    k = SYM:new(i.at,i.txt)
 215    for x,n in pairs(i.has) do k:add(x,n) end
 216    for x,n in pairs(j.has) do k:add(x,n) end
 217    ei, ej, ejk= i:div(), j:div(), k:div()
 218    if i.n==0 or j.n==0 or .99*ek <= (i.n*ei + j.n*ej)/k.n then
 219      return k end end
 220
 221
```

```lua
221  -- CLUSTER: recursively divides data by clustering towards two distant points
222  function CLUSTER.new(k,sample,top)
223    local i,enough,left,right
224    top    = top or sample
225    i      = new(k, {here=sample})
226    enough = (#top.rows)^the.enough
227    if #sample.rows >= 2*enough then
228      left, right, i.x, i.y, i.c, i.mid = sample:half(top)
229      if #left.rows < #sample.rows then
230        i.left = CLUSTER:new(left,   top)
231        i.right= CLUSTER:new(right, top) end end
232    return i end
233
234  function CLUSTER.show(i,pre,  here)
235    pre = pre or ""
236    here=""
237    if not i.left and not i.right then here= o(i.here:mid(i.here.cols.y)) end
238    print(fmt("%6s : %-30s %s",#i.here.rows, pre, here))
239    for _,kid in pairs{i.left, i.right} do
240      if kid then kid:show(pre .. "|.. ") end end end
```

```lua
241  --     ___   ___  __   __   __
242  --    |__)  |__  |  \ /  \ (_
243  --    |  \  |___ |__/ \__/ __)
244  --
245  fails=0
246  function asserts(test, msg)
247    print(test and "PASS: "or "FAIL: ",msg or "")
248    if not test then
249      fails=fails+1
250      if the.dump then assert(test,msg) end end end
251
252  function EGS.nothing() return true end
253  function EGS.the()     oo(the) end
254  function EGS.rand()    print(r()) end
255  function EGS.some(s,t)
256    s=SOME:new(100)
257    for i=1,100000 do s:add(i) end
258    for j,x in pairs(sort(s.all)) do
259      if (j % 10)==0 then print("") end
260      io.write(fmt("%6s",x))   end end
261
262  function EGS.clone( r,s)
263    r = ROWS:new(the.data)
264    s = r:clone()
265    for _,row in pairs(r.rows) do s:add(row) end
266    asserts(r.cols.x[1].lo==s.cols.x[1].lo,"clone.lo")
267    asserts(r.cols.x[1].hi==s.cols.x[1].hi,"clone.hi")
268    end
269
270  function EGS.data( r)
271    r = ROWS:new(the.data)
272    asserts(r.cols.x[1].hi == 8, "data.columns") end
273
274  function EGS.dist( r,rows,n)
275    r = ROWS:new(the.data)
276    rows = r.rows
277    n = NUM:new()
278    for _,row in pairs(rows) do n:add(r:dist(row, rows[1])) end
279    oo(r.cols.x[2]:sorted()) end
280
281  function EGS.many(   t)
282    t={}; for j=1,100 do push(t,j) end
283    print(oo(many(t, 10))) end
284
285  function EGS.far(    r,c,row1,row2)
286    r = ROWS:new(the.data)
287    row1   = r.rows[1]
288    c,row2 = r:far(r.rows[1], r.rows)
289    print(c,"\n",o(row1),"\n", o(row2)) end
290
291  function EGS.half(   r,c,row1,row2)
292    local lefts,rights,x,y,x
293    r = ROWS:new(the.data)
294    oo(r:mid(r.cols.y))
295    lefts,rights,x,y,c = r:half()
296    oo(lefts:mid(lefts.cols.y ))
297    oo(rights:mid(rights.cols.y))
298    end
299
300  function EGS.cluster(r)
301    r = ROWS:new(the.data)
302    CLUSTER:new(r):show() end
303
304  -- start-up
305  if arg[0] == "sl.lua" then
306    oo(the)
307    if the.help then print(help) else
308      local b4={}; for k,v in pairs(the) do b4[k]=v end
309      for _,todo in pairs(the.todo=="all" and slots(EGS) or {the.todo}) do
310        for k,v in pairs(b4) do the[k]=v end
311        math.randomseed(the.seed)
312        if type(EGS[todo])=="function" then EGS[todo]() end end
313    end
314    for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
315    os.exit(fails)
316  else
317    return {CLUSTER=CLUSTER, COLS=COLS, NUM=NUM, ROWS=ROWS,
318            SKIP=SKIP, SOME=SOME, SYM=SYM,the=the,oo=oo,o=o}
319  end
320  -- git rid of SOME for rows
321  -- nss  = NUM | SYM | SKIP
322  -- COLS = all:[nss]+, x:[nss]*, y:[nss]*, klass;col?
323  -- ROWS = cols:COLS, rows:SOME
```