```lua
1   local _=require("lib")
2   local the=_.settings[[
3   TINY2: a lean little learning library, in LUA
4   (c) 2022 Tim Menzies <timm@ieee.org> BSD-2 license
5
6   USAGE: lua l5.lua [OPTIONS]
7
8   OPTIONS:
9   -b --bins      max number of bins                = 8
10  -d --dump      on test failure, exit with stack dump  = false
11  -f --file      file with csv data                = ../data/auto93.csv
12  -F --Far       how far to look for poles (max=1) = .95
13  -g --go        start-up example                  = nothing
14  -h --help      show help                         = false
15  -m --min       min size. If<1 then t^min else min. = 10
16  -n --nums      number of nums to keep            = 512
17  -p --p         distance calculation coefficient  = 2
18  -r --rest      size of "rest" set                = 3
19  -s --seed      random number seed                = 10019
20  -S --Sample    how many numbers to keep          = 10000]]
21
22  local any,cli,copy,csv,lt,many,map = _.any,_.cli,_.copy,_.csv,_.lt,_.many,_.map
23  local o,obj,oo,per,pop,push        = _.o,_.obj,_.oo,_.per,_.pop,_.push
24  local rnd,rogues                   = _.rnd,_.rogues
25  local shallowCopy,shuffle,sort     = _.shallowCopy,_.shuffle,_.sort
26  local Data,Num,Row,Some,Sym = obj"Data",obj"Num",obj"Row",obj"Some",obj"Sym"
27
28  --[[ Type hints conventions:
29  | Function args  | Notes
30  :---------------------------------------------------------------
31  | 2 blanks       | 2 blanks denote optional arguments
32  | 4 blanks       | 4 blanks denote local arguments
33  | n              | prefix for numerics
34  | s              | prefix for strings
35  | is             | prefix for booleans
36  | fun            | prefix for functions
37  | suffix s       | list of thing (so names is list of strings)
38  | xy,row,col,data| for Xys, Rows, Num or Syms, Data objects
39
40  Another convention is that my code starts with a  help string (at top
41  of file) that is parsed to find the settings. Also my code ends with
42  lots of `go.x()` functions that describe various demos. To run
43  these, use 'lua tiny2.lua -go x'. --]]
44  -----------------------------------------------------------------
45  function Row:new(t) --- Hold one record
46     return {evaled=false,
47             cells=t,
48             cooled=shallowCopy(t)} end
49
50  function Sym:new(n,s) --- Summarize stream of symbols.
51     return {at=n or 0,
52             txt=s or "",
53             n=0,
54             has={}} end
55
56  function Some:new(n,s)  --- Keep at most the.Sample numbers
57     return {at=n or 0, txt=s or "",n=0, _has={},
58             isSorted=true } end
59
60  function Num:new(c,x) --- Summarize stream of numbers
61     return {at=c or 0,txt=x or "",n=0,
62             lo=1E32,hi=-1E32, n=0,
63             has=Some(),
64             w=(x or ""):find"-$" and -1 or 1} end
65
66  function Data:new(src) --- Store rows of data. Summarize the rows in `self.cols`
67     self.rows, self.cols = {}, {all={},x={},y={}}
68     if   type(src)=="string"
69     then csv(src,           function(row) self:add(row) end)
70     else map(src or {}, function(row) self:add(row) end) end end
```

```lua
71  -- ## Row      ----- ----- ----------------------------------------
72  -- ### sort
73  function Row:better(row2,data) --- order two rows
74     local row1= self
75     row1.evaled, row2.evaled = true,true
76     local s1,s2,d,n,x,y,ys=0,0,0,0
77     ys = data.cols.y
78     for _,col in pairs(ys) do
79       x,y= row1.cells[col.at], row2.cells[col.at]
80       x,y= col:norm(x), col:norm(y)
81       s1 = s1 - 2.71828^(col.w * (x-y)/#ys)
82       s2 = s2 - 2.71828^(col.w * (y-x)/#ys) end
83     return s1/#ys < s2/#ys end
84
85  function Row:betters(rows,data) --- order a whole list of rows
86     return sort(rows or self.rows,
87                 function(r1,r2) return r1:better(r2,data) end) end
88
89  -- #### dist
90  function Row:dist(row2,data,  tmp,n,d1) -- distance between rows
91     local row1=self
92     tmp,n = 0,0; for i,col in pairs(data.cols.x) do
93              d1    = col:dist(row1[col.at], row2[col.at],data)
94              n, tmp = n + 1,  tmp + d1^the.p end
95     return (tmp/n)^(1/the.p) end
96
97  function Row:dists(r1,rows,data) --- sort `rows` by distance to `r11.
98     return sort(map(rows,
99                 function(r2) return {r=r2,d=self:dist(r1,r2,data)} end),lt"d") end
100
101 function Row:far(row,rows,data) -- Find an item in `rows`, far from `row1.
102    return per(self:dists(row,rows,data),the.far).r end
103
104 -- ## Sym      ----- ----- ----------------------------------------
105 -- ### update
106 function Sym:add(s) --- Update.
107    if s~="?" then self.n =1+self.n;self.has[s]=1+(self.has[s] or 0) end end
108
109 -- ### dist
110 function Sym:dist(s1,s2) -- Gap between two symbols.
111    return  s1=="?" and s2=="?" and 1 or s1==s2 and 0 or 1 end
112
113 -- ### query
114 function Sym:entropy(   e,fun) -- Entropy
115    function fun(p) return p*math.log(p,2) end
116    e=0; for _,n in pairs(self.has) do if n>0 then e=e-fun(n/self.n) end end
117    return e end
118
119 -- ## Some  ----- ----- ----------------------------------------
120 -- ### update
121 function Some:add(x,    pos) --- update
122    if x~="?" then
123      self.n =  self.n+1
124      if #self._has < the.Sample then pos=1+(#self._has)
125      elseif math.random()<the.Sample/self.n then pos=math.rand(#self._has) end
126      if pos then self.isSorted=false
127                  self._has[pos]= x end end end
128
129 -- ### query
130 function Some:nums()
131    if not self.isSorted then table.sort(self._has) end
132    self.isSorted=true
133    return self._has end
134
135 -- ## Num     ----- ----- ----------------------------------------
136 -- ### update
137 function Num:add(n) --- update
138    if n~="?" then self.n = self.n+1
139                   self.lo = math.min(n,self.lo)
140                   self.hi = math.max(n,self.hi)
141                   self.has:add(n)   end end
142
143 -- ### query
144 function Num:norm(n,   lo,hi) --- convert 'n' to 0..1 for min..max
145    lo,hi=self.lo,self.hi
146    return n=="?" and n or (hi-lo < 1E-0 and 0 or  (n-lo)/(hi-lo + 1E-32)) end
147
148 function Num:pers(ns,    a) --- report a list over percentiles
149    a=self.has:nums()
150    return map(ns,function(p) return per(a,p) end) end
151
152 -- ### dist
153 function Num:dist(n1,n2) --- return 0..1. If unknowns, assume max distance.
154    if   n1=="?" and n2=="?" then return 1 end
155    n1,n2 = self:norm(n1), self:norm(n2)
156    if n1=="?" then n1 = n2<.5 and 1 or 0 end
157    if n2=="?" then n2 = n1<.5 and 1 or 0 end
158    return math.abs(n1-n2) end
```

```lua
159 -- ## Data     ----- ----- ----------------------------------------
160 -- ### create
161 function Data:body(row) --- Crete new row. Store in 'rows'. Update cols.
162    row = row.cell and row or Row(row) -- Ensure 'row' is a 'Row'.
163    push(self.rows, row)
164    for _,cols in pairs{self.cols.x, self.cols.y} do
165      for _,col in pairs(cols) do
166        col:add(row.cells[col.at]) end end end
167
168 function Data:clone( src,   data) --- Copy structure. Optionally, add in data.
169    data= Data( {map(self.all, function(col) return col.txt end)} )
170    map(src or {}, function (row) data:add(row) end)
171    return data end
172
173 function Data:header(row) --- Create the 'Num's and 'Sym's for the column headers
174    for n,s in pairs(row) do
175      local col = push(self.cols.all, (s:find"^[A-Z]" and Num or Sym)(n,s))
176      if not s:find":$" then
177        push(s:find"[!+-]" and self.cols.y or self.cols.x, col) end end end
178
179 -- ### udpate
180 function Data:add(row) --- the new row is either a header, or a data row
181    if #self.cols.all==0 then self:header(row) else self:body(row) end end
182
183 -- ### query
184 function Data:cheat(   ranks) --- return percentile ranks for rows
185    for i,row in pairs(self:betters()) do
186      row.rank = math.floor(.5+ 100*i/#self.rows) end
187    self.rows = shuffle(self.rows)
188    return self.rows end
189
190 -- ### cluster
191 function Data:half(  above, --- split data by distance to two distant points
192                 some,x,y,c,rxs,xs,ys)
193    some= many(self.rows, the.Sample)
194      x= above or self:far(any(some),some,data)
195      y= self:far(x,some,data)
196      c= self:dist(x,y,data)
197    rxs=function(r) return
198         {r=r,x=(self:dist(r,x,data)^2 + c^2 - self:dist(r,y,data)^2)/(2*c)} end
199    xs,ys= self:clone(), self:clone()
200    for j,rx in pairs(sort(map(self.rows,rxs),lt"x")) do
201      if j<=#self.rows/2 then xs:add(rx.r) else ys:add(rx.r) end end
202    return {xs=xs, ys=ys, x=x, y=y, c=c} end
203
204 function Data:best(  above,stop,evals) --- recursively divide, looking 4 best leaf
205    stop = stop or (the.min >=1 and the.min or (#self.rows)^the.min)
206    evals= evals or 2
207    if   #self.rows < stop
208    then return self,evals
209    else local node = self:half(above)
210         if    self:better(node.x,node.y)
211         then  return node.xs:best(node.x, stop, evals+1)
212         else  return node.ys:best(node.y, stop, evals+1) end end end
```

```lua
-- ## Demos/Tests  ----- ----- --------------------------------------------------
local go = {}
function go.the() oo(the); return true end

function go.num(  z)
  z=Num(); for i=1,100 do z:add(i) end; print(z); return true end

function go.sym(  z)
  z=Sym(); for _,x in pairs{1,1,1,1,2,2,3} do z:add(x) end;
  print(z); return true end

function go.eg( d)
  d=Data(the.file);  map(d.cols.x,print) return true end

function go.dist(     num,d,r1,r2,r3)
  d=Data(the.file)
  num=Num()
  for i=1,20 do
    r1= any(d.rows)
    r2= any(d.rows)
    r3= r1:far(d.rows,d)
    io.write(rnd(r1:dist(r1,r3,d))," ")
    num:add(rnd(r1:dist( r1,r2,d))) end
  oo(sort(num.has:nums()))
  print(#d.rows)
  return true end

function go.sort(     d,rows,ranks)
  d = Data(the.file)
  rows,ranks = d:cheat()
  for i=1,#d.rows,32 do print(i,ranks[rows[i][1]],o(rows[i])) end end

function go.clone(    d1,d2)
  d1 = Data(the.file)
  d2 = d1:clone(d1.rows)
  oo(d1.cols.x[2])
  oo(d2.cols.x[2]) end

function go.half( d,node)
  d=Data(the.file)
  node = d:half()
  print(#node.xs.rows, #node.ys.rows, d:dist(node.x, node.y,d))end

function go.best(      num)
  num=Num()
  for i=1,20 do
    local d=Data(the.file)
    local _,ranks = d:cheat()
    shuffle(d.rows)
    local leaf,evals = d:best()
    for _,row in pairs(leaf.rows) do num:add(ranks[ row[1] ]) end end
  print(o(num:pers(.1,.3,.5,.7,.9)))
end

function go.bests(      num,tmp)
  num=Num()
  for i=1,20 do
    local d = Data(the.file)
    d:cheat()
    shuffle(d.rows)
    tmp=d:best()
    map(tmp,function(row) num:add(row.rank) end) end
  print(#tmp,o(num:pers(.1,.3,.5,.7,.9)))
  return end

function go.discretize(   d)
  d=Data(the.file)
  print(d:xentropy()); return true end

function go.four(     num,d,some,evals,ranks)
  num=Num()
  for i=1,20 do
    d=Data(the.file)
    --_,ranks= d:cheat()
    some,evals = d:fours()
    _,ranks = d:cheat()
    print(#some)
    for _,row in pairs(some) do num:add(ranks[row[1]]) end end
  oo(num:pers(.1,.3,.5,.7,.9})
end
-- ## Start  ----- ----- --------------------------------------------------
local function on(settings,funs,    fails,old)
  fails=0
  old = copy(settings)
  for k,fun in pairs(funs) do
    if settings.go == "all" or settings.go == k then
      for k,v in pairs(old) do settings[k]=v end
      math.randomseed(settings.seed or 10019)
      print("\n>>>>",k)
      if not fun() then fails = fails+1 end end end
  rogues()
  os.exit(fails) end

the = cli(the)
on(the,go)
```