

```

1 local _require("lib")
2 local the=_settings[
3   TINY2: a lean little learning library, in LUA
4   (c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license
5
6 USAGE: lua 15.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,co,per,pop,push = _._o,_._obj,_._co,_._per,_._pop,_._push
24 local rnd,roques = _._rnd,_._roques
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 -- This code uses the following type hints for function arguments:
29
30 -- What      Notes
31 -- :-----:
32 -- 2 blanks  2 blanks denote optional arguments
33 -- 4 blanks  4 blanks denote local arguments
34 -- n         prefix for numerics
35 -- s         prefix for strings
36 -- is        prefix for booleans
37 -- fun       prefix for functions
38 -- suffix s  list of thing (so names is list of strings)
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 {evald=false,
47           cellst=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

```

```

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

```

```

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

```

```

211 -- ## Demos/Tests -----
212 local go = {}
213 function go.the() oo(the); return true end
214
215 function go.num( z)
216     z=Num(); for i=1,100 do z:add(i) end; print(z); return true end
217
218 function go.sym( z)
219     z=Sym(); for _,x in pairs{1,1,1,2,2,3} do z:add(x) end;
220     print(z); return true end
221
222 function go.eg( d)
223     d=Data(the.file); map(d.cols.x,print) return true end
224
225 function go.dist( num,d,r1,r2,r3)
226     d=Data(the.file)
227     num=Num()
228     for i=1,20 do
229         r1= any(d.rows)
230         r2= any(d.rows)
231         r3= d:far(r1,d.rows,d)
232         io.write(rnd(d:dist(r1,r3,d)), " ")
233         num:add(rnd(d:dist( r1,r2,d))) end
234     oo(sort(num.has:nums()))
235     print(#d.rows)
236     return true end
237
238 function go.sort( d,rows,ranks)
239     d = Data(the.file)
240     rows,ranks = d:cheat()
241     for i=1,#d.rows,32 do print(i,ranks[rows[i][1]],o(rows[i])) end end
242
243 function go.clone( d1,d2)
244     d1 = Data(the.file)
245     d2 = d1:clone(d1.rows)
246     oo(d1.cols.x[2])
247     oo(d2.cols.x[2]) end
248
249 function go.half( d,node)
250     d=Data(the.file)
251     node = d:half()
252     print(#node.xs.rows, #node.ys.rows, d:dist(node.x, node.y,d))end
253
254 function go.best( num)
255     num=Num()
256     for i=1,20 do
257         local d=Data(the.file)
258         local _,ranks = d:cheat()
259         shuffle(d.rows)
260         local leaf,evals = d:best()
261         for _,row in pairs(leaf.rows) do num:add(ranks[ row[1] ]) end end
262     print(o(num:pers{.1,.3,.5,.7,.9}))
263 end
264
265 function go.bests( num,tmp)
266     num=Num()
267     for i=1,20 do
268         local d = Data(the.file)
269         d:cheat()
270         shuffle(d.rows)
271         tmp=d:best()
272         map(tmp,function(row) num:add(row.rank) end) end
273     print(#tmp,o(num:pers{.1,.3,.5,.7,.9}))
274     return end
275
276 function go.discretize( d)
277     d=Data(the.file)
278     print(d:xentropy()); return true end
279
280 function go.four( num,d,some,evals,ranks)
281     num=Num()
282     for i=1,20 do
283         d=Data(the.file)
284         _,ranks= d:cheat()
285         some,evals = d:fours()
286         _,ranks = d:cheat()
287         print(#some)
288         for _,row in pairs(some) do num:add(ranks[row[1]]) end end
289     oo (num:pers{.1,.3,.5,.7,.9})
290 end
291 -- ## Start -----
292 local function on(settings,funs, fails,old)
293     fails=0
294     old = copy(settings)
295     for k,fun in pairs(funs) do
296         if settings.go == "all" or settings.go == k then
297             for k,v in pairs(old) do settings[k]=v end
298             math.randomseed(settings.seed or 10019)
299             print("\n>>>>",k)
300             if not fun() then fails = fails+1 end end end
301     roques()
302     os.exit(fails) end
303
304 the = cli(the)
305 on(the,go)

```