

```

1  #!/usr/bin/env lua
2  local _require=("lib")
3  local the=_settings[[
4
5  L5 : a lean little learning library, in LUA
6  (c) 2022 Tim Menzies <timmee.org> BSD-2 license
7
8  USAGE: lua l5.lua [OPTIONS]
9
10 OPTIONS:
11 -e --eg          start-up example          = nothing
12 -b --bins        max number of bins         = 8
13 -d --dump        on test failure, exit with stack dump = false
14 -f --file        file with csv data         = ../data/auto93.csv
15 -F --far         how far to look for poles (max=1) = .95
16 -h --help        show help                  = false
17 -m --min         min size. If<1 then t*min else min. = 10
18 -n --nums        number of nums to keep      = 512
19 -p --p           distance calculation coefficient = 2
20 -r --rest        size of "rest" set          = 3
21 -s --seed        random number seed         = 10019
22 -S --Sample      how many numbers to keep    = 10000 ]]
23
24 local any,cli,copy,csv,lt,many,map= _any,_cli,_copy,_csv,_lt,_many,_map
25 local o,obj,oo,per,push,rnd,rogues= _o,_obj,_oo,_per,_push,_rnd,_rogues
26 local shallowCopy,shuffle,sort = _sort,_shallowCopy,_sort
27 local Data,Num,Row,Some,Sym
28
29 -----
30 Sym=obj*"Sym"
31 function Sym:new(c,x) return {at=c or 0,txt=x or "",n=0,has={}} end
32 function Sym:add(x)
33   if x=="?" then self.n =1+self.n;self.has[x]=1+(self.has[x] or 0) end end
34 function Sym:discretize(x) return x end
35 function Sym:dist(v1,v2)
36   return v1=="?" and v2=="?" and 1 or v1==v2 and 0 or 1 end
37 function Sym:entropy(e,fun)
38   function fun(p) return p*math.log(p,2) end
39   e=0; for _,n in pairs(self.has) do if n>0 then e=e-fun(n/self.n) end end
40   return e end
41
42 -----
43 Some=obj*"Some"
44 function Some:new(c,x)
45   return {at=c or 0, txt=x or "",n=0,isSorted=true, _has={}} end
46 function Some:nums()
47   if not self.isSorted then table.sort(self._has) end
48   self.isSorted=true
49   return self._has end
50 function Some:add(v, pos)
51   if v=="?" then
52     self.n=self.n+1
53     if #self._has < the.Sample then pos=1+(#self._has)
54       elseif math.random()<the.Sample/self.n then pos=math.rand(#self._has) end
55     if pos then self.isSorted=false
56       self._has[pos]=v end end end
57
58 -----
59 Num=obj*"Num"
60 function Num:new(c,x)
61   return {at=c or 0,txt=x or "",lo=1E32,hi=-1E32, n=0, has=Some(),
62     w=(x or ""):find"$" and -1 or 1} end
63 function Num:add(x)
64   if x=="?" then self.n = self.n+1
65     self.lo = math.min(x,self.lo)
66     self.hi = math.max(x,self.hi)
67     self.has:add(x) end end
68 function Num:norm(n)
69   return n=="?" and n or (n-self.lo)/(self.hi-self.lo + 1E-32) end
70 function Num:pers(t, a)
71   a=self.has:nums()
72   return map(t,function(p) return per(a,p) end) end
73 function Num:discretize(x, tmp)
74   tmp = (self.hi - self.lo)/(the.bins - 1)
75   return self.lo == self.hi and 1 or math.floor(x/tmp+.5)*tmp end
76
77 function Num:dist(v1,v2)
78   if v1=="?" and v2=="?" then return 1 end
79   v1,v2 = self:norm(v1), self:norm(v2)
80   if v1=="?" then v1 = v2<.5 and 1 or 0 end
81   if v2=="?" then v2 = v1<.5 and 1 or 0 end
82   return math.abs(v1-v2) end
83
84 -----
85 Row=obj*"Row"
86 function Row:new(data,t) return {_data=data,cells=t} end
87 function Row:around(rows)
88   return sort(map(rows, function(r) return {row=r,d=self.r} end),lt"d") end
89 function Row:far(rows)
90   return per(self:around(rows),the.far).row end
91
92 function Row:__sub(row, d,n,dl,nl)
93   d,n = 0,0
94   for i,col in pairs(self._data.cols.x) do
95     dl= col:dist(self.cells[col.at], row.cells[col.at])
96     n = n + 1
97     d = d + dl*the.p end
98   return (d/n)^(1/the.p) end
99
100 function Row:__lt(row)
101   self.evald, row.evald = true,true
102   local s1,s2,d,n,x,y=0,0,0,0
103   local ys = self._data.cols.y
104   for _,col in pairs(ys) do
105     x,y= self.cells[col.at], row.cells[col.at]
106     x,y= col:norm(x), col:norm(y)
107     s1 = s1 - 2.71828*(col.w * (x-y)/#ys)
108     s2 = s2 - 2.71828*(col.w * (y-x)/#ys) end
109   return s1/#ys < s2/#ys end
110
111 function Row:discretize()
112   self.cooked=map(self._data.cols.all,
113     function(col) col:discretize(self.cells[col.at]) end) end
114
115
116

```

```

116 -----
117 Data=obj*"Data"
118 function Data:new(src)
119   self.rows, self.cols = {}, {all={},x={},y={}}
120   if type(src)=="string"
121     then csv(src, function(row) self:add(row) end)
122     else map(src or {}, function(row) self:add(row) end) end end
123
124 function Data:add(row, what)
125   function what(c,x)
126     return {x:find"[A-Z]" and Num or Sym}(c,x) end
127   if #self.cols.all==0
128     then for c,x in pairs(row) do
129       local col = push(self.cols.all, what(c,x))
130       if not x:find"$" then
131         push(x:find"[+-]" and self.cols.y or self.cols.x, col) end end
132     else row = row.cells and row or Row(self,row)
133     push(self.rows, row)
134     for _,cols in pairs(self.cols.x, self.cols.y) do
135       for _,col in pairs(cols) do col:add(row.cells[col.at]) end end end end
136
137 function Data:cheat()
138   for i,row in pairs(sort(self.rows)) do
139     row.rank = math.floor(.5+ 100*i/#self.rows)
140     row.evald = false end
141   self.rows = shuffle(self.rows) end
142
143 function Data:half(rows, above, some,x,y,c,rxs,xs,ys)
144   rows = rows or self.rows
145   some = many(rows, the.Sample)
146   x = above or any(some):far(some)
147   y = x:far(some)
148   c = x - y
149   rxs = function(r) return {r=r,x=(r-x)^2 + c^2 - (r-y)^2)/(2*c)} end
150   xs,ys= {},{}
151   for j,rx in pairs(sort(map(rows,rxs),lt"x")) do
152     push(js<#rows/2 and xs or ys, rx.r) end
153   return {xs=xs,ys=ys,x=x,y=y,c=c} end
154
155 function Data:best(rows, above,stop)
156   rows = rows or self.rows
157   stop = stop or (the.min >=1 and the.min or (#rows)*the.min)
158   if #rows < stop
159     then return rows
160     else local node = self:half(rows,above)
161       if node.x < node.y
162         then return self:best(node.xs, node.x, stop)
163         else return self:best(node.ys, node.y, stop) end end end
164
165 function Data:fours(rows, stop)
166   rows = rows or shallowCopy(self.rows)
167   stop = stop or (the.min >=1 and the.min or (#rows)*the.min)
168   if #rows < stop
169     then return rows
170     else rows = shuffle(rows)
171     fours = {}
172     for row1 in pairs(rows) do
173       four1 = sort(map(fours,function(row2)
174         t[four1_id] = t[four1_id] or {}
175         t[four1_id] = t[four1_id] or {}
176         push(t[four1_id], four1) end),lt"d")[1].r
177       self:fours(sort(fours)[1],stop) end
178     end
179
180 function Data:discretize()
181   for _,row in pairs(self.rows) do row:discretize() end end
182
183 function Data:xentropy(e,sym)
184   self:discretize()
185   e=0
186   for _,col in pairs(self.cols.x) do
187     sym = Sym()
188     for _,row in pairs(self.rows) do sym:add(row.cooked[col.at]) end
189     e = e + sym:entropy() end
190   return e end
191
192

```

```

192 -----
193 local eg = {}
194 local function egs( fails,cold)
195   the = cli(the)
196   fails=0
197   old = copy(the)
198   for k,fun in pairs(eg) do
199     if the.eg == "all" or the.eg == k then
200       for k,v in pairs(old) do the[k]=v end
201       math.randomseed(the.seed)
202       print("n>>>>>",k)
203       if not fun() then fails = fails+1 end end end
204   rogues()
205   os.exit(fails) end
206
207 function eg.the() oo(the); return true end
208
209 function eg.num( z)
210   z=Num(); for i=1,100 do z:add(i) end; print(z); return true end
211
212 function eg.sym( z)
213   z=Sym(); for _,x in pairs{1,1,1,1,2,2,3} do z:add(x) end;
214   print(z); return true end
215
216 function eg.data( d)
217   d=Data(the.file); map(d.cols.x,print) return true end
218
219 function eg.dist( num,d,r1,r2,r3)
220   d=Data(the.file)
221   num=Num()
222   for i=1,20 do
223     r1=1:any(d.rows)
224     r2=1:any(d.rows)
225     r3=r1:far(d.rows)
226     io.write(rnd(r3-r1)," ")
227     num:add(rnd(r2-r1)) end
228   oo(sort(num.has:nums()))
229   return true end
230
231 function eg.sort( d)
232   d = Data(the.file)
233   sort(d.rows)
234   for i=1,#d.rows,32 do print(i,o(d.rows[i].cells)) end end
235
236 function eg.half( num,tmp)
237   num=Num()
238   for i=1,20 do
239     local d = Data(the.file)
240     d:cheat()
241     tmp=d:best()
242     map(tmp,function(row) num:add(row.rank) end) end
243   print(#tmp,o(num:pers(.1,.3,.5,.7,.9)))
244   return end
245
246 function eg.discretize( d)
247   d=Data(the.file)
248   print(d:xentropy()); return true end
249
250 function eg.fours( d)
251   d=Data(the.file)
252   d:fours() end
253
254 egs()

```