```lua
1   -- If you understand "it", can you write "it" shorter?  Lets try.
2   -- E.G. how short can we write a multi-objective semi-supervised learner?<p>
3   -- &copy;2022 Tim Menzies.  [Github](http://github.com/timm/15) <hr>
4   --
5   -- One of my most productive days was throwing away 1,000 lines of code.
6   -- -- Ken Thompson<p>
7   -- It is vain to do with more what can be done with less.
8   -- -- William Of Occam<p>
9   -- Every block of stone has a statue inside it.
10  -- And it is the task of the sculptor to discover it.
11  -- -- Michelangelo.<p>
12  -- The more you have, the more you are occupied.
13  -- The less you have, the more free you are.  <br>-- Mother Teresa<p>
14  -- <img width=200 align=left src=cup.png>
15  -- Simplicity is the ultimate sophistication. <br>-- Leonardo da Vinci<p>
16  -- Simplicity is prerequisite for reliability.<br>-- Edsger W. Dijkstra<p>
17  -- Less, but better.                 <br>-- Dieter Rams<p>
18  -- less, plz                         <br>-- timm <p>
19  -- My heroes: [Jack Diederich](https://www.youtube.com/watch?v=o9pEzgHorH0)
20  -- | [Hilary Mason](https://boingboing.net/2017/06/30/next-level-regexp.html)<p>
21  local  help= [[
22  shorter.lua : a multi-objective semi-supervised learner.
23  (c)2022 Tim Menzies, timm@ieee.org
24
25  OPTIONS:
26   -b  --Bins   max number of bins       = 16
27   -F  --Few    only keep a "Few" numbers = 256
28   -k  --k      handle rare classes      = 1
29   -m  --m      handle rare attributes   = 2
30   -p  --p      distance coefficient     = 2
31   -S  --small  small leaf size          = .5
32   -w  --wait   wait before classifying  =
33
35  OPTIONS (other):
36   -f  --file   file         = ../../data/auto93.csv
36   -g  --go     start-up goal = nothing
37   -h  --help   show help     = false
38   -s  --seed   seed          = 10019]]
39  -----------------------------------------------------------------
40  -- ## Names
41  local _ = require"lib"
42  local argmax,big                     = _.argmax, _.big
43  local cli,csv,demos,klass,normpdf = _.cli,  _.csv,  _.demos,_.klass, _.normpdf
44  local oo,push,read,rnd,same,str    = _.oo,    _.push, _.read, _.rnd,_.same,_.str
45
46  -- 'THE' settings is parsed from any 'help' i
47  -- string lines that contain two dashes "'--'".
48  local THE={}
49  help:gsub(" [-|[-]([^%s]+)[^\n]*%s([^%s]+)",function(key,x) THE[key] = read(x) end)
50  -- ### Classes
51
52  -- (1) ROWS use COLS to make either NUMs or SYMs.
53  -- (2) ROWS holds data in ROWS, and summarizes columns in NUMs and SYMs.
54  -- (3) NUMs use SOMEs to store at most 'THE.Few' samples per numeric columns.
55  -- (4) RANGE objects track what 'y' values are seen between 'xlo' and 'xhi'.
56  local ROWS, COLS, NUM, SYM = klass"ROWS", klass"COLS", klass"NUM", klass"SYM"
57  local ROW = klass"ROW"
58  local SOME = klass"SOME"
59  local RANGE = klass"RANGE"
60  -----------------------------------------------------------------
61  -- ## class RANGE
62  function RANGE.new(i, xlo, xhi, ys) i.xlo,i.xhi,i.ys,i.rows = xlo,xhi,ys,{} end
63  function RANGE.add(i,x,y)
64    if x < i.xlo then i.xlo = x end -- works for string or num
65    if x > i.xhi then i.xhi = x end -- works for string or num
66    i.ys:add(y)  end
67
68  function RANGE.__tostring(i)
69    local x, lo, hi = i.ys.txt, i.xlo, i.xhi
70    if    lo == hi  then return fmt("%s==%s",x, lo)
71    elseif hi ==  big then return fmt("%s>%s",x, lo)
72    elseif lo == -big then return fmt("%s<=%s", x, hi)
73    else              return fmt("%s<%s<=%s",lo,x,hi) end end
74  -----------------------------------------------------------------
75  -- ## class SOME
76  function SOME.new(i) i.n,i.t,i.ok=0,{},true end
77  function SOME.has(i) i.t=i.ok and i.t or sort(i.t); i.ok=true; return i.t end
78  function SOME.add(i,x)
79    if x=="?" then return x end
80    i.n=i.n+1
81    if    #i.t  < THE.some    then i.ok=false; push(i.t,x)
82    elseif rand() < THE.some/i.n then i.ok=false; i.t[rand(#i.t)]=x end end
83  -----------------------------------------------------------------
84  -- ## class NUM
85  function NUM.new(i) i.n,i.mu,i.m2,i.w,i.lo,i.hi,i.some=0,0,0,1,big,-big,SOME() end
86  function NUM.mid(i,p)      return rnd(i.mu,p) end
87  function NUM.like(i,x,...) return normpdf(x, i.mu, i.sd) end
88  function NUM.bin(x)
89    b=(i.hi - i.lo)/THE.bins; return i.lo==i.hi and 1 or math.floor(x/b+.5)*b end
90
91  function NUM.add(i_NUM, v_number)
92    if v=="?" then return v end
93    i.some:add(v)
94    i.n = i.n + 1
95    local d  = v - i.mu
96    i.mu = i.mu + d/i.n
97    i.m2 = i.m2 + d*(v - i.mu)
98    i.sd = i.n<2 and 0 or (i.m2/(i.n-1))^0.5
99    i.lo = math.min(v, i.lo)
100   i.hi = math.max(v, i.hi) end
101
102 function NUM.merge(i,j,     k)
103   local k = NUM(i.at, i.txt)
104   for _,n in pairs(i.some.t) do k:add(x) end
105   for _,n in pairs(j.some.t) do k:add(x) end
106   return k end
107 -----------------------------------------------------------------
108 -- ## class SYM
109 function SYM.new(i)          i.n,i.syms,i.most,i.mode = 0,{},0,nil end
110 function SYM.mid(i,...)      return i.mode end
111 function SYM.like(i,x,prior) return ((i.syms[x] or 0)+THE.m*prior)/(i.n+THE.m) end
112 function SYM.bin(i,x)        return x end
113 function SYM.add(i,v,inc)
114   if v=="?" then return v end
115   inc=inc or 1
116   i.n = i.n + inc
117   i.syms[v] = inc + (i.syms[v] or 0)
118   if i.syms[v] > i.most then i.most,i.mode = i.syms[v],v end end
119
120 function SYM.merge(i,j,      k)
121   local k = SYM(i.at, i.txt)
122   for x,n in pairs(i.has) do k:add(x,n) end
123   for x,n in pairs(j.has) do k:add(x,n) end
124   return k end
125 -----------------------------------------------------------------
```

```lua
126 -- ## class COLS
127 local  is={}
128 function is.use(x)     return not x:find"$" end
129 function is.num(x)     return x:find"^[A-Z]" end
130 function is.goal(x)    return x:find"[!+-]$" end
131 function is.klass(x)   return x:find"!$" end
132 function is.dislike(x) return x:find"-$"  end
133
134 function COLS.new(i,t,      new,is)
135   i.xs, i.ys, i.names = {},{},{},t
136   for at,txt in pairs(t) do
137     new = (is.nump(txt) and NUM or SYM)(at,txt)
138     new.usep,  new.w = is.use(txt), is.dislike(txt) and -1 or 1
139     if new.usep then
140       if is.klass(new.txt) then i.klass=new end
141       push(is.goal(new.txt) and i.ys or i.xs, new) end end end
142
143 function COLS.add(i,t)
144   for _,cols in pairs(i.xs,i.ys) do
145     for _,col in pairs(cols) do col:add(t.cells[col.at]) end end
146   return t end
147 -----------------------------------------------------------------
148 -- ## class ROW
149 function ROW.new(i,of,t) i.of,i.cells,i.evaled  = of,t,false end
150 function ROW.klass(i)    return i.cells[i.of.cols.klass.at] end
151 function ROW.within(i,range)
152   local lo, hi, at = range.xlo, range.xhi, range.ys.at
153   local v = i.cells[at]
154   return v=="?" or (lo==hi and v==lo) or (lo<v and v<=hi) end
155 -----------------------------------------------------------------
156 -- ## class ROWS
157
158 -- __ROWS( 't' :[string] )__<br>constructor.
159 function ROWS.new(i,t) i.cols=COLS(t); i.rows={} end
160 -- __.add( 't' :(table|ROW) ) :ROW__<br>update with a table or ROW.
161 function ROWS.add(i,t) return push(i.rows, i.cols:add(t.cells and t or ROW(i,t))) end
162 -- __.mid( 'cols' :[NUM|SYM], 'p'=2 ) :table__
163 -- returns 'mid's of some columns; round numerics to 'p' decimal places.
164 function ROWS.mid(i, cols, p,     t)
165   t={};for _,col in pairs(cols or i.cols.ys) do t[col.txt]=col:mid(p) end;return t end
166 -- __.clone( ?'data' :(table|[ROW]) ) :ROWS__<br>copy this structure, maybe add data.
167 function ROWS.clone(i,data,  j)
168   j= ROWS(i.cols.names);for _,row in pairs(data or {}) do j:add(row) end; return j end
169 -- __.like( 'row' :ROWS, 'nklasses' :int; 'nrows' :int ) :number__
170 -- how likely is it that 'row' could live here?
171 function ROWS.like(i,row, nklasses, nrows,    prior,like,inc,x)
172   prior = (#i.rows + THE.k) / (nrows + THE.k * nklasses)
173   like  = math.log(prior)
174   for _,col in pairs(i.cols.xs) do
175     x = row.cells[col.at]
176     if x and x ~= "?" then
177       inc  = col:like(x,prior)
178       like = like + math.log(inc) end end
179   return like end
180
181 -- __doRows( ?'src' :(string|table), 'fun' :function( table|ROW )   )__
182 -- helper function for reading from tables or files. Case arg1 of ...
183 -- ... _table_ : call function for all items in table.
184 -- ... _string_ : call function on all rows from a file.
185 -- ... _nil_    : call function of all rows from standard input.
186 local function doRows(src, fun)
187   if type(src)=="table" then for  _,t in pairs(src) do fun(t) end
188                         else for  t in csv(src)  do fun(t) end end end
189 -----------------------------------------------------------------
190 -- ## class NB
191 -- (0) Use row1 to initial our 'overall' knowledge of all rows.
192 -- After that (1) add row to 'overall' and (2) ROWS about this row's klass.
193 -- (3) After 'wait' rows, classify row BEFORE updating training knowledge
195 function NB.new(i,src,report,      row)
196   report = report or print
197   i.overall, i.dict, i.list  = nil, {}, {}
198   doRows(src,   function(row,  k)
199     if not i.overall then i.overall = ROWS(row) else  -- (0) eat row1
200       row = i.overall:add(row)             -- add to overall
201       if #i.overall.rows > THE.wait then report(row:klass(), i:guess(row)) end
202       i:train(row) end end) end            -- add tp rows's klass
203
204 function NB.train(i,row,      k)
205   k=row:klass()
206   i.dict[k] = i.dict[k] or push(i.list,i.overall:clone()) -- klass is known
207   i.dict[k].txt = k                        -- each klass knows its name
208   i.dict[k]:add(row) end                   -- update klass with row
209
210 function NB.guess(i,row)
211   return argmax(i.dict,
212     function(klass) return klass:like(row,#i.list,#i.overall.rows) end) end
213 -- function TREE.new(i,listOfRows,gaurd)
214 --   i.gaurd, i.kids = gaurd, {}
215 --   of   = listOfRows[1][1].of
216 --   best = sort(map(of.cols.x,
217 --           function(col) i:bins(col,listOfRows) end),lt"div")[1]
218 --   i.kids = map(best.ranges, function(range)
219 --             listOfRows1 = {}
220 --   -- local function within(row)       return row:within(best) end
221 --   -- local function withins(rows)     return map(rows, within) end
222 --   -- map(listOrRanges, function(rows) return withins(rows) end) end
223 --     tmp= map(rows,withins)
224 --   -- if #tmp > stop then
225 --   -- end)
226 --
227 -----------------------------------------------------------------
228 -- ## class TREE
229 -- function decisionTree(listOfRows)
230 -- -- function tree(rows, xols, yklass,y, gaurd)
231 -- --   local function xranges(xcol) return i:ranges(rows,xcol,yklass,y) end
232 -- --   i.gaurd = gaurd
233 -- --   ranges = sort(map(xcols, xranges),lt"div")[1].ranges
234 -- --   for _,row in pairs(rows) do
235 -- --     for _,range in pairs(ranges) do
236 -- --       if row:within(range) then push(range.rows,row) end; break end end
237 -- --   i.kids = map(ranges,
238 -- --            function(range) return TREE(range.rows,xcols,yklass,y,range) end)
239 --   end
240 -- --   labels , all, xcols = {},{}
241 -- --   for label,rows in pairs(listofRows) do
242 -- --     for _,row in pairs(rows) do
243 -- --       xcols = row.of.cols.xs
244 -- --       labels[ push(all,row).id ] = label end end
245 -- --   return TREE(all, xcols, SYM, function(row) return labels[row.id] end) end
246 --
247 local _ranges, _merge
248 function _ranges(i,rows,xcol,yklass,y)
```

```lua
248 local n,list, dict = 0,{}, {}
250 for _,row in pairs(rows) do
251   local v = row.cells[xcol.at]
252   if v ~= "?" then
253     n = n + 1
254     local pos = xcol:bin(v)
255     dict[pos] = dict[pos] or push(list, RANGE(v,v, yklass(xcol.at, xcol.txt)))
256     dict[pos]:add(v, y(row)) end end
257 list = sort(list, lt"xlo")
258 list = xcol.is=="NUM" and _merge(list, n^THE.min) or list
259 return {ranges  = list,
260         div    = sum(list,function(z) return z.ys:div()*z.ys.n/n end)} end
```

page 3

```lua
function _merge(b4,min)
  local j,t a,b,c,ay,by,cy = 1,{}
  while j <= #b4 do
    a, b = b4[j], b4[j+1]
    if b then
      ay,by,cy = a.ys, b.ys, a.ys:merge(b.ys)
      if ay.n<min or by.n<min or cy:div() <= (ay.n*ay:div()+by.n*by:div())/cy.n
      then a = raNGE(a.xlo, b.xhi, cy)
        j = j + 1 end end -- skip one, since it has just been merged
    t[#t+1] = a
    j = j + 1 end
  if #t < #b4 then return _merge(t,min) end
  for j=2,#t do t[j].xlo = t[j-1].xhi end
  t[1].xlo, t[#t].xhi = -big, big
  return t end
-------------------------------------------------------------------------------
-- ## TESTS
local no,go = {},{}
function go.the()  print(1); print(THE); return type(THE.p)=="number" and THE.p==2 end

function go.argmax(  t,fun)
  fun=function(x) return -x end
  t={50,40,0,40,50}
  return 3 == argmax(t,fun) end

function go.num(n) n=NUM(); for x=1,100 do n:add(x) end; return n.mu==50.5 end

function go.sym(s)
  s=SYM(); for _,x in pairs{"a","a","a","a","b","b","c"} do s:add(x) end
  return s.mode=="a" end

function go.csv(    n,s)
  n,s=0,0; for row in csv(THE.file)  do n=n+1; if n>1 then s=s+row[1] end end
  return rnd(s/n,3) == 5.441  end

function go.rows(    rows)
  doRows(THE.file,function(t) if rows then rows:add(t)  else rows=ROWS(t) end end)
  return rnd(rows.cols.ys[1].sd,0)==847 end

function go.nb()
  return 268 == #NB("../../data/diabetes.csv").dict["positive"].rows  end

local function _classify(file)
  local Abcd=require"abcd"
  local abcd=Abcd()
  NB(file, function(got,want) abcd:add(got,want) end)
  abcd:pretty(abcd:report())
  return true end

function go.soybean() return _classify("../../data/soybean.csv") end
function go.diabetes() return _classify("../../data/diabetes.csv") end
-------------------------------------------------------------------------------
-- ## START
if    pcall(debug.getlocal, 4, 1)
then  return {ROW=ROW, ROWS=ROWS, NUM=NUM, SYM=SYM, THE=THE,lib=lib}
else  THE = cli(THE,help)
      demos(THE,go) end
--         ---------
--        |         |
--     -= |_____| =-
--
--       |___) = (___|
--        ---   ---
--             ###
--          #   =   #    "This ain't chemistry.
--          #######     This is art."
```