```lua
 1  -------------------------------------------------------------------------
 2  --- vim: ts=2 sw=2 et :
 3  local b4,help = {},[[
 4  CHOP: best or rest multi-objective optimization.
 5  (c) 2022 Tim Menzies, timm@ieee.org
 6  "I think the highest and lowest points are the important ones.
 7   Anything else is just...in between." ~ Jim Morrison
 8
 9  USAGE: lua chop.lua [OPTIONS]
10
11  OPTIONS:
12   -m  --min    exponent of min size     = .5
13   -b  --bins   max bins                 = 16
14   -s  --seed   random number seed       = 10019
15   -S  --some   number of nums to keep   = 256
16   -p  --p      exponent of distance     = 2
17
18
19  OPTIONS (other):
20   -f  --file  where to find data        = ../etc/data/auto93.csv
21   -h  --help  show help                 = false
22   -r  --rnd   rounding rules            = %5.2f
23   -g  --go    start up action           = nothing
24
25  Usage of the works is permitted provided that this instrument is
26  retained with the works, so that any entity that uses the works is
27  notified of this instrument. DISCLAIMER:THE WORKS ARE WITHOUT WARRANTY. ]]
28  -------------------------------------------------------------------------
29  -- ## Namespace
30  local the={}
31  local _,big,clone,csv,demos,discretize,dist,eg,entropy,fmt,gap,is,like,lt
32  local map,merge,mid,mode,mu,norm,num,o,oo,pdf,per,push,rand,range
33  local rnd,rnds,rowB4,slice,sort,some,same,sd,string2thing,sym,these
34  local NUM,SYM,RANGE,EGS,COLS,ROW
35  for k,__ in pairs(_ENV) do b4[k]=k end -- At end, use 'b4' to find rogue vars.
36  -------------------------------------------------------------------------
37  -- ## Coding Conventions
38  -- - _Separate policy from mechanism:_
39  --   All "magic parameters" that control code behavior should be part
40  --   of that help text. Allow for '-h' on the command line to print
41  --   help. Parse that string to set the options.
42  -- - _Dialogue independence_: Isolate and separate operating system interaction.
43  -- - _Test-driven development_: The 'go' functions store tests.
44  --   Tests should be silent unless they -- fail. ~tests can be
45  --   disabled by renaming from 'go.fun' to 'no.fun'. Tests should
46  --   return 'true' if the test passes. On exit, return number of
47  --   failed tests.
48  -- - _Less is more:_ Code 80 chars wide, or less. Functions in 1 line,
49  --   if you can. Indent with two spaces. Divide code into 120 line (or
50  --   less) pages. Use 'i' instead of 'self'. Use '_' to denote the
51  --   last created class/ Use '__' for anonymous variable.s Minimize
52  --   use of local (exception: define all functions as local at top of
53  --   file).
54  -- - _Encapsulation:_ Use polymorphism but no inheritance (simpler
55  --   debugging). All classes get a 'new' constructor.
56  --   Use UPPERCASE for class names.
57  --
58  -- ## About the Learning
59  -- - Data is stored in ROWs.
60  -- - Beware missing values (marked in "?") and avoid them
61  -- - Where possible all learning should be  incremental.
62  -- - Standard deviation and entropy generalized to 'div' (diversity);
63  -- - Mean and mode generalized to 'mid' (middle);
64  -- - Rows are created once and shared between different sets of
65  --   examples (so we can accumulate statistics on how we are progressing
66  --   inside each row).
67  -- - When a row is first created, it is assigned to a 'base'; i.e.
68  --   a place to store the 'lo,hi' values for all numerics.
69  -- - XXX tables very sueful
70  -- - XXX table have cols. cols are num, syms. ranges
71
72  -------------------------------------------------------------------------
73  -- ## Utils
74  -- Misc
75  big=math.huge
76  rand=math.random
77  fmt=string.format
78  same = function(x) return x end
79
80  -- Sorting
81  function sort(t,f)      table.sort(#t>0 and t or map(t,same), f); return t end
82  function lt(x)          return function(a,b) return a[x] < b[x] end end
83
84  -- Query and update
85  function map(t,f, u)  u={};for k,v in pairs(t) do u[1+#u]=f(v) end; return u end
86  function push(t,x)      t[1+#t]=x; return x end
87  function slice(t,i,j,k,    u)
88    i,j = (i or 1)//1, (j or #t)//1
89    k   = (k and (j-i)/k or 1)//1
90    u={}; for n=i,j,k do u[1+#u] = t[n] end return u end
91
92  -- "Strings 2 things" coercion.
93  function string2thing(x)
94    x = x:match"^%s*(.-)%s*$"
95    if x=="true" then return true elseif x=="false" then return false end
96    return math.tointeger(x) or tonumber(x) or x   end
97
98  function csv(csvfile)
99    csvfile = io.input(csvfile)
100   return function(line, row)
101     line=io.read()
102     if not line then io.close(csvfile) else
103       row={}; for x in line:gmatch("([^,]+)") do push(row,string2thing(x)) end
104       return row end end end
105
106 -- "Things 2 strings" coercion.
107 function oo(t)  print(o(t)) end
108 function o(t,    u)
109   if #t>0 then return "{"..table.concat(map(t,tostring),", ")..."}" else
110     u={}; for k,v in pairs(t) do u[1+#u] = fmt(":%s %s",k,v) end
111     return (t.is or "").."{"..table.concat(sort(u),", ").."}" end end
112
113 function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
114 function rnd(x,f)
115   return fmt(type(x)=="number" and (x~=x//1 and f or the.rnd) or "%s",x) end
116
117 -- Polymorphic objects.
118 function is(name,    t,new)
119   function new(kl,...)
120     local x=setmetatable({},kl); kl.new(x,...); return x end
121   t = {__tostring=o, is=name or ""}; t.__index=t
122   _ = t
123   return setmetatable(t, {__call=new}) end
124 -------------------------------------------------------------------------
125 -- ## Objects
126
127 -- ### NUM
128 -- - For a stream of 'add'itions, incrementally maintain 'mu,sd'.
129 -- - 'Norm'alize data for distance and discretization calcs
130 --   (see 'dist' and 'range').
131 -- - Comment on 'like'lihood that something belongs to this distribution.
132 NUM="NUM"
133 function _.new(i,at,txt)
134   i.at=at or 0; i.txt=txt or ""; i.lo,i.hi=big, -big
135   i.n,i.mu,i.m2,i.sd = 0,0,0,0,0;  i.w=(txt or""):find"-$" and -1 or 1 end
136
137 function _.add(i,x,   d)
138   if x=="?" then return x end
139   i.n  = i.n + 1
140   d    = x - i.mu
141   i.mu = i.mu + d/i.n
142   i.m2 = i.m2 + d*(x - i.mu)
143   i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n - 1))^0.5)
144   i.lo = math.min(i.lo,x)
145   i.hi = math.max(i.hi,x) end
146
147 function _.range(i,x,n,  b)  b=(i.hi-i.lo)/n;  return math.floor(x/b+0.5)*b end
148 function _.mid(i)  return i.mu end
149
150 function _.norm(i,x)  return i.hi-i.lo<1E-9 and 0 or (x-i.lo)/(i.hi-i.lo+1/big)end
151
152 function _.dist(i,  x,y)
153   if    x=="?" and y=="?"  then return 1 end
154   if    x=="?"             then y = i:norm(y); x = y<.5 and 1 or 0
155   elseif y=="?"            then x = i:norm(x); y = x<.5 and 1 or 0
156   else x,y = i:norm(x), i:norm(y) end
157   return math.abs(x - y) end
158
159 function _.like(i,x,__,          e)
160   return (x < i.mu - 4*i.sd and 0 or x > i.mu + 4*i.sd and 0 or
161     2.7183^(-(x - i.mu)^2 / (z + 2*i.sd^2))/(z + (math.pi*2*i.sd^2).5)) end

162 -------------------------------------------------------------------------
163 -- ### SYM
164 -- - For a stream of 'add'itions, incrementally maintain count of 'all' symbols.
165 -- - Using that info, report 'dist', mode ('mid') symbol, and entropy
166 --    ('div') of this distribution.
167 -- - Comment on 'like'lihood that something belongs to this distribution.
168 -- - Discretization of a symbol just returns that sym ('range').
169 SYM-is"SYM"
170 function _.new(i,at,txt) i.at=at or 0; i.txt=txt or ""; i.n,i.all = 0,{} end
171 function _.add(i,x,n)
172   if x=="?" then return x end
173   n = n or 1
174   i.n=i.n+n; i.all[x] = n + (i.all[x] or 0) end
175
176 function _.range(i,x,__)  return x end
177 function _.dist(i,x,y)  return (a==b and 0 or 1) end
178
179 function _.mid(i)
180   m=0; for y,n in pairs(i.all) do if n>m then m,x=n,y end end; return x end
181
182 function _.div(i,    n,e)
183   e=0; for k,n in pairs(i.all) do e=e-n/i.n*math.log(n/i.n,2) end ;return e end
184
185 function _.like(i,x,prior) return ((c.all[x] or 0) + the.m*prior)/(c.n+the.m) end
186 -------------------------------------------------------------------------
187 -- ### RANGE
188 -- - For a stream of 'add'itions, incrementally maintain counts of 'x' and 'y'.
189 -- - Summarize 'x' as the 'lo,hi' seen so far and summarize 'y' in 'SYM' counts
190 --   in 'y.all' (and get counts there using 'of').
191 -- - Support range sorting ('__lt') and printing ('__tostring').
192 -- - Check if this range's 'x' values 'select's for a particular row.
193 -- - 'Merge' adjacent ranges if the entropy of the whole is less than the parts.
194 RANGE=is"RANGE"
195 function _.new(i,col,lo,hi,y)
196   i.col, i.x, i.y = col, (lo=lo or big, hi=hi or -big), (y or  SYM()) end
197
198 function _.add(i,x,y)
199   if x=="?" then return x end
200   i.x.lo = math.min(i.x.lo,x)
201   i.x.hi = math.max(i.x.hi,x)
202   i.y:add(y) end
203
204 function _.__lt(i,j)  return i.x.lo < j.x.lo end
205 function _.of(i,x)  return i.y.all[x] or 0 end
206
207 function _.selects(i,t,    x)
208   t = t.cells and t.cells or t
209   x = t[i.at]
210   return x=="?" or (i.x.lo==i.x.hi and i.x.lo==x) or (i.x.lo<=x and x<i.x.hi)end
211
212 function _.__tostring(i)
213   local x, lo, hi = i.col.txt, i.x.lo, i.x.hi
214   if    lo ==  hi  then return fmt("%s == %s",x, lo)
215   elseif hi ==  big then return fmt("%s >= %s",x, lo)
216   elseif lo == -big then return fmt("%s < %s", x, hi)
217   else              return fmt("%s <= %s < %s",lo,x,hi) end end
218
219 function _.merge(i,j,n0,    k)
220   k = SYM(i.col.at, i.col.txt)
221   for x,n in pairs(i.y.all) do k:add(x,n) end
222   for x,n in pairs(j.y.all) do k:add(x,n) end
223   if i.y.n<(n0 or 0) or j.y.n<(n0 or 0) or (
224     (i.y:div(i)*i.y.n + j.y:div(j)*j.y.n)/k.n >= .99*k:div())
225   then return RANGE(i.col, i.x.lo, j.x.hi, k) end end
226
227 -- ### ROW
228 -- - Using knowledge 'of' the geometry of the data, support distance calcs
229 -- i ('__sub' and 'around') as well as multi-objective ranking ('__lt').
230 ROW=is"ROWS"
231 function _.new(i,eg, cells) i.of,i.cells = eg,cells end
232 function _.__lt(i,j,       s1,s2,e,y,a,b)
233   y = i.of.cols.y
234   s1, s2, e = 0, 0,  math.exp(1)
235   for __,col in pairs(y) do
236     a = col:norm(i.cells[col.at])
237     b = col:norm(j.cells[col.at])
238     s1 = s1 - e^(col.w * (a - b) / #y)
239     s2 = s2 - e^(col.w * (b - a) / #y) end
240   return s1/#y < s2/#y end
241
242 function _.__sub(i,j)
243   for __,col in pairs(i.of.cols.x) do
244     a,b = i.cells[col.at], j.cells[col.at]
245     inc = a=="?" and b=="?" and 1 or col:dist(a,b)
246     d = d + inc^the.p end
247   return (d / (#i.of.cols.x)) ^ (1/the.p) end
248
249 function _.around(i,rows)
250   return sort(map(rows or i.of.rows, function(j) return {dist=i-j,row=j} end),
251        lt"dist") end
```

```
252  ----------------------------------------------------------------
253  -- ### COLS
254  -- - Factory for converting column 'names' to 'NUM's ad 'SYM's.
255  -- - Store all columns in -- 'all', and for all columns we are not skipping,
256  --   store the independent and dependent columns distributions in 'x' and 'y'.
257  COLS=is"COLS"
258  function _.new(i,names,      head,row,col)
259    i.names=names; i.all={}; i.y={}; i.x={}
260    for at,txt in pairs(names) do
261      col       = push(i.all, (txt:find"^[A-Z]" and NUM or SYM)(at, txt))
262      col.goalp = txt:find"[!+-]$" and true or false
263      if not txt:find"$" then
264        if txt:find"$" then i.klass=col end
265        push(col.goalp and i.y or i.x, col) end end end
266  ----------------------------------------------------------------
267  -- ### EGS
268  -- - For a stream of 'add'itions, incrementally store rows, summarized in 'cols'.
269  -- - When 'add'ing, build new rows for new data. Otherwise reuse rows across
270  --   multiple sets of examples.
271  -- - Supporting 'copy'ing of this structure, without or without rows of data.
272  -- - Report how much this set of examples 'like' a new row.
273  -- - Discretize columns as 'ranges' that distinguish two sets of rows
274  --   (merging irrelevant distinctions).
275  -- - Summarize the 'mid'point of these examples.
276  EGS=is"EGS"
277  function _.new(i,names) i.rows,i.cols = {}, COLS(names) end
278  function _.load(f,  i)
279    for row in csv(the.file) do if i then i:add(row) else i=EGS(row) end end
280    return i end
281
282  function _.add(i,row,    cells)
283    cells = push(i.rows, row.cells and row or ROW(i,row)).cells
284    for n,col in pairs(i.cols.all) do col:add(cells[n]) end end
285
286  function _.mid(i,cols)
287    return map(cols or i.cols.y, function(c) return c:mid() end) end
288
289  function _.copy(i,rows,  j)
290    j=EGS(i.cols.names); for __,r in pairs(rows or {}) do j:add(r) end;return j end
291
292  function _.like(i,t,overall, nHypotheses,    c)
293    prior = (#i.rows + the.k) / (overall + the.k * nHypotheses)
294    like  = math.log(prior)
295    for at,x in pairs(t) do
296      c=i.cols.all.at[at]
297      if x~="?" and not c.goalp then
298        like = math.log(col:like(x)) + like end end
299    return like end
300
301  local _merge, _xpand, _ranges
302  function _.ranges(i,one,two,    t)
303    t={}; for __,c in pairs(i.cols.x) do t[c.at]=_ranges(c,one,two) end;return t end
304
305  function _ranges(col,yes,no,    out,x,bin)
306    out = {}
307    for __,what in pairs({rows=yes, klass=true}, {rows=no, klass=false}) do
308      for __,row in pairs(what.rows) do x = row.cells[col.at]; if x~="?" then
309        bin       = col:range(x,the.bins)
310        out[bin]  = out[bin] or RANGE(col,x,x)
311        out[bin]:add(x, what.klass) end end end
312    return _xpand(_merge(sort(map(out,same)),
313                        .9*(#yes+#no)^the.min)) end
314
315  function _merge(b4,min,        a,b,c,j,n,tmp)
316    j,n,tmp = 1,#b4,{}
317    while j<=n do
318      a, b = b4[j], b4[j+1]
319      if b then c = a:merge(b,min); if c then a,j = c,j+1 end end
320      tmp[#tmp+1] = a
321      j = j+1 end
322    return #tmp==#b4 and tmp or _merge(tmp,min) end
323
324  function _xpand(t)
325    for j=2,#t do t[j].lo=t[j-1].hi end
326    t[1].x.lo, t[#t].x.hi= -big,big
327    return t end
```

```
328  ----------------------------------------------------------------
329  function nasa93dem()
330    local vl,l,n,h,vl,xh=1,2,3,4,5,6; return {
331    {"id:","center","Year","prec","flex","resl","team","pmat","rely","data","cplx",
332                   "ruse","docu","time","stor","pvol","acap","pcap","pcon",
333                   "apex","plex","ltex","tool","site","sced","Kloc",
334                   "Effort-","Defects-","Months-"},
335    {1,2,1979,h,h,h,vh,h,h,h,n,n,n,n,n,n,n,n,n,n,n,n,n,1,25.9,117.6,808,15.3},
336    {2,2,1979,h,h,h,vh,h,h,h,n,n,n,n,n,n,n,n,n,n,n,n,n,1,24.6,117.6,767,15},
337    {3,2,1979,h,h,h,vh,h,h,l,n,n,n,n,n,n,n,n,n,n,n,n,n,1,7.7,31.2,240,10.1},
338    {4,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,n,n,n,n,n,n,n,n,1,8.2,36,256,10.4},
339    {5,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,n,n,n,n,n,n,n,n,1,9.7,25.2,302,11},
340    {6,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,n,n,n,n,n,n,n,n,1,2.2,8.4,69,6.6},
341    {7,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,n,n,n,n,n,n,n,n,1,3.5,10.8,109,7.8},
342    {8,2,1982,h,h,h,vh,h,h,l,n,n,n,n,n,n,n,n,n,n,n,n,n,1,66.6,352.8,2077,21},
343    {9,1,1980,h,h,h,vh,h,h,n,n,n,n,xh,xh,l,h,h,h,h,h,n,n,n,7.5,72,226,13.6},
344    {10,1,1980,h,h,h,vh,n,n,l,h,n,n,n,n,l,h,vh,n,vh,n,n,n,n,20,72,566,14.4},
345    {11,1,1984,h,h,h,vh,n,n,l,h,n,n,n,n,l,h,h,n,vh,n,n,n,n,6,24,188,9.9},
346    {12,1,1980,h,h,h,vh,n,n,l,h,n,n,n,n,l,h,vh,n,vh,n,n,n,n,100,360,2832,25.2},
347    {13,1,1985,h,h,h,vh,n,n,l,h,n,n,n,n,l,h,n,vh,n,l,n,n,n,11.3,36,456,12.8},
348    {14,1,1980,h,h,h,vh,n,n,l,h,n,n,n,n,h,h,h,h,l,vl,n,n,n,100,215,5434,30.1},
349    {15,1,1983,h,h,h,vh,n,n,l,h,n,n,n,n,l,h,h,n,vh,n,n,n,n,20,48,626,15.1},
350    {16,1,1982,h,h,h,vh,n,n,l,h,n,n,n,n,n,n,n,n,vl,n,n,n,100,360,4342,28},
351    {17,1,1980,h,h,h,vh,n,n,l,h,n,n,n,xh,l,h,vh,n,h,n,n,n,150,324,4868,32.5},
352    {18,1,1984,h,h,h,vh,n,n,l,h,n,n,n,n,l,h,h,n,h,n,n,n,31.5,60,986,17.6},
353    {19,1,1983,h,h,h,vh,n,n,l,h,n,n,n,n,h,h,vh,n,h,n,n,n,15.48,470,13.6},
354    {20,1,1984,h,h,h,vh,n,n,l,h,n,n,n,xh,l,h,n,n,h,n,n,n,32.5,60,1276,20.8},
355    {21,2,1985,h,h,h,vh,h,h,l,n,n,n,n,n,n,n,n,n,n,n,n,1,19.7,60,614,13.9},
356    {22,2,1985,h,h,h,vh,h,h,l,h,n,n,n,n,n,n,n,n,n,n,n,1,66.6,300,2077,21},
357    {23,2,1985,h,h,h,vh,h,h,l,h,n,n,n,n,n,n,n,n,n,n,n,1,29.5,120,726,19.5},
358    {24,2,1986,h,h,h,n,n,n,n,n,n,n,n,n,n,n,n,n,n,n,n,n,15,90,575,15.2},
359    {25,2,1985,h,h,h,vh,h,h,n,n,n,n,n,n,n,n,n,n,n,n,n,38,210,1553,21.3},
360    {26,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,n,n,10.48,427,12.4},
361    {27,2,1982,h,h,h,vh,h,n,vh,h,n,n,vh,l,vh,n,n,n,h,l,h,n,n,l,15.4,70,765,14.5},
362    {28,2,1982,h,h,h,vh,h,n,vh,h,n,n,vh,vh,n,n,n,l,n,n,l,48.5,239,2409,21.4},
363    {29,2,1982,h,h,h,vh,n,n,vh,h,n,n,vh,vh,l,vh,n,n,n,l,h,n,n,l,16.3,82,810,14.8},
364    {29,2,1982,h,h,h,vh,h,n,vh,h,n,n,vh,vh,l,vh,n,n,n,l,h,n,n,l,12.8,62,636,13.6},
365    {31,2,1982,h,h,h,vh,n,n,vh,h,n,n,vh,vh,l,vh,n,n,n,l,h,n,n,l,32.6,170,1619,18.7},
366    {32,2,1982,h,h,h,vh,n,n,vh,h,n,n,vh,vh,l,vh,n,n,n,l,h,n,n,l,35.5,192,1763,19.3},
367    {33,2,1985,h,h,h,vh,h,l,h,n,n,n,n,l,n,n,n,n,n,n,1,5.5,18,172,9.1},
368    {34,2,1987,h,h,h,vh,h,l,h,n,n,n,n,l,n,n,n,n,n,n,1,10.4,50,324,11.2},
369    {35,2,1987,h,h,h,vh,h,l,h,n,n,n,n,l,n,n,n,n,n,n,1,14.60,437,12.4},
370    {36,2,1986,h,h,h,n,n,n,h,n,n,n,n,n,n,n,n,n,n,n,n,6.5,42,290,12},
371    {37,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,n,13,60,683,14.8},
372    {38,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,n,90,444,3343,26.7},
373    {39,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,n,8,42,420,12.5},
374    {40,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,n,16,114,887,16.4},
375    {41,2,1980,h,h,h,vh,h,n,h,n,n,n,vh,h,l,h,h,n,n,l,h,n,n,l,177.9,1248,7998,31.5},
376    {42,6,1975,h,h,h,vh,h,l,n,n,n,n,l,n,h,n,n,n,n,n,302,2400,8543,38.4},
377    {43,5,1982,h,h,h,vh,n,n,l,n,n,n,n,h,n,n,n,n,n,n,282.1,1368,9820,37.3},
378    {44,5,1982,h,h,h,vh,h,l,h,n,n,n,n,n,n,n,n,n,n,n,284.7,973,8518,38.1},
379    {45,5,1982,h,h,h,vh,n,h,h,n,n,n,n,n,l,n,n,h,n,n,n,79,400,2327,26.9},
380    {46,5,1977,h,h,h,vh,l,l,n,n,n,n,n,n,l,h,vh,n,l,h,n,n,n,423,2400,18447,41.9},
381    {47,5,1977,h,h,h,vh,h,n,n,n,n,n,n,l,l,h,vh,n,vh,l,h,n,n,n,190,420,5092,30.3},
382    {48,5,1984,h,h,h,vh,n,n,h,n,n,n,h,n,n,n,n,n,n,n,1,47.5,252,2007,22.3},
383    {49,5,1980,h,h,h,vh,l,vh,n,xh,n,n,h,l,n,n,n,n,n,n,1,21,107,1058,21.3},
384    {50,5,1983,h,h,h,vh,l,n,h,h,n,n,vh,n,n,n,n,h,h,n,n,n,78,571.4,4815,30.5},
385    {51,5,1984,h,h,h,vh,l,n,h,h,n,n,vh,n,n,h,h,n,n,n,n,11.4,98.8,704,15.5},
386    {52,5,1985,h,h,h,vh,l,n,h,n,n,vh,n,n,h,h,n,n,n,n,n,19.3,155,1191,18.6},
387    {53,5,1979,h,h,h,vh,l,h,n,vh,n,n,h,l,h,n,n,n,n,n,n,101,750,4840,32.4},
388    {54,5,1979,h,h,h,vh,l,h,n,n,n,n,h,h,n,n,n,n,n,n,n,219,2120,11761,42.8},
389    {55,5,1979,h,h,h,vh,l,h,n,n,h,n,n,h,l,n,n,n,n,n,n,50,370,2685,25.4},
390    {56,2,1979,h,h,h,vh,h,vh,h,n,n,vh,vh,n,vh,n,vh,n,vh,n,h,n,n,1,227,1181,6293,33},
391    {57,2,1977,h,h,h,vh,h,n,h,n,n,vh,vh,h,n,n,n,n,l,n,n,n,1,70,278,2950,20.2},
392    {58,2,1979,h,h,h,vh,h,l,h,n,n,n,n,l,n,n,n,n,n,n,1,0.9,8.4,28,4.9},
393    {59,6,1974,h,h,h,vh,l,vh,l,xh,n,n,xh,vh,l,h,n,vh,vl,h,n,n,n,980,4560,50961,96},
394    {60,6,1975,h,h,h,vh,n,n,l,h,n,n,n,n,l,vh,vh,n,n,h,n,n,n,350,720,8447,39.5},
395    {61,5,1976,h,h,h,vh,h,n,xh,n,n,h,h,l,h,n,n,h,h,n,n,70,458,2404,27.5},
396    {62,5,1976,h,h,h,vh,h,n,xh,n,n,h,l,h,n,n,n,h,h,n,n,271,2460,9308,43.4},
397    {63,5,1971,h,h,h,vh,h,n,n,n,n,n,n,n,n,n,n,n,n,n,90,162,2743,25},
398    {64,5,1980,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,40,150,1219,18.9},
399    {65,5,1979,h,h,h,vh,l,l,n,n,n,n,n,n,n,n,n,n,n,n,137,636,4210,32.2},
400    {66,5,1977,h,h,h,vh,h,n,n,h,n,n,n,h,h,h,n,n,n,n,n,150,882,5848,36.2},
401    {67,5,1976,h,h,h,vh,n,vh,n,n,n,n,h,l,h,h,n,n,n,n,n,339,444,8477,45.9},
402    {68,5,1983,h,h,h,n,n,l,h,l,n,n,n,n,n,n,n,n,n,n,240,192,10313,37.1},
403    {69,5,1978,h,h,h,vh,l,h,n,n,n,n,vh,l,h,h,h,h,h,n,n,n,1,144,576,6129,28.8},
404    {70,5,1979,h,h,h,vh,l,n,l,n,n,n,n,vh,l,h,h,h,h,n,n,n,1,151,432,6136,26.2},
405    {71,5,1979,h,h,h,vh,l,n,l,h,n,n,n,vh,l,h,h,h,h,n,n,n,1,34,72,1555,16.2},
406    {72,5,1979,h,h,h,vh,l,n,n,n,n,n,vh,l,h,h,h,h,n,n,n,1,98,300,4907,24.4},
407    {73,5,1979,h,h,h,vh,l,n,l,n,n,n,vh,l,h,h,h,h,n,n,n,1,85,300,4256,23.2},
408    {74,5,1982,h,h,h,vh,l,n,l,n,n,n,vh,l,h,h,h,h,n,n,n,1,20,240,813,12.8},
409    {75,5,1978,h,h,h,vh,l,n,l,n,n,n,vh,l,h,h,h,h,n,n,n,1,111,600,4511,23.5},
410    {76,5,1978,h,h,h,vh,l,h,vh,h,n,n,n,n,h,h,h,h,n,n,n,1,162,756,7553,32.4},
411    {77,5,1978,h,h,h,vh,l,h,vh,n,n,n,vh,l,h,h,h,h,n,n,n,1,352,1200,17597,42.9},
412    {78,5,1979,h,h,h,vh,l,h,n,vh,n,n,n,n,h,h,h,h,n,n,n,1,165,97,7867,31.5},
413    {79,5,1984,h,h,h,vh,h,h,n,vh,n,n,h,l,h,n,n,n,h,h,n,n,60,409,2004,24.9},
414    {80,5,1984,h,h,h,vh,h,h,n,vh,n,n,h,h,n,n,n,n,n,100,703,3340,29.6},
415    {81,2,1980,h,h,h,vh,h,n,vh,vh,n,n,xh,xh,h,n,n,h,n,n,l,l,n,n,32,1350,2984,33.6},
416    {82,2,1980,h,h,h,vh,h,h,h,h,n,n,vh,xh,h,h,h,h,h,n,n,n,53,480,2227,28.8},
417    {83,3,1977,h,h,h,vh,h,h,l,vh,n,n,vh,xh,l,vh,vh,n,vh,vl,h,n,n,41,599,1594,23},
418    {84,3,1977,h,h,h,vh,h,h,l,vh,n,n,vh,xh,l,vh,vh,n,vh,vl,vl,h,n,n,24,430,933,19.2},
419    {85,5,1977,h,h,h,vh,h,vh,h,n,n,n,xh,xh,n,h,h,h,h,h,n,n,n,165,4178.2,6266,47},
420    {86,5,1977,h,h,h,vh,h,vh,h,n,n,n,xh,xh,n,h,h,h,h,n,n,n,65,1772.5,2468,34.5},
421    {87,5,1977,h,h,h,vh,h,vh,h,n,n,n,xh,xh,n,h,h,h,h,n,n,n,70,1645.9,2658,35.4},
422    {88,5,1977,h,h,h,vh,h,vh,h,n,n,n,xh,xh,n,h,h,h,h,n,n,n,50,1924.5,2102,34.2},
423    {89,5,1982,h,h,h,vh,l,vh,l,vh,n,n,vh,xh,l,h,n,n,l,vl,l,h,n,n,7.25,648,406,15.6},
424    {90,5,1980,h,h,h,vh,h,vh,h,vh,n,n,xh,xh,n,h,h,h,h,n,n,n,233,8211,8848,53.1},
425    {91,2,1983,h,h,h,vh,n,h,n,vh,n,n,vh,n,n,n,n,n,l,l,n,n,n,16.3,480,1253,21.5},
426    {92,2,1983,h,h,h,vh,n,h,n,vh,n,n,vh,vh,h,n,n,n,n,l,l,n,n,n,6.2,12,477,15.4},
427    {93,2,1983,h,h,h,vh,n,h,n,vh,n,n,vh,vh,h,n,n,n,n,l,l,n,n,n,3,38,231,12}} end
```

```
428  ----------------------------------------------------------------
429  -- ## DEMOS
430  local go,no={},{}
431
432  -- Convert help string to a table. Check command line for any updates.
433  function these(f1,f2,k,x)
434    for n,flag in ipairs(arg) do if flag==f1 or flag==f2 then
435      x = x=="false" and"true" or x=="true" and"false" or arg[n+1] end end
436    the[k] = string2thing(x) end
437
438  -- Run the demos, resetting settings and random number see before each.
439  -- Return number of failures.
440  function demos(    fails,names,defaults,status)
441    fails=0     -- this code will return number of failures
442    names, defaults = {},{}
443    for k,f in pairs(go) do if type(f)=="function" then push(names,k) end end
444    for k,v in pairs(the) do defaults[k]=v end
445    for __,one in pairs(sort(names)) do    -- for all we want to do
446      for k,v in pairs(defaults) do the[k]=v end   -- set settings to defaults
447      math.randomseed(the.seed or 10019)           -- reset random number seed
448      io.stderr:write(".")
449      status = go[one]()                           -- run demo
450      if status ~= true then
451        print("-- Error",one,status)
452        fails = fails + 1 end end                   -- update fails
453    return fails end                               -- return total failure count
454
455
456  -- Simple stuff
457  function go.the()       return type(the.bins)=="number" end
458  function go.sort( t) return 0==sort({100,3,4,2,10,0})[1] end
459  function go.slice( t,u)
460    t = {10,20,30,40,50,60,70,80,90,100,110,120,130,140}
461    u = slice(t,3,#t,3)
462    t = slice(t,3,5)
463    return #t==3 and #u==4 end
464
465  function go.num(       n,mu,sd)
466    n, mu, sd = NUM(), 10, 1
467    for i=1,10^4 do
468      n:add(mu+sd*math.sqrt(-2*math.log(rand()))*math.cos(2*math.pi*rand())) end
469    return math.abs(n.mu - mu) < 0.05 and math.abs(n.sd - sd) < 0.5 end
470
471  -- Can we read rows off the disk?
472  function go.rows( n,m)
473    m,n=0,0; for row in csv(the.file) do m=m+1; n=n+#row; end; return n/m==8 end
474
475  -- Can we turn a list of names into columns?
476  function go.cols(  i)
477    i=COLS("name","Age","ShoeSize")
478    return i.y[1].w == -1 end
479
480  -- Can we read data, summarized as columns?
481  function go.egs( it)
482    it = EGS.load(the.file); return math.abs(2970 - it.cols.y[1].mu) < 1 end
483
484  -- Can we discretize
485  function go.ranges( it,n,best,rest,min)
486    it = EGS.load(the.file)
487    print("all",o(rnds(it:mid())))
488    it.rows = sort(it.rows)
489    for j,row in pairs(sort(it.rows)) do row.klass = 1+j//(#it.rows*.35/6) end
490    n = (#it.rows)^.5
491    best,rest = slice(it.rows,1,n), slice(it.rows, n+1, #it.rows, 3*n)
492    print("best",#best,o(rnds(it:copy(best):mid())))
493    print("rest",#rest,o(rnds(it:copy(rest):mid())))
494    for __,ranges in pairs(it:ranges(best,rest)) do
495      print""
496      for at,range in pairs(ranges) do
497        print(range,o(range.y.all)) end end
498    --oo(a:mid())
499    --oo(b:mid())
500    return math.abs(2970 - it.cols.y[1].mu) < 1 end
```

```
501  -------------------------------------------------------------------------     527
502  -- ## Main                                                                    528
503                                                                                529
504  -- - Parse help text for flags and defaults, check CLI for updates.           530
505  -- - Maybe print the help (with some pretty colors).
506  -- - Run the demos.
507  -- - Check for rogue vars.
508  -- - Exit, reporting number of failures.
509  help:gsub("\n ([-][^%s]+)[%s]+([-][-]([^%s]+))[^\n]*%s([^%s]+)",these)
510  if the.help then
511    print(help:gsub("%u%u+", "\27[31m%1\27[0m")
512               :gsub("(%s)([-][-]?[^%s]+)(%s)","%1\27[33m%2\27[0m%3"),"")
513  else
514    local fails = demos()
515    for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end  end
516    os.exit(fails) end
517  --- function SOME() return {all={}, ok=false, n=0} end
518  --- function some(i,x)
519  ---   if x=="?" then return x end
520  ---   i.n = 1 + i.n
521  ---   if     #i.all < the.some     then i.ok=false; push(i.all, x)
522  ---   elseif rand() < the.some/i.n then i.ok=false; i.all[rand(#i.all)]=x end end
523  ---
524  --- function per(i,p)
525  ---   i.all = i.ok and i.all or sort(i.all); i.ok=true
526  ---   return i.all[math.max(1, math.min(#i.all, (p or .5)*#i.all//1))] end
```