

```

1  --- vim: ts=2 sw=2 et :
2  local b4,help = {},{}
3  CHOP: best or rest multi-objective optimization.
4  (c) 2022 Tim Menzies, tim@ieee.org
5  "I think the highest and lowest points are the most important ones.
6  Anything else is just...in between." - Jim Morrison
7
8  USAGE: lua chop.lua [OPTIONS]
9
10 OPTIONS:
11 -how --how good or bad or novel = good
12 -m --min exponent of min size = .5
13 -b --bins max bins = 16
14 -s --seed random number seed = 10019
15 -S --size number of num to keep = 256
16 -p --p exponent of distance = 2
17
18 OPTIONS (other):
19 -f --file where to find data = ../etc/data/auto93.csv
20 -h --help show help = false
21 -r --rrnd rounding rules = %5.2f
22 -g --go start up action = nothing
23
24 Usage of the works is permitted provided that this instrument is
25 retained with the works, so that any entity that uses the works is
26 notified of this instrument. DISCLAIMER:THE WORKS ARE WITHOUT WARRANTY. ]]
27
28 --- ## Namespace
29 local the={}
30 local big,copy,csv,demos,discretize,dist,eg,entropy,fill_in,the,fmt,gap,is,like,lt
31 local map,merge,mid,mode,mu,nasa93dem,notrm,num,o,oo,pdf,per,push,rand,range
32 local rnd,rnds,rowB4,slice,sort,some,same,sd,string2thing,sym
33 local NUM,SYM,RANGE,EGS,COLS,ROW
34 for k, _ in pairs(_ENV) do b4[k]=k end -- At end, use 'b4' to find rogue vars.
35
36 --- ## Coding Conventions
37 --- Separate policy from mechanism:
38 --- All "magic parameters" that control code behavior should be part
39 of that help text. Allow for 'h' on the command line to print
40 help. Parse that string to set the options.
41 --- Dialogue independence.: Isolate and separate operating system interaction.
42 --- Test-driven development.: The 'go' functions store tests.
43 --- Tests should be silent unless they -- fail. -tests can be
44 disabled by renaming from 'go.fun' to 'no.fun'. Tests should
45 return 'true' if the test passes. On exit, return number of
46 failed tests.
47 --- Code as specification.: "I have made this letter longer since I
48 did not have time to make it shorter". Good code is short code.
49 --- Lots of short functions. Methods listed alphabetically.
50 --- Code 80 chars wide, or less. Functions in 1 line,
51 if you can. Indent with two spaces. Divide code into 120 line (or
52 less) pages. Use 'i' instead of 'and'.
53 --- Minimize use of local (exception: define all functions
54 local at top of file).
55 --- Encapsulation: Use polymorphism but no inheritance (simpler
56 debugging). All classes get a 'new' constructor.
57 --- Use UPPERCASE for class names.
58
59 --- ## About the Learning
60 --- Data is stored in ROWs.
61 --- Beware missing values (marked in "?") and avoid them
62 --- Where possible all learning should be incremental.
63 --- Standard deviation and entropy generalized to 'div' (diversity);
64 --- Mean and mode generalized to 'mid' (middle);
65 --- Rows are created once and shared between different sets of
66 examples (so we can accumulate statistics on how we are progressing
67 inside each row).
68 --- When a row is first created, it is assigned to a 'base'; i.e.
69 a place to store the 'lo,hi' values for all numerics.
70 --- XXX tables very useful.
71 --- XXX table have cols. cols are num, syms. ranges
72

```

```

73
74
75 function nasa93dem()
76 local vl,1,n,h,vh,xhsl,2,3,4,5,6; return {
77 ["id","center","Year","prec","flex","resl","team","pmat","rely","data","cpix",
78 "ruse","docu","time","stor","pvol","acap","pcap","pcon",
79 "apex","plex","flex","tool","site","sced","kloc",
80 "effort","Defects","Months"],
81
82 {1,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,1,25,9,117,6,808,15.3},
83 {2,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,24,6,117,6,767,15},
84 {3,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,7,7,31,2,240,10.1},
85 {4,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,8,2,36,256,10.4},
86 {5,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,9,7,25,2,302,11},
87 {6,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,2,2,8,4,69,6.6},
88 {7,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,3,5,10,8,109,7.8},
89 {8,2,1982,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,66,6,352,8,2077,21},
90 {9,1,1980,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,7,5,72,2,240,15.6},
91 {10,1,1980,h,h,h,vh,n,n,1,h,n,n,n,n,1,h,vh,n,vh,n,n,n,20,72,566,14.4},
92 {11,1,1984,h,h,h,vh,n,n,1,h,n,n,n,n,1,h,h,n,vh,n,h,n,n,6,24,188,9.9},
93 {12,1,1980,h,h,h,vh,n,n,1,h,n,n,n,n,1,h,vh,n,vh,n,n,n,100,360,2832,25.2},
94 {13,1,1985,h,h,h,vh,n,n,1,h,n,n,n,n,1,h,n,n,n,1,n,n,n,1,3,36,456,12.8},
95 {14,1,1980,h,h,h,vh,n,n,1,h,n,n,n,n,1,h,n,n,n,100,215,5438,30.1},
96 {15,1,1983,h,h,h,vh,n,n,1,h,n,n,n,n,1,h,h,n,vh,n,h,n,n,20,48,626,15.1},
97 {16,1,1982,h,h,h,vh,n,n,1,h,n,n,n,n,1,h,n,n,n,vl,n,n,n,100,360,4342,28},
98 {17,1,1980,h,h,h,vh,n,n,1,h,n,n,n,xh,l,h,vh,n,vh,n,h,n,n,150,324,4868,32.5},
99 {18,1,1984,h,h,h,vh,n,n,1,h,n,n,n,n,1,h,n,n,n,h,n,n,n,31,5,60,986,17.6},
100 {19,1,1983,h,h,h,vh,n,n,1,h,n,n,n,n,1,h,h,n,vh,n,h,n,n,15,48,470,13.6},
101 {20,1,1984,h,h,h,vh,n,n,1,h,n,n,n,xh,l,h,n,n,n,h,n,h,n,n,32,5,60,1276,20.8},
102 {21,2,1985,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,13,7,6,14,13,9},
103 {22,2,1985,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,66,6,352,8,2077,21},
104 {23,2,1985,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,29,5,120,920,16},
105 {24,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,h,n,n,n,h,n,n,n,15,90,575,15.2},
106 {25,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,h,n,n,n,h,n,n,n,38,210,1553,21.3},
107 {26,2,1982,h,h,h,vh,h,h,l,h,n,n,n,n,1,h,n,n,n,h,n,n,n,10,48,470,13.6},
108 {27,2,1982,h,h,h,vh,n,n,vh,h,n,n,vh,vh,l,vh,n,n,h,l,h,n,n,1,15,4,70,765,14.5},
109 {28,2,1982,h,h,h,vh,h,n,vh,h,n,n,vh,vh,l,vh,n,n,h,l,h,n,n,1,48,5,239,2409,21.4},
110 {29,2,1982,h,h,h,vh,h,n,vh,h,n,n,vh,vh,l,vh,n,n,h,l,h,n,n,1,16,3,82,810,14.8},
111 {30,2,1982,h,h,h,vh,h,n,vh,h,n,n,vh,vh,l,vh,n,n,h,l,h,n,n,1,12,3,68,680,14.8},
112 {31,2,1982,h,h,h,vh,h,n,vh,h,n,n,vh,vh,l,vh,n,n,h,l,h,n,n,1,32,6,170,619,18.7},
113 {32,2,1982,h,h,h,vh,h,n,vh,h,n,n,vh,vh,l,vh,n,n,h,l,h,n,n,1,35,5,192,1763,19.3},
114 {33,2,1985,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,5,5,18,172,9.1},
115 {34,2,1987,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,10,4,324,11.2},
116 {35,2,1987,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,14,60,457,12.4},
117 {36,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,6,5,42,290,12},
118 {37,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,13,60,683,14.8},
119 {38,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,90,444,3343,26.7},
120 {39,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,42,42,420,22.3},
121 {40,2,1986,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,16,114,887,16.4},
122 {41,2,1980,h,h,h,vh,h,n,h,n,h,n,n,vh,h,l,h,n,n,n,1,177,9,1248,7998,31.5},
123 {42,6,1975,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,h,n,n,n,n,n,302,2400,8543,38.4},
124 {43,6,1975,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,h,n,n,n,n,n,180,420,9040,39.3},
125 {44,5,1982,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,h,n,n,n,n,n,284,7,973,8518,38.1},
126 {45,5,1982,h,h,h,vh,n,n,n,n,n,n,n,n,n,n,n,n,n,n,79,400,2327,26.9},
127 {46,5,1977,h,h,h,vh,l,1,n,n,n,n,n,n,1,h,vh,n,h,n,h,n,n,423,2400,18447,41.9},
128 {47,5,1977,h,h,h,vh,l,1,n,n,n,n,n,n,1,h,vh,n,h,n,h,n,n,180,420,9040,39.3},
129 {48,5,1984,h,h,h,vh,h,n,n,n,n,n,n,n,n,h,n,h,n,n,n,47,5,252,2007,22.3},
130 {49,5,1980,h,h,h,vh,l,1,vh,n,xh,n,n,h,h,l,n,n,n,h,n,n,n,21,107,1058,21.3},
131 {50,5,1983,h,h,h,vh,l,1,n,h,n,n,n,vh,n,n,n,h,n,h,n,n,n,78,571,4,4815,30.5},
132 {51,5,1978,h,h,h,vh,l,1,n,h,n,n,n,vh,n,n,n,h,n,h,n,n,n,1,4,98,9,705,15.6},
133 {52,5,1985,h,h,h,vh,n,n,n,n,n,vh,n,n,n,h,n,h,n,n,n,19,3,155,1191,18.6},
134 {53,5,1979,h,h,h,vh,l,1,n,n,vh,n,n,h,h,l,h,n,n,n,h,h,n,n,101,750,4840,32.4},
135 {54,5,1979,h,h,h,vh,l,1,n,n,n,n,n,h,h,l,n,n,n,h,n,n,n,219,2120,11761,42.8},
136 {55,5,1979,h,h,h,vh,l,1,n,n,n,n,n,h,h,l,n,n,n,h,n,n,n,50,370,2685,25.4},
137 {56,2,1979,h,h,h,vh,h,h,n,n,n,n,vh,vh,l,vh,n,n,vh,h,n,n,1,227,1181,6293,33},
138 {57,2,1977,h,h,h,vh,h,n,h,n,h,vh,n,n,n,1,h,vh,n,n,1,70,278,2950,20.2},
139 {58,2,1979,h,h,h,vh,h,h,l,h,n,n,n,n,1,n,n,n,n,h,n,n,1,0.9,8,4,28,4.9},
140 {59,6,1974,h,h,h,vh,l,1,vh,l,xh,n,n,xh,vh,l,h,h,n,n,vh,vl,h,n,n,n,980,4560,50961,96},
141 {60,6,1975,h,h,h,vh,n,n,1,h,n,n,n,n,1,vh,vh,n,n,h,n,n,n,n,350,720,8547,35.7},
142 {61,5,1976,h,h,h,vh,h,h,n,xh,n,n,h,h,l,h,n,n,n,h,h,h,n,70,458,2404,27.5},
143 {62,5,1979,h,h,h,vh,h,h,n,xh,n,n,h,h,l,h,n,n,n,h,h,h,n,271,2460,9308,43.4},
144 {63,5,1971,h,h,h,vh,n,n,n,n,n,n,n,n,1,h,h,n,h,n,n,n,n,90,162,2743,25},
145 {64,5,1980,h,h,h,vh,n,n,n,n,n,n,n,n,1,h,h,n,h,n,n,n,n,40,150,1219,18.9},
146 {65,5,1979,h,h,h,vh,n,n,n,n,n,n,n,n,1,h,h,n,h,n,n,n,n,137,636,4210,32.2},
147 {66,5,1977,h,h,h,vh,n,n,n,n,n,n,n,n,h,h,h,h,n,n,n,n,150,882,5848,36.2},
148 {67,5,1976,h,h,h,vh,n,n,n,n,n,n,n,n,1,h,h,n,h,n,n,n,n,339,444,8477,45.9},
149 {68,5,1983,h,h,h,vh,n,l,1,h,n,n,n,n,h,h,h,n,h,n,n,n,240,192,10313,37.1},
150 {69,5,1978,h,h,h,vh,l,1,n,n,n,n,n,n,vh,l,1,h,n,h,h,h,n,n,1,144,576,6129,28.8},
151 {70,5,1979,h,h,h,vh,l,1,n,n,n,n,n,n,vh,l,1,h,n,h,h,h,n,n,1,151,432,6136,26.2},
152 {71,5,1979,h,h,h,vh,l,1,n,h,n,n,n,n,n,1,h,n,n,h,h,h,n,n,1,34,72,1555,16.2},
153 {72,5,1979,h,h,h,vh,l,1,n,n,n,n,n,n,vh,l,1,h,n,n,h,h,h,n,n,1,98,300,4907,24.4},
154 {73,5,1979,h,h,h,vh,l,1,n,n,n,n,n,n,vh,l,1,h,n,h,h,h,n,n,1,85,300,4256,32.2},
155 {74,5,1982,h,h,h,vh,l,1,n,1,n,n,n,n,vh,l,1,h,n,n,h,h,h,n,n,1,20,240,813,12.8},
156 {75,5,1978,h,h,h,vh,l,1,n,1,n,n,n,n,vh,l,1,h,n,n,h,h,h,n,n,1,111,600,4511,23.5},
157 {76,5,1978,h,h,h,vh,l,1,h,vh,n,n,n,vh,l,1,h,n,h,h,h,n,n,1,162,756,7553,32.4},
158 {77,5,1978,h,h,h,vh,l,1,h,vh,n,n,n,vh,l,1,h,n,h,h,h,n,n,1,352,1200,17597,42.9},
159 {78,5,1979,h,h,h,vh,l,1,h,n,vh,n,n,n,vh,l,1,h,n,h,h,h,n,n,1,165,97,7867,31.5},
160 {79,5,1984,h,h,h,vh,h,h,n,vh,n,n,n,h,h,l,h,n,n,n,h,h,n,n,60,409,2004,24.9},
161 {80,5,1984,h,h,h,vh,h,h,n,vh,n,n,n,h,h,l,h,n,n,n,h,h,n,n,100,703,3340,29.6},
162 {81,2,1980,h,h,h,vh,h,h,n,vh,n,n,n,xh,xh,h,n,n,n,32,1350,2984,33.6},
163 {82,2,1980,h,h,h,vh,h,h,n,vh,n,n,n,vh,xh,h,h,n,h,h,n,n,53,480,2227,28.8},
164 {83,3,1977,h,h,h,vh,h,h,l,1,vh,n,n,vh,xh,l,vh,vh,vl,vl,h,n,n,41,599,1594,23},
165 {84,3,1977,h,h,h,vh,h,h,l,1,vh,n,n,vh,xh,l,vh,vh,vh,vl,vl,h,n,n,24,430,933,19.2},
166 {85,5,1977,h,h,h,vh,h,h,n,xh,n,n,xh,xh,n,h,h,n,n,165,78,78,64,405,15.6},
167 {86,5,1977,h,h,h,vh,h,h,n,xh,n,n,xh,xh,n,h,h,n,n,65,1772,5,2468,34.5},
168 {87,5,1977,h,h,h,vh,h,h,n,xh,n,n,xh,xh,n,h,h,n,n,70,1645,9,2658,35.4},
169 {88,5,1977,h,h,h,vh,h,h,n,xh,n,n,xh,xh,n,h,h,n,n,50,1924,5,2102,34.2},
170 {89,5,1977,h,h,h,vh,h,h,n,xh,n,n,xh,xh,n,h,h,n,n,25,64,405,15.6},
171 {90,5,1980,h,h,h,vh,h,h,n,vh,n,n,n,xh,xh,n,h,h,n,n,233,8211,8848,53.1},
172 {91,2,1983,h,h,h,vh,n,h,n,vh,n,n,vh,vh,n,n,n,1,1,n,n,16,3,480,1253,21.5},
173 {92,2,1983,h,h,h,vh,n,h,n,vh,n,n,vh,vh,n,n,n,1,1,n,n,6,2,12,477,15.4},
174 {93,2,1983,h,h,h,vh,n,h,n,vh,n,n,vh,vh,n,n,n,1,1,n,n,3,38,231,121} end

```

```

175 --- ## Utils
176 --- Misc
177 big=math.huge
178 rand=math.random
179 fmt=string.format
180 same = function(x) return x end
181
182 --- Sorting
183 function sort(t,f) table.sort(#t>0 and t or map(t,same), f); return t end
184 function lt(x) return function(a,b) return a[x] < b[x] end end
185
186 --- Query and update
187 function map(t,f, u) u={} ;for k,v in pairs(t) do u[l+#u]=f(v) end; return u end
188 function push(t,x) t[l+#t]=x; return x end
189 function slice(t,i,j,k,u) u={} ;for k,v in pairs(t) do if i<=k and k<=j then u[k]=v end end
190 i,j = (i or 1)//1, (j or l)//1
191 k = (k and (j-i)/k or 1) //1
192 u={} ;for n=i,j,k do u[l+#u] = t[n] end return u end
193
194 --- "Strings 2 things" coercion.
195 function string2thing(x)
196 x = x:match("%s*(-)%s%*")
197 if x=="true" then return true elseif x=="false" then return false end
198 return math.tointeger(x) or tonumber(x) or x end
199
200 function csv(csvfile)
201 csvfile = io.input(csvfile)
202 return function(line, row)
203 line=io.read()
204 if not line then io.close(csvfile) else
205 row={} ;for x in line:gmatch("%([^\n])") do push(row,string2thing(x)) end
206 return row end end
207
208 --- "Things 2 strings" coercion.
209 function oo(t) print(o(t)) end
210 function o(t)
211 if #t>0 then return {""..table.concat(map(t,tostring),"")..""} else
212 u={} ;for k,v in pairs(t) do u[l+#u] = fmt("%s%s",k,v) end
213 return (t.is or "").."{""..table.concat(sort(u),"")..""}" end end
214
215 function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
216 function rnd(x,f)
217 return fmt(type(x)=="number" and (x-x)//1 and f or the.rnd) or "%s",x) end
218
219 --- Convert help string to a table. Check command line for any updates.
220 function fill_in_the(shortflag,longflag,slo,t,x)
221 for n,flag in ipairs(arg) do
222 if flag==shortflag or flag==longflag then
223 if flag=="true" and "true" or arg[n+1] end end
224 the[slot] = string2thing(x) end
225
226 --- Run demos, each time resetting settings and random seed. Return #failures.
227 local go,no=1,{} -- place to store enabled and disabled tests
228 function demos( fails,names,defaults,status)
229 fails=0 -- this code will return number of failures
230 names, defaults = {},{}
231 for k,f in pairs(go) do if type(f)=="function" then push(names,k) end end
232 for k,v in pairs(the) do defaults[k]=v end
233 if go[the.go] then names[the.go] end
234 for _,one in pairs(sort(names)) do
235 for k,v in pairs(defaults) do defaults[k]=v end -- for all we want to do
236 math.randomseed(the.seed or 10019) -- set settings to defaults
237 -- reset random number seed
238 io.stderr:write("\n")
239 status = go[one]()
240 if status == true then
241 print("== Enro, one,status)
242 fails = fails + 1 end end
243 return fails end
244
245 --- Polymorphic objects.
246 function is(name, t,new)
247 function new(kl,...)
248 local x=metatable({},{kl}; kl.new(x,...)) return x end
249 t = {__tostring=, __name= or ""; t.__index=
250 return setmetatable(t, {__call=new}) end

```

```

251 -----
252 -- ## Objects
253 COLS,EGS,NUM,RANGE,ROW,SYM=is"COLS",is"EGS",is"NUM",is"RANGE",is"SYM",is"ROW"
254
255 -- ### NUM
256 -- - For a stream of 'add'itions, incrementally maintain 'mu,sd'.
257 -- - 'Norm'alize data for distance and discretization calcs
258 -- (see 'dist' and 'range').
259 -- - Comment on 'like'lihood that something belongs to this distribution.
260 function NUM.new(i,at,txt)
261   i.at=at or 0; i.txt=txt or ""; i.lo,i.hi=big, -big
262   i.n,i.mu,i.sd,i.sd = 0,0,0,0; i.w=(txt or ""):find"--$ and -1 or 1 end
263
264 function NUM.add(i,x, d)
265   if x=="?" then return x end
266   i.n = i.n + 1
267   d = x - i.mu
268   i.mu = i.mu + d/i.n
269   i.m2 = i.m2 + d*(x - i.mu)
270   i.sd = (i.m2/0 or i.n-2) and 0 or ((i.m2/(i.n - 1))^0.5)
271   i.lo = math.min(i.lo,x)
272   i.hi = math.max(i.hi,x) end
273
274 function NUM.dist(i,x,y)
275   if x=="?" and y=="?" then return 1 end
276   if x=="?" then y = i:norm(y); x = y<.5 and 1 or 0
277   elseif y=="?" then x = i:norm(x); y = x<.5 and 1 or 0
278   else x,y = i:norm(x), i:norm(y) end
279   return math.abs(x - y) end
280
281 function NUM.like(i,x,_, e)
282   return (x < i.mu - 4*i.sd and 0 or x > i.mu + 4*i.sd and 0 or
283     2.7183*(-(x - i.mu)^2 / (z + 2*i.sd^2))/(z + (math.pi^2*i.sd^2)^.5)) end
284
285 function NUM.merge(i,ranges,min, a,b,c,j,n,tmp)
286 function expand(t)
287   if t<2 then return {} end
288   for j=2,#t do t[j].lo=t[j-1].hi end
289   t[1].x.lo, t[#t].x.hi = -big,big
290   return t
291 end
292 j,n,tmp = 1,ranges,{}
293 while j<#n do
294   a, b = ranges[j], ranges[j+1]
295   if b then c = a:merge(b,min); if c then a,j = c,j+1 end end
296   tmp[#tmp+1] = a
297   j = j+1 end
298   return #tmp==#ranges and expand(tmp) or i:merge(tmp,min) end
299
300 function NUM.mid(i) return i.mu end
301
302 function NUM.norm(i,x)
303   return i.hi-i.lo<1E-9 and 0 or (x-i.lo)/(i.hi-i.lo+1/big) end
304
305 function NUM.range(i,x,n, b) b=(i.hi-i.lo)/n; return math.floor(x/b+0.5)*b end
306 -----
307 -- ### SYM
308 -- - For a stream of 'add'itions, incrementally maintain count of 'all' symbols.
309 -- - Using that info, report 'dist', mode ('mid') symbol, and entropy
310 -- ('div') of this distribution.
311 -- - Comment on 'like'lihood that something belongs to this distribution.
312 -- Discretization of a symbol just returns that sym ('range').
313 function SYM.new(i,at,txt) i.at=at or 0; i.txt=txt or ""; i.n,i.all = 0,{} end
314
315 function SYM.add(i,x,n)
316   if x=="?" then return x end
317   n = n or 1
318   i.n=i.n+n; i.all[x] = n + (i.all[x] or 0) end
319
320 function SYM.dist(i,x,y) return (a==b and 0 or 1) end
321
322 function SYM.div(i, n,e)
323   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
324
325 function SYM.like(i,x,prior) return ((c.all[x] or 0)+the.m*prior)/(c.n+the.m) end
326
327 function SYM.merge(i,ranges,min) return ranges end
328
329 function SYM.mid(i)
330   m=0; for y,n in pairs(i.all) do if n>m then m,x=n,y end end; return x end
331
332 function SYM.range(i,x,_) return x end

```

```

333 -----
334 -- ### RANGE
335 -- - For a stream of 'add'itions, incrementally maintain counts of 'x' and 'y'.
336 -- - Summarize 'x' as the 'lo,hi' seen so far and summarize 'y' in 'SYM' counts
337 -- in 'y.all' (and get counts there using 'of').
338 -- - Support range sorting ('_lt') and printing ('__tostring').
339 -- - Check if this range's 'x' values 'select's for a particular row.
340 -- - Merge 'adjacent' ranges if the entropy of the whole is less than the parts.
341 function RANGE.new(i,col,lo,hi,y)
342   i.col, i.x, i.y = col, (lo==0 or big, hi==hi or -big), (y or SYM()) end
343
344 function RANGE.__lt(i,j) return i.x.lo < j.x.lo end
345
346 function RANGE.__tostring(i)
347   local x, lo, hi = i.col.txt, i.x.lo, i.x.hi
348   if lo == hi then return fmt("%s==%s",x, lo)
349   elseif hi == big then return fmt("%s>=%s",x, lo)
350   elseif lo == -big then return fmt("%s<=%s", x, hi)
351   else return fmt("%s<=%s<%s",lo,x,hi) end end
352
353 function RANGE.add(i,x,y)
354   if x=="?" then return x end
355   i.x.lo = math.min(i.x.lo,x)
356   i.x.hi = math.max(i.x.hi,x)
357   i.y:add(y) end
358
359 function RANGE.merge(i,j,n0, k)
360   k = SYM(i.col.at, i.col.txt)
361   for x,n in pairs(i.y.all) do k:add(x,n) end
362   for x,n in pairs(j.y.all) do k:add(x,n) end
363   if i.y.n<(n0 or 0) or j.y.n<(n0 or 0) or (
364     (i.y:div(i)*i.y.n + j.y:div(j)*j.y.n)/k.n >= .99*k:div())
365   then return RANGE(i.col, i.x.lo, j.x.hi, k) end end
366
367 function RANGE.of(i,x) return i.y.all[x] or 0 end
368
369 function RANGE.score(i,goal,B,R, how)
370   how={}
371   how.good= function(b,r) return ((b<r or b+r < .05) and 0) or b^2/(b+r) end
372   how.bad= function(b,r) return ((r<b or b+r < .05) and 0) or r^2/(b+r) end
373   how.noel=function(b,r) return 1/(b+r) end
374   b, r, z = 0, 0, 1/big
375   for x,n in pairs(i.y.all) do
376     if x==goal then b = b+n else r=r+n end end
377   return how[the.how or "good"] (b/(B+z), z/(R+z)) end
378
379 function RANGE.selects(i,t, x)
380   t = t.cells and t.cells or t
381   x = t[i.at]
382   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
383 -----
384 -- ### ROW
385 -- - Using knowledge 'of' the geometry of the data, support distance calcs
386 -- i ('_sub' and '_around') as well as multi-objective ranking ('__lt').
387 function ROW.new(i,eg, cells) i.of,i.cells = eg,cells end
388
389 function ROW.__lt(i,j, sl,s2,e,y,a,b)
390   y = i.of.cols.y
391   sl, s2, e = 0, 0, math.exp(1)
392   for _,col in pairs(y) do
393     a = col:norm(i.cells[col.at])
394     b = col:norm(j.cells[col.at])
395     s1 = s1 - e^(col.w * (a - b) / #y)
396     s2 = s2 - e^(col.w * (b - a) / #y) end
397   return sl/#y < s2/#y end
398
399 function ROW.__sub(i,j)
400   for _,col in pairs(i.of.cols.x) do
401     a,b = i.cells[col.at], j.cells[col.at]
402     inc = a=="?" and b=="?" and 1 or col:dist(a,b)
403     d = d + inc*the.p end
404   return (d / (#i.of.cols.x)) ^ (1/the.p) end
405
406 function ROW.around(i,rows)
407   return sort(map(rows or i.of.rows, function(j) return (dist=i-j,row=j) end),
408     i.txt") end

```

```

409 -----
410 -- ### COLS
411 -- - Factory for converting column 'names' to 'NUM's ad 'SYM's.
412 -- - Store all columns in -- 'all', and for all columns we are not skipping,
413 -- store the independent and dependent columns distributions in 'x' and 'y'.
414 function COLS.new(i,names, head,row,col)
415   i.names=names; i.all={}; i.y={}; i.x={}
416   for at,txt in pairs(names) do
417     col = push(i.all, (txt:find"^[A-Z]" and NUM or SYM) (at, txt))
418     col.goalp = txt:find"^[+]=S" and true or false
419     if not txt:find"$" then
420       if txt:find"$" then i.klass=col end
421       push(col.goalp and i.y or i.x, col) end end end
422 -----
423 -- ### EGS
424 -- - For a stream of 'add'itions, incrementally store rows, summarized in 'cols'.
425 -- - When 'adding', build new rows for new data. Otherwise reuse rows across
426 -- multiple sets of examples.
427 -- - Supporting 'copy'ing of this structure, without or without rows of data.
428 -- - Replaces how much this set of examples 'like' a new row.
429 -- - Discretize columns as 'ranges' that distinguish two sets of rows
430 -- (merging irrelevant distinctions).
431 -- - Summarize the 'mid'point of these examples.
432 function EGS.new(i,names) i.rows,i.cols = {}, COLS(names) end
433
434 function EGS.add(i,row, cells)
435   cells = push(i.rows, row.cells and row or ROW(i,row)).cells
436   for n,col in pairs(i.cols.all) do col:add(cells[n]) end end
437
438 function EGS.copy(i,rows, j)
439   j=EGS(i.cols.names); for _,r in pairs(rows or {}) do j:add(r) end;return j end
440
441 function EGS.like(i,t,overall, nHypotheses, c)
442   prior = (#i.rows + the.k) / (overall + the.k * nHypotheses)
443   like = math.log(prior)
444   for at,x in pairs(t) do
445     col=i.cols.all[at]
446     if x=="?" and not c.goalp then
447       like = math.log(col:like(x)) + like end end
448   return like end
449
450 function EGS.load(src, i)
451   if src=nil or type(src)=="string"
452   then for row in csv(src) do if i then i:add(row) else i=EGS(row) end end
453   else for _,row in pairs(src) do if i then i:add(row) else i=EGS(row) end end end
454   return i end
455
456 function EGS.mid(i,cols)
457   return map(cols or i.cols.y, function(c) return c:mid(i) end) end
458
459 function EGS.ranges(i,yes,no, out,x,bin,tmp,score)
460   out={}
461   for _,col in pairs(i.cols.x) do
462     tmp = {}
463     for _,what in pairs((rows=yes, klass=true), (rows=no, klass=false)) do
464       for _,row in pairs(what.rows) do x = row.cells[col.at]
465         if x=="?" then
466           bin = col:range(x,the.bins)
467           tmp[bin] = tmp[bin] or RANGE(col,x,x)
468           tmp[bin]:add(x, what.klass) end end end
469       tmp = map(tmp,same) -- a hack. makes tmp sortable (has consecutive indexes)
470       for _,range in pairs(col:merge(sort(tmp), (#yes+#no)*the.min)) do
471         push(out,range) end end
472     score = function(range) return range:score(true,#yes,#no) end
473     return sort(out,score) end

```

```

474 -----
475 -- ## Code for tests and demos
476
477 -- Simple stuff
478 function go.the()      return type(the.bins)=="number" end
479 function go.sort( t)  return 0==sort((100,3,4,2,10,0))[1] end
480 function go.slice( t,u)
481   t = {10,20,30,40,50,60,70,80,90,100,110,120,130,140}
482   u = slice(t,3,#t,3)
483   t = slice(t,3,5)
484   return #t==3 and #u==4 end
485
486 function go.num(      n,mu,sd)
487   n, mu, sd = NUM(), 10, 1
488   for i=1,10^4 do
489     n=add(mu*sd*math.sqrt(-2*math.log(rand()))*math.cos(2*math.pi*rand())) end
490     return math.abs(n.mu - mu) < 0.05 and math.abs(n.sd - sd) < 0.5 end
491
492 -- Can we read rows off the disk?
493 function go.rows( n,m)
494   m,n=0,0; for row in csv(the.file) do m=m+1; n=n+#row; end; return n/m==8 end
495
496 -- Can we turn a list of names into columns?
497 function go.cols( i)
498   i=COLS("Name","Age","ShoeSize-")
499   return i.y[1].w == -1 end
500
501 -- Can we read data, summarized as columns?
502 function go.egs( it)
503   it=EGS.load(nasa93dem())
504   return math.abs(it.cols.y[1].mu - 624) < 1 end
505   --for _,row in pairs(nasa93dem())do oo(row) end end
506   --it = EGS.load(the.file); return math.abs(2970 - it.cols.y[1].mu) < 1 end
507
508 -- Does discretization work?
509 function go.ranges( it,n,best,rest,min)
510   --it = EGS.load(the.file)
511   print(the.how)
512   it=EGS.load(nasa93dem())
513   print("all",o(rnds(it:mid()))))
514   it.rows = sort(it.rows)
515   for j,row in pairs(sort(it.rows)) do row.klass = 1+j//(#it.rows*.35/6) end
516   n = (#it.rows)^.5
517   best,rest = slice(it.rows,1,n), slice(it.rows, n+1, #it.rows, 3*n)
518   print("best",#best,o(rnds(it:copy(best):mid()))))
519   print("rest",#rest,o(rnds(it:copy(rest):mid()))))
520   tmp={}; for _,ranges in pairs(it:ranges(best,rest)) do
521     for at,range in pairs(ranges) do
522       push(tmp,range).val= range:score(true,#best,#rest) end end
523   for _,range in pairs(sort(tmp,lt*val")) do print(range.val, range) end
524   --oo(a:mid())
525   --oo(b:mid())
526   return math.abs(2970 - it.cols.y[1].mu) < 1 end

```

```

527 -----
528 -- ## Main
529
530 -- - Parse help text for flags and defaults, check CLI for updates.
531 -- - Maybe print the help (with some pretty colors).
532 -- - Run the demos.
533 -- - Check for rogue vars.
534 -- - Exit, reporting number of failures.
535 help:gsub("un ([-|^%s+)([%s]+)([-|~|(|(%s+)|^n)*%s([^\s+)*",fill_in_the)
536 if the.help then
537   print(help:gsub("%u%u+", "%127[31m%127[0m")
538         :gsub("(%s)([-|~|?][^\s+)(%s)", "%127[33m%2127[0m%3")", ""))
539 else
540   local fails = demos()
541   for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
542   os.exit(fails) end

```