

```

1 require"lib"
2 -- modules start with an Upper case letter
3 -- class methods are in Module.UPPERCASE (e.g. Module.NEW for constructors)
4 -- instance methods are in Module.method(i,...)
5 -----
6 the={ min = .5,
7       bins = 16,
8       some = 256,
9       seed = 10019,
10      file = "../data/aut093.csv"}
11
12 -----
13 -- data model
14 local Is={}
15 function Is.GOAL(x) return (x or ""):find("[4-5]" end
16 function Is.NUM(x) return (x or ""):find("[A-Z]" end
17 function Is.KLASS(x) return (x or ""):find"$" end
18 function Is.SKIP(x) return (x or ""):find"$" end
19 function Is.WEIGHT(x) return (x or ""):find"$" and -1 or 1 end
20
21 -----
22 local Col={}
23 function Col.NEW(at,txt)
24   return {n = 0,
25          at = at or 0,
26          txt = txt or "",
27          ok = false,
28          log = {},
29          div = 0,
30          mid = 0} end
31
32 function Col.NUM(at,txt)
33   return adds(Col.NEW(at,txt),{
34     nump=true,
35     w = Is.WEIGHT(txt),
36     lo = big,
37     hi = -big}) end
38
39 function Col.add(i,v,inc)
40   inc = inc or 1
41   if v == ""
42   then i.n = i.n + r
43       if i.nump
44       then for _,l,inc do
45         i.lo = math.min(v, i.lo)
46         i.hi = math.max(v, i.hi)
47         if #i.log < the.some then i.ok=false; push(i.log,v)
48         elseif R() < the.some/i.n then i.ok=false; i.log[ R(#i.log) ]=v end
49       end -- end for
50       else i.ok = false
51       i.log[v] = inc + (i.log[v] or 0) end end
52   return i end
53
54 function Col.ok(i)
55   if not i.ok then
56     i.ok = true
57     i.div, i.mid = 0, 0
58     if i.nump
59     then i.log = sort(i.log)
60         i.mid = per(i.log, .5)
61         i.div = (per(i.log, .9) - per(i.log, .1)) / 2.56
62     else local most = -1
63         for x,n in pairs(i.log) do
64           if n > most then most, i.mid = n, x end
65           i.div = i.div - n/i.n * math.log( n/i.n, 2) end end end end
66
67 function Col.div(i) Col.ok(i); return i.div end
68 function Col.mid(i) Col.ok(i); return i.mid end
69 function Col.mids(i,t)
70   t={};for _,c in pairs(i.y) do t[c.txt]=Col.mid(c) end;return t end
71
72 function Col.norm(i,x) return i.hi-i.lo < 1E-9 and 0 or (x-i.lo)/(i.hi-i.lo) end
73
74 function Col.bin(i,x)
75   if i.nump
76   then b=(i.hi - i.lo)/the.bins; return i.lo==i.hi and 1 or math.floor(x/b+.5)*b
77   else return x end end
78 -----
79 local Row={}
80 function Row.NEW(of,cells) return {of=of,cells=cells,evald=false} end
81
82 function Row.better(i,j)
83   local s1, s2, n = 0, 0, #i.of.y
84   for _,c in pairs(i.of.y) do
85     local x,y = Col.norm(c, i.cells[c.at]), Col.norm(c, j.cells[c.at])
86     s1 = s1 - 2.7183^(c.w * (x-y)/n)
87     s2 = s2 - 2.7183^(c.w * (y-x)/n) end
88   return s1/n < s2/n end
89
90 -----
91 local Data={}
92 function Data.READ(src,fun)
93   if type(src)=="table" then for _,t in pairs(src) do fun(t) end
94   else for t in csv(src) do fun(t) end end end
95
96 function Data.NEW(names)
97   local i={x={}, y={}, rows={}, names=names,klass=nil}
98   for at,txt in pairs(names) do
99     local new = Is.NUM(txt) and Col.NUM(at,txt) or Col.NEW(at,txt)
100    if not Is.SKIP(txt) then
101      push(Is.GOAL(txt) and i.y or i.x, new)
102      if Is.KLASS(txt) then i.klass=new end end end
103   return i end
104
105 function Data.clone(i,init, j)
106   j=Data.NEW(i.names)
107   for _,t in pairs(init or {}) do Data.add(j,t) end; return j end
108
109 function Data.add(i,t)
110   t = t.cells and t or Row.NEW(i,t)
111   push(i.rows, t)
112   for _,cols in pairs(i.x, i.y) do
113     for _,c in pairs(cols) do Col.add(c, t.cells[c.at]) end end end
114
115 -----
116 local Bin={}
117 function Bin.new(xlo, xhi, ys) return {lo=xlo, hi=yhi, ys=ys} end
118 function Bin.add(i,x,y)
119   i.lo = math.min(i.lo, x)
120   i.hi = math.max(i.hi, x)
121   Col.add(i.ys, y) end
122
123 function Bin.merge(i,j, min)
124   local k = Col.NEW(i.at, i.txt)
125   for x,n in pairs(i.ys.log) do Col.add(k,x,n) end

```

```

126 for x,n in pairs(j.ys.log) do Col.add(k,x,n) end
127 if i.n<min or j.n<min or div(k) <= (i.n*div(i) + j.n*div(j)) / k.n then
128   return {lo=i.lo, hi=j.hi, ys=k} end end
129
130 function Bin.RANGES(listOfRows,col,y)
131   local n,list, dict = 0, {}, {}
132   for label,rows in pairs(listOfRows) do
133     for _,row in pairs(rows) do
134       local v = row[col.at]
135       if v ~= "" then
136         n = n + 1
137         local pos = Col.bin(col,v)
138         dict[pos] = dict[pos] or push(list, Bin.new(v,v,Col.new(col.at,col.txt)))
139         Bin.add(dict[pos], v, label) end end end
140   list = sort(list, function(a,b) return a.lo < b.lo end)
141   list = col.nump and Bin.MERGES(list, n*the.min) or list
142   return {bins= list,
143         div = sum(list,function(z) return div(z.ys)*z.ys.n/n end)} end
144
145 function Bin.MERGES(b4, min)
146   local j,now = 1, {}
147   while j <= #b4 do
148     local merged = j<#b4 and Bin.merge(b4[j], b4[j+1], min)
149     now[#now+1] = merged or b4[j]
150     j = j + (merged and 2 or 1) end
151   if #now < #b4
152   then return Bin.MERGES(now,min)
153   else for j=2,#now do now[j].lo = now[j-1].hi end
154         now[1].lo, now[#t].hi = -big, big
155         return now end end
156
157 -----
158 -- test suite
159 Go,No = {}, {}
160
161 function Go.THE() oo(the) end
162
163 function Go.ROWS( i)
164   i=rows(the.file)
165   oo(i.x[1]) end
166
167 function Go.STATS()
168   oo(summarize(rows(the.file) ))
169   end
170
171 function Go.ORDER( i,t)
172   i= rows(the.file)
173   t= orders(i, i.xy)
174   left = clone(i,splice(i.xy,1,30))
175   right= clone(i,splice(i.xy,360))
176   print("first",o(mids(left)))
177   print("last", o(mids(right)))
178   print("all", o(mids(i)))
179   end
180
181 math.randomseed(the.seed)
182 if arg[1]=="-g" and type(Go[arg[2]])=="function" then Go[arg[2]]() end

```