

```

1 require"lib"
2
3 -----
4 -- config
5 the={ min = .5,
6       bins = 16,
7       some = 256,
8       seed = 10019,
9       file = "../data/auto93.csv"}
10
11 -----
12 -- data model
13 function goalp(x) return (x or ""):find("[+-]$") end
14 function nump(x) return (x or ""):find("[A-Z]") end
15 function klassp(x) return (x or ""):find"$" end
16 function skip(x) return (x or ""):find"$" end
17 function weight(x) return (x or ""):find"$" and -1 or 1 end
18
19 -----
20 -- col create
21 function col(at,txt)
22     return
23     {
24         n = at or 0,
25         txt = txt or "",
26         v = weight(txt),
27         ok = false,
28         log = {},
29         div = 0,
30         mid = 0
31     } end
32
33 function num(at,txt)
34     local i = col(at,txt)
35     i.nump = true
36     i.v = i.txt:find"$" and -1 or 1
37     i.lo = big
38     i.hi = -big
39     return i end
40
41 -- col update
42 function cell(i,v,r)
43     r = r or 1
44     if v ~= "?"
45     then i.n = i.n + r
46         if i.nump
47         then i.lo = math.min(v, i.lo)
48             i.hi = math.max(v, i.hi)
49             if #i.log < the.some then i.ok=false; push(i.log,v)
50                 elseif R(i) < the.some/i.n then i.ok=false; i.log[ R(#i.log) ]=v end
51             else r = r or 1
52                 i.ok = false
53                 i.log[v] = r + (i.log[v] or 0) end end
54         return i end
55     end
56
57 function ok(i)
58     if not i.ok then
59         i.div, i.mid = 0, 0
60         if i.nump
61         then i.log = sort(i.log)
62             i.mid = per(i.log, .5)
63             i.div = (per(i.log, .9) - per(i.log, .1)) / 2.56
64             local most = -1
65             for x,n in pairs(i.log) do if v>0 then
66                 if n > most then most, i.mid = n, x end
67                 i.div = i.div - n/i.n * math.log( n/i.n, 2) end end end end
68             i.ok = true
69             return i end
70     end
71
72 -- col query
73 function norm(i,x)
74     return i.hi - i.lo < 1E-9 and 0 or (x-i.lo)/(i.hi-i.lo) end
75
76 -----
77 -- data create
78 function data(names)
79     local i={x={}, y={}, xy={}, names=names,klass=nil}
80     for at,txt in pairs(names) do
81         local new = txt:find("[A-Z]") and num(at,txt) or col(at,txt)
82         if not skip(txt) then
83             push(goalp(txt) and i.y or i.x, new)
84             if klassp(txt) then i.klass=new end end end
85         return i end
86     end
87
88 function rows(src)
89     local i
90     if type(src)=="table" then for _,t in pairs(src) do i=row(t,i) end
91         else for t in csv(src) do i=row(t,i) end end
92     return i end
93
94 function clone(i,init, j)
95     j=row(i.names); for _,t in pairs(init or {}) do j=row(t,j) end; return j end
96
97 function row(t,i)
98     if not i then return data(t) end
99     push(i.xy, t)
100     for _,cols in pairs(i.x, i.y) do
101         for _,c in pairs(cols) do cell(c, t[c.at]) end end end
102
103 -- data query
104 function div(i) if not i.ok then ok(i) end; return i.div end
105 function mid(i) if not i.ok then ok(i) end; return i.mid end
106 function mids(i, t) t={};for _,c in pairs(i.y) do t[c.txt]=mid(c) end;return t end
107 function divs(i, t) t={};for _,c in pairs(i.y) do t[c.txt]=div(c) end;return t end
108 function bin(i,x)
109     if i.nump
110     then b=(i.hi - i.lo)/the.bins; return i.lo==i.hi and 1 or math.floor(x/b+.5)*b
111     else return x end end
112
113 -----
114 -- row sort
115 function orders(i,t)
116     local function first(t1,t2)
117         local s1, s2, n, e = 0, 0, #i.y, math.exp(1)
118         for _,c in pairs(i.y) do
119             local x,y = norm(c, t1[c.at]), norm(c, t2[c.at])
120             s1 = s1 - e^(c.w * (x-y)/n)
121             s2 = s2 - e^(c.w * (y-x)/n) end
122         return s1/n < s2/n
123     end
124     return sort(t, first) end
125
126 -----
127 -- discretization
128 function ranges(listOfRows,xcol,yklass,y)
129     local n,list,dict = 0,{}, {}
130
131     for label,rows in pairs(listOfRows) do
132         for _,row in pairs(rows) do
133             local v = row[xcol.at]
134             if v ~= "?" then
135                 n = n + 1
136                 local pos = bin(v)
137                 dict[pos] = dict[pos] or push(list, {lo=v,hi=v,ys=col(xcol.at, xcol.txt)})
138                 dict[pos].lo = math.min(v, dict[pos].lo)
139                 dict[pos].hi = math.max(v, dict[pos].hi)
140                 cell(dict[pos].ys, label)
141             end end end
142         list = sort(list, lt"lo")
143         list = xcol.nump and _xpad(_merges(list, n^the.min)) or list
144         return (ranges= list,
145             div = sum(list,function(z) return div(z.ys)*z.ys.n/n end)) end
146
147 function merge(i,j, min)
148     local k = col(i.at, i.txt)
149     for x,n in pairs(i.ys.log) do cell(k,x,n) end
150     for x,n in pairs(j.ys.log) do cell(k,x,n) end
151     if i.n<min or j.n<min or div(k) <= (i.n*div(i) + j.n*div(j)) / k.n then
152         return {lo=i.lo, hi=j.hi, ys=k} end end
153
154 function _merges(b4, min)
155     local j,now = 1,{}
156     while j <= #b4 do
157         local merged
158         if j<#b4 then merged = merge(b4[j], b4[j+1], min) end
159         now[#now+1] = merged and merged or b4[j]
160         j = j + (merged and 2 or 1) end-- skip to next and next if we found a me
161     rge
162     return #now == #b4 and now or _merges(now,min) end -- hunt for more merges
163
164 function _xpad(ranges)
165     for j=2,#ranges do ranges[j].lo = ranges[j-1].hi end
166     ranges[1].lo, ranges[#t].hi = -big, big
167     return ranges end
168
169 -----
170 -- test suite
171 go,no = {},{}
172
173 function go.the() oo(the) end
174
175 function go.rows( i)
176     i=rows(the.file)
177     oo(i.x[1]) end
178
179 function go.stats()
180     oo(summarize(rows(the.file) ))
181     end
182
183 function go.order( i,t)
184     i= rows(the.file)
185     t= orders(i, i.xy)
186     left = clone(i, splice(i.xy,1,30))
187     right= clone(i, splice(i.xy,360))
188     print("fin", o(mids(left)))
189     print("last", o(mids(right)))
190     print("all", o(mids(i)))
191     end
192
193 math.randomseed(the.seed)
194 if arg[1]=="-s" and type(go[arg[2]])=="function" then go[arg[2]]() end

```