```lua
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   -- don't say self, say "i" (shorter)
6   -- where p-- osible, if looking at two instances, use "i,j"
7   -- types = int,real,str,tab
8   ------------------------------------------------------------------------------
9   local the={ min  = .5,
10              bins = 16,
11              some = 256,
12              seed = 10019,
13              file = "../../data/auto93.csv"}
14
15  ------------------------------------------------------------------------------
16  local Col={}
17  function Col.GOAL(x)   return (x or ""):find"[!+-]$" end
18  function Col.NUMP(x)   return (x or ""):find"^[A-Z]" end
19  function Col.KLASS(x)  return (x or ""):find"!$"  end
20  function Col.SKIP(x)   return (x or ""):find"$"  end
21  function Col.WEIGHT(x) return (x or ""):find"-$" and -1 or 1 end
22
23  --> .COLS(names:[str]) :COLS -> constructor
24  function Col.COLS(names)
25    local i={x={}, y={}, names=names, klass=nil}
26    for at,txt in pairs(names) do
27      local new = Col.NUMP(txt) and Col.NUM(at,txt) or Col.NEW(at,txt)
28      if not Col.SKIP(txt) then
29        push(Col.GOAL(txt) and i.y or i.x, new)
30        if Col.KLASS(txt) then i.klass=new end end end
31    return i end
32
33  --> .NEW(at:?int, txt:?str) :COL -> constructor of numbers
34  function Col.NEW(at,txt)
35    return {n  =0,     at=at or 0, txt=txt or "",
36            ok =false, kept={},
37            div=0,     mid=0} end
38
39  function Col.NUM(at,txt,some)
40    i    = Col.NEW(at,txt)
41    i.w  = Col.WEIGHT(txt)
42    i.nums= some or the.some -- if non-nil the i.nums is a numeric
43    return i end
44
45  function Col.add(i,v,inc)
46    inc = inc or 1
47    if   v ~= "?"
48    then i.n = i.n + inc
49        if i.nums
50        then for _=1,inc do
51               if        #i.kept < i.nums then i.ok=false;push(i.kept,v)
52               elseif R() < i.nums/i.n then i.ok=false;i.kept[R(#i.kept)]=v end end
53        else i.ok = false
54             i.kept[v] = inc + (i.kept[v] or 0) end end
55    return i end
56
57  function Col.ok(i)
58    if not i.ok then
59      i.div, i.mid = 0, 0
60      if   i.nums
61      then i.kept = sort(i.kept)
62           i.mid  = per(i.kept, .5)
63           i.div  = (per(i.kept, .9) - per(i.kept, .1)) / 2.56
64      else local most = -1
65           for x,n in pairs(i.kept) do
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
68    i.ok = true end
69
70  function Col.lo(i)  Col.ok(i); return i.kept[1] end
71  function Col.hi(i)  Col.ok(i); return i.kept[#i.kept] end
72  function Col.div(i) Col.ok(i); return i.div end
73  function Col.mid(i) Col.ok(i); return i.mid end
74  function Col.norm(i,x)
75    local lo,hi = Col.lo(i), Col.hi(i)
76    return hi-lo < 1E-9 and 0 or (x-lo)/(hi-lo) end
77
78  function Col.bin(i,x)
79    if   i.nums then
80      local lo,hi = Col.lo(i), Col.hi(i)
81      local b=(hi - lo)/the.bins
82      x = lo==hi and 1 or math.floor(x/b+.5)*b end
83    return x end
84  ------------------------------------------------------------------------------
85  local Row={}
86  function Row.NEW(of,cells) return {of=of,cells=cells,evaled=false} end
87
88  function Row.better(i,j)
89    local s1, s2, n = 0, 0, #i.of.y
90    for _,c in pairs(i.of.y) do
91      local x,y =  i.cells[c.at], j.cells[c.at]
92      x,y = Col.norm(c, x), Col.norm(c, y)
93      s1  = s1 - 2.7183^(c.w * (x-y)/n)
94      s2  = s2 - 2.7183^(c.w * (y-x)/n) end
95    return s1/n < s2/n  end
96
97  ------------------------------------------------------------------------------
98  local Data={}
99  function Data.NEW(t) return {rows={}, cols=Col.COLS(t)} end
100
101 function Data.ROWS(src,fun)
102   if type(src)=="table" then for  _,t in pairs(src) do fun(t) end
103                         else for   t in csv(src)   do fun(t) end end end
104
105 function Data.clone(i,inits,   j)
106   j=Data.NEW(i.names)
107   for _,t in pairs(inits 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.cols.x, i.cols.y) do
113     for _,c in pairs(cols) do Col.add(c, t.cells[c.at]) end end end
114
115 function Data.mids(i,cols, t)
116   t={}
117   for _,c in pairs(cols or i.cols.y) do t[c.txt] = Col.mid(c)  end;return t end
118 ------------------------------------------------------------------------------
119
120 local Bin={}
121 function Bin.new(xlo, xhi, ys) return {lo=xlo, hi=yhi, ys=ys} end
122 function Bin.add(i,x,y)
123   i.lo = math.min(i.lo, x)
124   i.hi = math.max(i.hi, x)
125   Col.add(i.ys, y) end
```

```lua
126 ------------------------------------------------------------------------------
127 function Bin.merge(i,j, min)
128   local k = Col.NEW(i.at, i.txt)
129   for x,n in pairs(i.ys.kept) do Col.add(k,x,n) end
130   for x,n in pairs(j.ys.kept) do Col.add(k,x,n) end
131   if i.n<min or j.n<min or Col.div(k) <= (i.n*Col.div(i) + j.n*Col.div(j)) / k.n
132   then return {lo=i.lo, hi=j.hi, ys=k} end end
133
134 function Bin.BINS(listOfRows,col,y)
135   local n,list, dict = 0,{}, {}
136   for label,rows in pairs(listOfRows) do
137     for _,row in pairs(rows) do
138       local v = row[col.at]
139       if v ~= "?" then
140         n = n + 1
141         local pos = Col.bin(col,v)
142         dict[pos] = dict[pos] or push(list, Bin.new(v,v,Col.new(col.at,col.txt)))
143         Bin.add(dict[pos], v, label) end end end
144   list = sort(list, lt"lo")
145   list = col.nums and Bin.MERGES(list, n^the.min) or list
146   return {bins= list,
147           div = sum(list,function(z) return Col.div(z.ys)*z.ys.n/n end)} end
148
149 function Bin.MERGES(b4, min)
150   local j,now = 1,{}
151   while j <= #b4 do
152     local merged = j<#b4 and Bin.merge(b4[j], b4[j+1], min)
153     now[#now+1]  = merged or b4[j]
154     j            = j + (merged and 2 or 1)  end
155   if   #now < #b4
156   then return Bin.MERGES(now,min) -- loop to look for other merges
157   else -- stretch the bins to cover minus infinity to plus Infinity
158     for j=2,#now do now[j].lo = now[j-1].hi end
159     now[1].lo, now[#now].hi = -big, big
160     return now end end
161 ------------------------------------------------------------------------------
162 Go,No = {},{}
163
164
165 function Go.THE() oo(the) end
166
167 function Go.ROWS(  d)
168   Data.ROWS(the.file,function(row)
169     if not d then d=Data.NEW(row) else
170        Data.add(d,row) end end)
171   oo(Data.mids(d)) end
172
173 function Go.STATS()
174   oo(summarize(rows(the.file) ))
175 end
176
177 function Go.ORDER(  i,t)
178   i= rows(the.file)
179   t= orders(i, i.xy)
180   left = clone(i,splice(i.xy,1,30))
181   right= clone(i,splice(i.xy,360))
182   print("first",o(mids(left)))
183   print("last", o(mids(right)))
184   print("all",  o(mids(i)))
185   end
186 ------------------------------------------------------------------------------
187 math.randomseed(the.seed)
188 if arg[1]=="-g" and type(Go[arg[2]])=="function" then Go[arg[2]]() end
```