

```

1  USAGE: lua eg.lua [OPTIONS]

2  OPTIONS:
3  -b --bins      number of bins      = 8
4  -c --cohen     small effect         = .35
5  -e --eg       start-up example     = nothing
6  -F --far       far away            = .95
7  -f --file      file with csv data   = ../docs/auto93.csv
8  -h --help      show help            = false
9  -m --min       min size = n*(the.min) = .5
10 -n --nums      how many numbers to keep = 256
11 -p --p         distance coefficient  = 2
12 -s --seed      random number seed   = 10019
13 -S --sample    how many rows to search = 512]]]

```

Commonly used lib functions.

```

14 local lt,o,o,map  = 1.lt,1.o,1.o,1.map
15 local per,push,sort = 1.per, 1.push,1.sort

16 -----
17 local Data,Cols,Sym,Num,Row

```

Holder of rows and their summaries (in cols).

```

18 function Data() return {_is = "Data",
19                        cols= nil, -- summaries of data
20                        rows= {} -- kept data
21                        } end

```

Holds of summaries of columns. Columns are created once, then may appear in multiple slots.

```

22 function Cols() return {
23   _is = "Cols",
24   names={}, -- all column names
25   all={}, -- all the columns (including the skipped ones)
26   klass=nil, -- the single dependent klass column (if it exists)
27   x={}, -- independent columns (that are not skipped)
28   y={} -- depedent columns (that are not skipped)
29 } end

```

Summarizers a stream of symbols.

```

30 function Sym(c,s)
31   return {_is= "Sym",
32           n=0, -- items seen
33           at=c or 0, -- column position
34           names=s or "", -- column name
35           _has={} -- kept data
36   } end

```

Summarizes a stream of numbers.

```

37 function Num(c,s)
38   return {_is="Nums",
39           n=0,at=c or 0, name=s or "", _has={}, -- as per Sym
40           isNum=true, -- mark that this is a number
41           lo=-math.huge, -- lowest seen
42           hi=-math.huge, -- highest seen
43           isSorted=true, -- no updates since last sort of data
44           w=(s or ""):find"$" and -1 or 1 -- minimizing if w=-1
45   } end

```

Holds one record

```

46 function Row(t) return {_is="Row",
47                        cells=t, -- one record
48                        cooked=1,copy(t), -- used if we discretize data
49                        isEvaluated=false -- true if y-values evaluated.
50   } end

```

```

51 -----
52 local add,adds,clone,div,mid,norm,nums,record,records,stats

```

----- Create Generate rows from some src. If src is a string, read rows from file; else read rows from a src table. When reading, use row1 to define columns.

```

53 function records(src, data,head,body)
54   function head(sNames)
55     local cols = Cols()
56     cols.names = names
57     for c,s in pairs(sNames) do
58       local col = push(cols.all, -- Numerics start with Uppercase.
59                        (s:find"[A-Z]" and Num or Sym)(c,s))
60       if not s:find"$" then -- some columns are skipped
61         push(s:find"[!-]" and cols.y or cols.x, col) -- some cols are goal cols
62         if s:find"$" then cols.klass=col end end end
63     return cols
64   end -----
65   function body(t) -- treat first row differently (defines the columns)
66     if data.cols then record(data,t) else data.cols=head(t) end
67   end -----
68   data = Data()
69   if type(src)=="string" then l.csv(src, body) else
70     for _,t in pairs(src or {}) do body(t) end end
71   return data end

```

Return a new data with same structure as data1. Optionally, oad in rows.

```

72 function clone(data1, rows)
73   data2=Data()
74   data2.cols = _head(data1.cols.names)
75   for _,row in pairs(rows or {}) do record(data2,row) end
76   return data2 end

77 -----

```

Add one thing to col. For Num, keep at most nums items.

```

78 function add(col,v)
79   if v=="?" then
80     col.n = col.n + 1
81     if not col.isNum then col._has[v] = 1 + (col._has[v] or 0) else
82       col.lo = math.min(v, col.lo)
83       col.hi = math.max(v, col.hi)
84     local pos
85     if #col._has < the.nums then pos = 1 + (#col._has)
86     elseif math.random() < the.nums/col.n then pos = math.random(#col._has) end
87     if pos then col.isSorted = false
88     col._has[pos] = tonumber(v) end end end end

```

Add many things to col

```

89 function adds(col,t) for _,v in pairs(t) do add(col,v) end; return col end

```

Add a row to data. Calls add() to updatie the cols with new values.

```

90 function record(data,xs)
91   local row= push(data.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
92   for _,todo in pairs(data.cols.x, data.cols.y) do
93     for _,col in pairs(todo) do
94       add(col, row.cells[col.at]) end end end

```

----- Query

Return kept numbers, sorted.

```

96 function nums(num)
97   if not num.isSorted then num._has = sort(num._has); num.isSorted=true end
98   return num._has end

```

Normalized numbers 0.1. Everything else normalizes to itself.

```

99 function norm(col,n)
100   return x=="?" or not col.isNum and x or (n-col.lo)/(col.hi-col.lo + 1E-32) end

```

Diversity (standard deviation for Nums, entropy for Syms)

```

101 function div(col)
102   if col.isNum then local a=nums(col); return (per(a,.9)-per(a,.1))/2.58 else
103     local function fun(p) return p*math.log(p,2) end
104     local e=0
105     for _,n in pairs(col._has) do if n>0 then e=e+fun(n/col.n) end end
106     return e end end

```

Central tendency (median for Nums, mode for Syms)

```

107 function mid(col)
108   if col.isNum then return per(nums(col),.5) else
109     local most,mode = -1
110     for k,v in pairs(col._has) do if v>most then mode,most=k,v end end
111     return mode end end

```

For showCols (default=data.cols.x) in data, report fun (default=mid).

```

112 function stats(data, showCols,fun, t)
113   showCols, fun = showCols or data.cols.y, fun or mid
114   t={}; for _,col in pairs(showCols) do t[col.name]=fun(col) end; return t end

```

```

115 -----
116 local bins,cook,divs

```

Find ranges within a num (unsupervised).

```

117 function bins(num)
118   local a, epsilon = nums(num), the.cohen*div(num)
119   local enough = #a*the.min
120   local one = {lo=a[1], hi=a[1], n=0}
121   local t = {one}
122   for i,x in pairs(a) do
123     if i < #a-enough and x ~= a[i+1] and n > enough and hi-lo > epsilon then
124       one = push(t, {lo=one.hi, hi=a[i], n=0}) end
125     one.hi = a[i]
126     one.n = 1 + one.n end
127   t[1].lo = -math.huge
128   t[#t].ho = math.huge
129   return t end

```

Fill in discretized values (in cooked).

```

130 function cook(data)
131   for _,num in pairs(data.cols.x) do
132     if num.isNum then local t = bins(num)
133       for _,row in pairs(data.rows) do
134         local v = row.cells[num.at]
135         if v ~= "?" then
136           for _,bin in pairs(t) do
137             if v > bin.lo and v <= bin.hi then
138               row.cooked[col.at] = bin.lo
139             break end end end end end

```

Sum the entropy of the cooked independent columns.

```

140 function divs(data,rows)
141   local n = 0
142   for _,col in pairs(data.cols.x) do
143     local sym= Sym()
144     for _,row in pairs(rows or data.rows) do
145       v = row.cooked[col.at]
146       if v ~= "?" then add(s, v) end end
147     n = n + div(sym) end
148   return n end

149 -----
150 local around, dist, far, half, halves, tree

```

Distance between rows (returns 0.1). For unknown values, assume max distance.

```

151 function dist(data,t1,t2)
152   local function fun(col, v1,v2)
153     if v1=="?" and v2=="?" then return 1 end
154     if not col.isNum then return v1==v2 and 0 or 1 end
155     v1,v2 = norm(col,v1), norm(col,v2)
156     if v1=="?" then v1 = v2<.5 and 1 or 0 end
157     if v2=="?" then v2 = v1<.5 and 1 or 0 end
158     return math.abs(v1-v2)
159   end -----
160   local d = 0
161   for _,col in pairs(data.cols.x) do
162     d = d + fun(col, t1.cells[col.at], t2.cells[col.at])*the.p end
163   return (d/#data.cols.x)^(1/the.p) end

```

Sort rows (default=data.rows) by distance to row1.

```

164 function around(data,row1, rows, fun)
165   function fun(row2) --print("r2",#row2);
166     return {row=row2, dist=dist(data,row1,row2)} end
167   print("a")
168   return sort(map(rows or data.rows,fun),lt"dist") end

```

Return the row that is the far to max distance away from row.

```

169 function far(data,row, rows)
170   print("f",data._is,#rows,o(row))
171   return per(around(data,row,rows), the.far).row end

```

Split rows (default=data.rows) in half by distance to 2 distant points.

```

172 function half(data,rows, rowAbove)
173   local rows = rows or data.rows
174   local some = 1.many(rows, the.sample)
175   local left = rowAbove or far(data, 1.any(some),some)
176   print(4,data._is,o(left),#some)
177   local right = far(data, left,some)
178   print(5)
179   local c = dist(data,left,right)
180   local lefts,rights = {},{}
181   local function fun(row)
182     local a = dist(data,row,left)
183     local b = dist(data,row,right)
184     return {row=rows, d=(a^2 + c^2 - b^2) / (2*c)} end
185   for i,rowd in pairs(sort(map(rows, fun), lt"d")) do
186     push(i < (#rows)/2 and lefts or rights, rowd.row) end
187   return left,right,lefts,rights,c end

```

```

188 function halves(data,rows, stop,rowAbove)
189   rows = rows or data.rows
190   stop = stop or (#rows)*the.min
191   if #rows <= stop then return {node=rows} end
192   local left,right,lefts,rights,_ = half(data,rows,rowAbove)
193   return {node=rows, kids={halves(data,lefts,stop,left),
194                               halves(data,rights,stop,right)}} end

```

```

195 function tree(x, nodeFun, pre)
196   nodeFun = nodeFun or io.write
197   pre = pre or "... "
198   print(pre,nodeFun(x,node))
199   for _,kid in pairs(x.kids or {}) do tree(kid, nodeFun, pre..".. ") end end

```

That's all folks.

“lua return { the=the, Data=Data, Cols=Cols, Sym=Sym, Num=Num, Row=Row, add=add, adds=adds, around=around, bin=bin,clone=clone,cook=cook,dist=dist, div=div, divs=divs, far=far, half=half, halves=halves, mid=mid, nums=nums, records=records, record=record, stats=stats}