```lua
1   local b4={}; for k,v in pairs(_ENV) do b4[k]=v end -- LUA trivia. Ignore.
2   local help=[[
3   CSV : summarized csv file
4   (c) 2022 Tim Menzies <timm@ieee.org> BSD-2 license
5
6   USAGE: lua seen.lua [OPTIONS]
7
8   OPTIONS:
9    -e  --eg        start-up example              = nothing
10   -d  --dump      on test failure, exit with stack dump = false
11   -f  --file      file with csv data            = ../data/auto93.csv
12   -h  --help      show help                     = false
13   -n  --nums      number of nums to keep        = 512
14   -s  --seed      random number seed            = 10019
15   -S  --seperator feild seperator               = ,]]
16
17  -- Function argument conventions:
18  -- 1. two blanks denote optionas, four blanls denote locals:
19  -- 2. prefix n,s,is,fun denotes number,string,bool,function;
20  -- 3. suffix s means list of thing (so names is list of strings)
21  -- 4. c is a column index (usually)
22
23  -- ## Misc routines
24  -- ### Handle Settings
25  local the,coerce,cli
26  -- Parse 'the' config settings from 'help'.
27  function coerce(s,      fun)
28    function fun(s1)
29      if s1=="true"  then return true end
30      if s1=="false" then return false end
31      return s1 end
32    return math.tointeger(s) or tonumber(s) or fun(s:match"^%s*(.-)%s*$") end
33
34  -- Create a 'the' variables
35  the={}
36  help:gsub("\n[-][%S]+[%s]+[-][-]([%S]+)[^\n]+=([%S]+)",
37                   function(k,x) the[k]=coerce(x) end)
38
39  -- Update settings from values on command-line flags. Booleans need no values
40  -- (we just flip the defeaults).
41  function cli(t)
42    for slot,v in pairs(t) do
43      v = tostring(v)
44      for n,x in ipairs(arg) do
45        if x=="-"..(slot:sub(1,1)) or x=="--"..slot then
46          v = v=="false" and "true" or v=="true" and "false" or arg[n+1] end end
47      t[slot] = coerce(v) end
48    if t.help then os.exit(print("\n"..help.."\n")) end
49    return t end
50
51  -- ### Linting code
52  -- ### Lists
53  local copy,per,push,csv
54  -- deepcopy
55  function copy(t,      u)
56    if type(t) ~= "table" then return t end
57    u={}; for k,v in pairs(t) do u[k] = copy(v) end
58    return setmetatable(u,getmetatable(t))   end
59
60  -- Return the 'p'-th thing from the sorted list 't'.
61  function per(t,p)
62    p=math.floor(((p or .5)*#t)+.5); return t[math.max(1,math.min(#t,p))] end
63
64  -- Add to 't', return 'x'.
65  function push(t,x) t[1+#t]=x; return x end
66
67  -- ## Call 'fun' on each row. Row cells are divided in 'the.seperator'.
68  function csv(fname,fun,       sep,src,s,t)
69    sep = "([^" .. the.seperator .. "]+)"
70    src = io.input(fname)
71    while true do
72      s = io.read()
73      if not s then return io.close(src) else
74        t={}
75        for s1 in s:gmatch(sep) do t[1+#t] = coerce(s1) end
76        fun(t) end end end
77
78  -- ### Strings
79  local o,oo
80  -- 'o' is a telescopt and 'oo' are some binoculars we use to exam stucts.
81  -- 'o':  generates a string from a nested table.
82  function o(t,     show,u)
83    if type(t) ~= "table" then return tostring(t) end
84    function show(k,v)
85      if not tostring(k):find"^_"  then
86        v = o(v)
87        return #t==0 and string.format(":%s %s",k,v) or tostring(v) end end
88    u={}; for k,v in pairs(t) do u[1+#u] = show(k,v) end
89    if #t==0 then table.sort(u) end
90    return "["..table.concat(u,"")..."]" end
91
92  -- 'oo': prints the string from 'o'.
93  function oo(t) print(o(t)) return t end
94
95  -- ### Misc
96  local rogues, rnd,  obj
97  --- Find rogue locals.
98  function rogues()
99    for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end end
100
101 -- ### Maths
102 function rnd(x, places)
103   local mult = 10^(places or 2)
104   return math.floor(x * mult + 0.5) / mult end
105
106 -- obj("Thing") enables a constructor Thing:new() ... and a pretty-printer
107 -- for Things.
108 function obj(s,      t,i,new)
109   function new(k,...) i=setmetatable({},k);
110                       return setmetatable(t.new(i,...) or i,k) end
111   t={__tostring = function(x) return s..o(x) end}
112   t.__index = t;return setmetatable(t,{__call=new}) end
```

```lua
113 -- ---------------------------------------
114 -- ## Objects
115 local Cols,Data,Num,Row,Sym=obj"Cols",obj"Data",obj"Num",obj"Rows",obj"Sym"
116
117 -- 'Sym's summarize a stream of symbols.
118 function Sym:new(c,s)
119   return {n=0,          -- items seen
120           at=c or 0,    -- column position
121           name=s or "", -- column name
122           _has={}       -- kept data
123          } end
124
125 -- 'Num' ummarizes a stream of numbers.
126 function Num:new(c,s)
127   return {n=0,at=c or 0, name=s or "", _has={}, -- as per Sym
128           lo= math.huge,   -- lowest seen
129           hi= -math.huge,  -- highest seen
130           isSorted=true,   -- no updates since last sort of data
131           w = ((s or "):find"-$" and -1 or 1)
132          } end
133
134 -- 'Columns' Holds of summaries of columns.
135 -- Columns are created once, then may appear in  multiple slots.
136 function Cols:new(names)
137   self.names=names -- all column names
138   self.all={}      -- all the columns (including the skipped ones)
139   self.klass=nil   -- the single dependent klass column (if it exists)
140   self.x={}        -- independent columns (that are not skipped)
141   self.y={}        -- depedent columns (that are not skipped)
142   for c,s in pairs(names) do
143     local col = push(self.all, -- Numerics start with Uppercase.
144                   (s:find"^[A-Z]*" and Num or Sym)(c,s))
145     if not s:find"$" then -- some columns are skipped
146       push(s:find"[!+-]" and self.y or self.x, col) -- some cols are goal cols
147       if s:find"!$" then self.klass=col end end end end
148
149 -- 'Row' holds one record
150 function Row:new(t) return {cells=t,         -- one record
151                             cooked=copy(t),  -- used if we discretize data
152                             isEvaled=false   -- true if y-values evaluated.
153                            } end
154
155 -- 'Data' is a holder of 'rows' and their sumamries (in 'cols').
156 function Data:new(src)
157   self.cols = nil -- summaries of data
158   self.rows = {}  -- kept data
159   if   type(src) == "string"
160   then csv(src, function(row) self:add(row) end)
161   else for _,row in pairs(src or {}) do self:add(row) end end end
```

```lua
162 -- ---------------------------------------
163 -- ## Sym
164 -- Add one thing to 'col'. For Num, keep at most 'nums' items.
165 function Sym:add(v)
166   if v~="?" then self.n=self.n+1; self._has[v] = 1 + (self._has[v] or 0) end end
167
168 function Sym:mid(col,     most,mode)
169   most = -1; for k,v in pairs(self._has) do if v>most then mode,most=k,v end end
170   return mode end
171
172 function Sym:div(     e,fun)
173   function fun(p) return p*math.log(p,2) end
174   e=0; for _,n in pairs(self._has) do if n>0 then e=e - fun(n/self.n) end end
175   return e end
176
177 -- ---------------------------------------
178 -- ## Num
179 -- Return kept numbers, sorted.
180 function Num:nums()
181   if not self.isSorted then table.sort(self._has); self.isSorted=true end
182   return self._has end
183
184 -- Reservoir sampler. Keep at most 'the.nums' numbers
185 -- (and if we run out of room, delete something old, at random).,
186 function Num:add(v,      pos)
187   if v~="?" then
188     self.n = self.n + 1
189     self.lo = math.min(v, self.lo)
190     self.hi = math.max(v, self.hi)
191     if      #self._has < the.nums           then pos = 1 + (#self._has)
192     elseif math.random() < the.nums/self.n then pos = math.random(#self._has) end
193     if pos then self.isSorted = false
194                 self._has[pos] = tonumber(v) end end end
195
196 -- Diversity (standard deviation for Nums, entropy for Syms)
197 function Num:div(     a)  a=self:nums(); return (per(a,.9)-per(a,.1))/2.58 end
198
199 -- Central tendancy (median for Nums, mode for Syms)
200 function Num:mid()  return per(self:nums(),.5) end
```

```
200
201  -- ---------------------------------------
202  -- ## Data
203  -- Add a 'row' to 'data'. Calls 'add()' to  updatie the 'cols' with new values.
204  function Data:add(xs,    row)
205    if   not self.cols
206    then self.cols = Cols(xs)
207    else row= push(self.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
208         for _,todo in pairs{self.cols.x, self.cols.y} do
209           for _,col in pairs(todo) do
210             col:add(row.cells[col.at]) end end end end
211
212  -- For 'showCols' (default='data.cols.x') in 'data', report 'fun' (default='mid'),
213  -- rounding numbers to 'places' (default=2)
214  function Data:stats(  places,showCols,fun,    t,v)
215    showCols, fun = showCols or self.cols.y, fun or "mid"
216    t={}; for _,col in pairs(showCols) do
217           v=fun(col)
218           v=type(v)=="number" and rnd(v,places) or v
219           t[col.name]=v; return t end
```

```
220
221  -- ---------------------------------
222  -- ## Test Engine
223  local eg, fails = {},0
224
225  -- 1. reset random number seed before running something.
226  -- 2. Cache the detaults settings, and...
227  -- 3. ... restore them after the test
228  -- 4. Print error messages or stack dumps as required.
229  -- 5. Return true if this all went well.
230  local function runs(k,    old,status,out,msg)
231    if not eg[k] then return end
232    math.randomseed(the.seed) -- reset seed [1]
233    old={}; for k,v in pairs(the) do old[k]=v end --  [2]
234    if the.dump then -- [4]
235      status,out = true, eg[k]()
236    else
237      status,out = pcall(eg[k])  -- pcall means we do not crash and dump on errror
238    end
239    for k,v in pairs(old) do the[k]=v end -- restore old settings [3]
240    msg = status and ((out==true and "PASS") or "FAIL") or "CRASH" -- [4]
241    print("!!!!!!!", msg, k, status)
242    return out or err end
243
244  -- ---------------------------------
245  -- ## Tests
246  -- Test that the test  happes when something crashes?
247  function eg.BAD() print(eg.dont.have.this.field) end
248
249  -- Sort all test names.
250  function eg.LIST(   t)
251    t={}; for k,_ in pairs(eg) do t[1+#t]=k end; table.sort(t); return t end
252
253    -- List test names.
254  function eg.LS()
255    print("\nExamples lua csv -e ...")
256    for _,k in pairs(eg.LIST()) do print(string.format("\t%s",k)) end
257    return true end
258
259  -- Run all tests
260  function eg.ALL()
261    for _,k in pairs(eg.LIST()) do
262      if k ~= "ALL" then
263        print("\n---------------------------------"
264        if not runs(k) then fails=fails+ 1 end end end
265    return true end
266
267  -- Settings come from big string top of "sam.lua"
268  -- (maybe updated from comamnd line)
269  function eg.the() oo(the); return true end
270
271  -- The middle and diversity of a set of symbols is called "mode"
272  -- and "entropy" (and the latter is zero when all the symbols
273  -- are the same).
274  function eg.sym(   sym,entropy,mode)
275    sym= Sym()
276    for _,x in pairs{"a","a","a","a","b","b","c"} do sym:add(x) end
277    mode, entropy = sym:mid(), sym:div()
278    entropy = (1000*entropy)//1/1000
279    oo({mid=mode, div=entropy})
280    return mode=="a" and 1.37 <= entropy and entropy <=1.38 end
281
282  -- The middle and diversity of a set of numbers is called "median"
283  -- and "standard deviation" (and the latter is zero when all the nums
284  -- are the same).
285  function eg.num(   num,mid,div)
286    num=Num()
287    for i=1,100 do num:add(i) end
288    mid,div = num:mid(), num:div()
289    print(mid ,div)
290    return 50<= mid and mid<= 52 and 30.5 <div and div<32 end
291
292  -- Nums store only a sample of the numbers added to it (and that storage
293  -- is done such that the kept numbers span the range of inputs).
294  function eg.bignum(  num)
295    num=Num()
296    the.nums = 32
297    for i=1,1000 do num:add(i) end
298    oo(num:nums())
299    return 32==#num._has; end
300
301  -- Show we can read csv files.
302  function eg.csv(   n)
303    n=0
304    csv("../data/auto93.csv",function(row)
305      n=n+1; if n> 10 then return else oo(row) end end); return true end
306
307  -- Can I load a csv file into a Data?.
308  function eg.data(   d)
309    d = Data("../data/auto93.csv")
310    for _,col in pairs(d.cols.y) do oo(col) end
311    return true
312  end
313
314  -- Print some stats on columns.
315  function eg.stats(   data,mid,div)
316    data = Data("../data/auto93.csv")
317    div=function(col) return col:div() end
318    mid=function(col) return col:mid() end
319    print("xmid", o( data:stats(2,data.cols.x, mid)))
320    print("xdiv", o( data:stats(3,data.cols.x, div)))
321    print("ymid", o( data:stats(2,data.cols.y, mid)))
322    print("ydiv", o( data:stats(3,data.cols.y, div)))
323    return true
324  end
325
326  -- ---------------------------------
327  the = cli(the)
328  runs(the.eg)
329  rogues()
330  os.exit(fails)
```