```
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   -k  --k         bayes low frequency factor = 2
14   -m  --m         bayes low frequency factor = 1
15   -n  --nums      number of nums to keep     = 512
16   -p  --p         distance calculation coefficient = 2
17   -s  --seed      random number seed        = 10019
18   -S  --seperator feild seperator           = , ]]
```

Function argument conventions:

1. two blanks denote optionas, four blanls denote locals:
2. prefix n,s,is,fun denotes number,string,bool,function;
3. suffix s means list of thing (so names is list of strings)
4. c is a column index (usually)

## Misc routines

### Handle Settings

```
19   local the,coerce,cli
```

Parse the config settings from help.

```
20   function coerce(s,        fun)
21     function fun(s1)
22       if s1=="true"  then return true end
23       if s1=="false" then return false end
24       return s1 end
25     return math.tointeger(s) or tonumber(s) or fun(s:match"^%s*(.-)%s*$") end
```

Create a the variables

```
26   the={}
27   help:gsub("\n [-]%S]+[%s]+[-]-([%S]+)[^\n]+= ([%S]+)",
28                function(k,x) the[k]=coerce(x) end)
```

Update settings from values on command-line flags. Booleans need no values (we just flip the defeaults).

```
29   function cli(t)
30     for slot,v in pairs(t) do
31       v = tostring(v)
32       for n,x in ipairs(arg) do
33         if x=="-".. (slot:sub(1,1)) or x=="--"..slot then
34           v = v=="false" and "true" or v=="true" and "false" or arg[n+1] end end
35       t[slot] = coerce(v) end
36     if t.help then os.exit(print("\n"..help.."\n")) end
37     return t end
```

### Lists

```
38   local copy,per,push,csv,sort,map,lt
```

Deepcopy

```
39   function copy(t,       u)
40     if type(t) ~= "table" then return t end
41     u={}; for k,v in pairs(t) do u[k] = copy(v) end
42     return setmetatable(u,getmetatable(t))   end
```

Return the p-th thing from the sorted list t.

```
43   function per(t,p)
44     p=math.floor(((p or .5)*#t)+.5); return t[math.max(1,math.min(#t,p))] end
```

Add to t, return x.

```
45   function push(t,x) t[1+#t]=x; return x end
```

Function, return a sorted list.

```
46   function sort(t,f) table.sort(t,f); return t end
```

Sorting function

```
47   function lt(x) return function(t1,t2) return t1[x] < t2[x] end end
```

Map a function over a list

```
48   function map(t1,fun,      t2)
49     t2={}; for _,v in pairs(t1) do t2[1+#t2] = fun(v) end; return t2 end
```

Call fun on each row. Row cells are divided in the seperator.

```
50   function csv(fname,fun,      sep,src,s,t)
51     sep = "([^".. the.seperator .. "]+)"
52     src = io.input(fname)
53     while true do
54       s = io.read()
55       if not s then return io.close(src) else
56         t={}
57         for s1 in s:gmatch(sep) do t[1+#t] = coerce(s1) end
58         fun(t) end end end
```

### Strings

```
59   local o,oo
```

o is a telescopt and oo are some binoculars we use to exam stucts. o: generates a string from a nested table.

```
60   function o(t,        show,u)
61     if type(t) ~=  "table" then return tostring(t) end
62     function show(k,v)
63       if not tostring(k):find"^_"  then
64         v = o(v)
65         return #t==0 and string.format(":%s %s",k,v) or tostring(v) end end
66     u={}; for k,v in pairs(t) do u[1+#u] = show(k,v) end
67     if #t==0 then table.sort(u) end
68     return "{"..table.concat(u," ").."}" end
```

oo: prints the string from o.

```
69   function oo(t) print(o(t)) return t end
```

### Misc

```
70   local rogues, rnd,  obj
```

Find rogue locals.

```
71   function rogues()
72     for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end end
```

### Maths

```
73   function rnd(x, places)
74     local mult = 10^(places or 2)
75     return math.floor(x * mult + 0.5) / mult end
```

obj("Thing") enables a constructor Thing:new() ... and a pretty-printer for Things.

```
76   function obj(s,     t,i,new)
77     function new(k,...) i=setmetatable({},k);
78                         return setmetatable(t.new(i,...) or i,k) end
79     t={__tostring = function(x) return s..o(x) end}
80     t.__index = t;return setmetatable(t,{__call=new}) end
```

## Objects

```
82   local Cols,Data,Num,Row,Sym=obj"Cols",obj"Data",obj"Num",obj"Rows",obj"Sym"
```

Syms summarize a stream of symbols.

```
83   function Sym:new(c,s)
84     return {n=0,               -- items seen
85             at=c or 0,         -- column position
86             name=s or "",      -- column name
87             _has={}            -- kept data
88            } end
```

Num ummarizes a stream of numbers.

```
89   function Num:new(c,s)
90     return {n=0,at=c or 0, name=s or "", _has={}, -- as per Sym
91             lo= math.huge,    -- lowest seen
92             hi= -math.huge,   -- highest seen
93             isSorted=true,    -- no updates since last sort of data
94             w = ((s or ""):find"-$" and -1 or 1)
95            } end
```

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

```
96   function Cols:new(names)
97     self.names=names  -- all column names
98     self.all={}       -- all the columns (including the skipped ones)
99     self.klass=nil    -- the single dependent klass column (if it exists)
100    self.x={}         -- independent columns (that are not skipped)
101    self.y={}         -- depedent columns (that are not skipped)
102    for c,s in pairs(names) do
103      local col = push(self.all, -- Numerics start with Uppercase.
104                       (s:find"^[A-Z]*" and Num or Sym)(c,s))
105      if not s:find":$" then -- some columns are skipped
106        push(s:find"[!+-]" and self.y or self.x, col) -- some cols are goal cols
107        if s:find"!$" then self.klass=col end end end end
```

Row holds one record

```
108   function Row:new(t) return (cells=t,           -- one record
109                        cooked=copy(t),  -- used if we discretize data
110                        isEvaled=false   -- true if y-values evaluated.
111                        } end
```

Data is a holder of rows and their sumamries (in cols).

```
112   function Data:new(src)
113     self.cols = nil -- summaries of data
114     self.rows = {}  -- kept data
115     if   type(src) == "string"
116     then csv(src, function(row) self:add(row) end)
117     else for _,row in pairs(src or {}) do self:add(row) end end end
```

### Sym

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

```
118   function Sym:add(v)
119     if v~="?" then self.n=self.n+1; self._has[v]= 1+(self._has[v] or 0) end end
```

```
120   function Sym:mid(col,     most,mode)
121     most=-1; for k,v in pairs(self._has) do if v>most then mode,most=k,v end end
122     return mode end
```

distance between two values.

```
123   function Sym:dist(v1,v2)
124     return  v1=="?" and v2=="?" and 1 or v1==v2 and 0 or 1 end
```

Diversity measure for symbols = entropy.

```
125   function Sym:div(     e,fun)
126     function fun(p) return p*math.log(p,2) end
127     e=0; for _,n in pairs(self._has) do if n>0 then e=e - fun(n/self.n) end end
128     return e end
```

```
129   -- Return how much `x` might belong to `self`.
130   function SYM:like(x,prior)
131     return ((self._has[x] or 0)+the.m*prior) / (self.n+the.m) end
```

### Num

Return kept numbers, sorted.

```
132   function Num:nums()
133     if not self.isSorted then table.sort(self._has); self.isSorted=true end
134     return self._has end
```

Reservoir sampler. Keep at most the nums numbers (and if we run out of room, delete something old, at random).,

```
135   function Num:add(v,     pos)
136     if v~="?" then
137       self.n  = self.n + 1
138       self.lo = math.min(v, self.lo)
139       self.hi = math.max(v, self.hi)
140       if      #self._has < the.nums        then pos=1 + (#self._has)
141       elseif math.random() < the.nums/self.n then pos=math.random(#self._has) end
142       if pos then self.isSorted = false
143                   self._has[pos] = tonumber(v) end end end
```

distance between two values.

```
144   function Num:dist(v1,v2)
145     if     v1=="?" and v2=="?" then return 1 end
146     v1,v2 = self:norm(v1),  self:norm(v2)
147     if v1=="?" then v1 = v2<.5 and 1 or 0 end
148     if v2=="?" then v2 = v1<.5 and 1 or 0 end
149     return math.abs(v1-v2) end
```

Return middle

```
150   function Num:mid()  return per(self:nums(), .5) end
```

Return diversity

```
151    function Num:div() return (per(self:nums(), .9) - per(self:nums(),.1))/2.58 end
```

Normalized numbers 0..1. Everything else normalizes to itself.

```
152    function Num:norm(n)
153      return x=="?" and x or (n-self.lo)/(self.hi-self.lo + 1E-32) end
```

Return the likelihood that x belongs to i. <

```
154    function NUM:like(x,...)
155      local sd,mu=self:div(), self:mid()
156      if sd<=0 then return x==mu and 1 or 1/big end
157      return math.exp(-.5*((x - mu)/sd)^2) / (sd*((2*math.pi)^0.5)) end
```

## Data

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

```
158    function Data:add(xs,     row)
159      if   not self.cols
160      then self.cols = Cols(xs)
161      else row= push(self.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
162           for _,todo in pairs(self.cols.x, self.cols.y) do
163             for _,col in pairs(todo) do
164               col:add(row.cells[col.at]) end end end end
```

Return a new Data that mimics structure of self. Add src to the clone.

```
165    function Data:clone(  src,     out)
166      out = Data()
167      out:add(self.cols.name)
168      for _,row in pairs(src or {}) do out:add(row) end
169      return out end
```

For showCols (default=data.cols.x) in data, show fun (default=mid), rounding numbers to places (default=2)

```
170    function Data:stats(  places,showCols,fun,     t,v)
171      showCols, fun = showCols or self.cols.y, fun or "mid"
172      t={}; for _,col in pairs(showCols) do
173        v=fun(col)
174        v=type(v)=="number" and rnd(v,places) or v
175        t[col.name]=v end; return t end
```

Distance between rows (returns 0..1). For unknown values, assume max distance.

```
177    function Data:dist(row1,row2)
178      local d = 0
179      for _,col in pairs(self.cols.x) do
180        d = d + col:dist(row1.cells[col.at], row2.cells[col.at])^the.p end
181      return (d/#self.cols.x)^(1/the.p) end
```

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

```
182    function Data:around(row1,  rows,    fun)
183      function fun(row2) return {row=row2, dist=self:dist(row1,row2)} end
184      return sort(map(rows or self.rows, fun),lt"dist") end
```

Return P(H)\*P(E1|H))\*p(E2|H)... . Work in logs (to cope with small nums)

```
185    function Data:like(row, nklasses, nrows)
186      local prior,like,inc,x
187      prior = (#self.rows + the.k) / (nrows + the.k * nklasses)
188      like  = math.log(prior)
189      row = row.cells and row.cells or row
190      for _,col in pairs(self.cols.x) do
191        x = row[col.at]
192        if x ~= nil and x ~= "?" then
193          inc  = col:like(x,prior)
194          like = like + math.log(inc) end end
195      return like end
```

## Test Engine

```
196    local eg, fails = {},0
```

1.  reset random number seed before running something.
2.  Cache the defaults settings, and…
3.  … restore them after the test
4.  Print error messages or stack dumps as required.
5.  Return true if this all went well.

```
197    local function runs(k,      old,status,out,msg)
198      if not eg[k] then return end
199      math.randomseed(the.seed) -- reset seed [1]
200      old={}; for k,v in pairs(the) do old[k]=v end --  [2]
201      if the.dump then -- [4]
```

```
202        status,out=true, eg[k]()
203      else
204        status,out=pcall(eg[k]) -- pcall means we do not crash and dump on errror
205      end
206      for k,v in pairs(old) do the[k]=v end -- restore old settings [3]
207      msg = status and ((out==true and "PASS") or "FAIL") or "CRASH" -- [4]
208      print("!!!!!!!", msg, k, status)
209      return out or err end
```

## Tests

Test that the test happes when something crashes?

```
210    function eg.BAD() print(eg.dont.have.this.field) end
```

Sort all test names.

```
211    function eg.LIST(    t)
212      t={}; for k,_ in pairs(eg) do t[1+#t]=k end; table.sort(t); return t end
```

```
213    -- List test names.
214    function eg.LS()
215      print("\nExamples lua csv -e ...")
216      for _,k in pairs(eg.LIST()) do print(string.format("\t%s",k)) end
217      return true end
```

Run all tests

```
218    function eg.ALL()
219      for _,k in pairs(eg.LIST()) do
220        if k ~= "ALL" then
221          print"\n--------------------------------"
222          if not runs(k) then fails=fails+ 1 end end end
223      return true end
```

Settings come from big string top of "sam.lua" (maybe updated from comamnd line)

```
225    function eg.the() oo(the); return true end
```

The middle and diversity of a set of symbols is called "mode" and "entropy" (and the latter is zero when all the symbols are the same).

```
226    function eg.sym(  sym,entropy,mode)
227      sym= Sym()
228      for _,x in pairs{"a","a","a","a","b","b","c"} do sym:add(x) end
229      mode, entropy = sym:mid(), sym:div()
230      entropy = (1000*entropy)//1/1000
231      oo{{mid=mode, div=entropy})
232      return mode=="a" and 1.37 <= entropy and entropy <=1.38 end
```

The middle and diversity of a set of numbers is called "median" and "standard deviation" (and the latter is zero when all the nums are the same).

```
233    function eg.num(  num,mid,div)
234      num=Num()
235      for i=1,100 do num:add(i) end
236      mid,div = num:mid(), num:div()
237      print(mid ,div)
238      return 50<= mid and mid<= 52 and 30.5 <div and div<32 end
```

Nums store only a sample of the numbers added to it (and that storage is done such that the kept numbers span the range of inputs).

```
239    function eg.bignum(  num)
240      num=Num()
241      the.nums = 32
242      for i=1,1000 do num:add(i) end
243      oo(num:nums())
244      return 32==#num._has; end
```

Show we can read csv files.

```
245    function eg.csv(    n)
246      n=0
247      csv("../data/auto93.csv",function(row)
248        n=n+1; if n> 10 then return else oo(row) end end); return true end
```

Can I load a csv file into a Data?.

```
249    function eg.data(    d)
250      d = Data("../data/auto93.csv")
251      for _,col in pairs(d.cols.y) do oo(col) end
252      return true
253    end
```

Print some stats on columns.

```
254    function eg.stats(    data,mid,div)
255      data = Data("../data/auto93.csv")
256      div  = function(col) return col:div() end
257      mid  = function(col) return col:mid() end
258      print("xmid", o( data:stats(2, data.cols.x, mid)))
259      print("xdiv", o( data:stats(3, data.cols.x, div)))
260      print("ymid", o( data:stats(2, data.cols.y, mid)))
261      print("ydiv", o( data:stats(3, data.cols.y, div)))
262      return true
263    end
```

distance functions

```
264    function eg.around(    data,around)
265      data = Data("../data/auto93.csv")
266      around = data:around(data.rows[1] )
267      for i=1,380,40 do print(around[i].dist, o(around[i].row.cells)) end
268      return true end
```

Start up

```
269    the = cli(the)
270    runs(the.eg)
271    rogues()
272    os.exit(fails)
```

That's all folks.