

Aug 27, 22 12:28

csv.lua

Page 1/3

```

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 <tim@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 --separator feild separator = ,
16
17 -- Function argument conventions:
18 -- 1. two blanks denote optionas, four blans 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 coerced(s:match("^%s*(-)%s*$") end
33
34 -- Create a 'the' variables
35 the={}
36 help:gsub("\n[-+%$]+[-+][-+](%$)+[^\n]+=(%$)+",
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 defaults).
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.."")) end
49   return t end
50
51 -- ## Linting code
52 local roques
53 -- Find roque locals.
54 function roques()
55   for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end end
56
57 -- ## lists
58 local copy,per,push,csv
59 -- deepcopy
60 function copy(t, u)
61   if type(t) ~= "table" then return t end
62   u={}; for k,v in pairs(t) do u[k] = copy(v) end
63   return setmetatable(u, getmetatable(t)) end
64
65 -- Return the 'p'-th thing from the sorted list 't'.
66 function per(t,p)
67   p=math.floor((p or .5)*#t+.5); return t[math.max(1,math.min(#t,p))] end
68
69 -- Add to 't', return 'x'.
70 function push(t,x) t[#t+1]=x; return x end
71
72 -- ## Call 'fun' on each row. Row cells are divided in 'the.separator'.
73 function csv(fname,fun, sep,src,s,t)
74   sep = "(\\" .. the.separator .. "\")+)"
75   src = io.input(fname)
76   while true do
77     s = io.read()
78     if not s then return io.close(src) else
79       t={}
80       for sl in s:gmatch(sep) do t[#t+1] = coerce(sl) end
81       fun(t) end end end
82
83 -- ## Strings
84 local o,oo
85 -- 'o' is a telescope and 'oo' are some binoculars we use to exam stuots.
86 -- 'o': generates a string from a nested table.
87 function o(t, show,u)
88   if type(t) ~= "table" then return tostring(t) end
89   function show(k,v)
90     if not tostring(k):find"^_" then
91       v = o(v)
92       return #t==0 and string.format("%s%s",k,v) or tostring(v) end end
93   u={}; for k,v in pairs(t) do u[#u+1] = show(k,v) end
94   if #t==0 then table.sort(u) end
95   return "{"..table.concat(u, ", ")}" end
96
97 -- 'oo': prints the string from 'o'.
98 function oo(t) print(o(t)) return t end
99
100 -- ## OO
101 local ako,obj
102 -- obj("Thing") enables a constructor Thing:new() ... and a pretty-printer
103 -- for Things.
104 ako=setmetatable
105 function obj(name, t)
106   t={}
107   t.__index,t.__tostring = t, function(x) return name .. o(x) end
108   return ako(t, {__call=function(x,...)
109     x=ako({},k); return ako(x.new(t,...) or t,x) end}) end

```

Aug 27, 22 12:28

csv.lua

Page 2/3

```

110 -- ## Objects
111 -- Coils,Data,Num,Rows,Sym=obj"Cols", obj"Data", obj"Num", obj"Rows", obj"Sym"
112 -- 'Data' is a holder of 'rows' and their summaries (in 'cols').
113 function Data:new() return { cols=nil, -- summaries of data
114   rows={} -- kept data
115 } end
116
117 -- 'Columns' Holds of summaries of columns.
118 -- Columns are created once, then may appear in multiple slots.
119 function Coils:new() return {
120   names={}, -- all column names
121   all={}, -- all the columns (including the skipped ones)
122   klass=nil, -- the single dependent klass column (if it exists)
123   x={}, -- independent columns (that are not skipped)
124   y={} -- dependent columns (that are not skipped)
125 } end
126
127 -- 'Sym's summarize a stream of symbols.
128 function Sym:new(c,s)
129   return {n=0, -- items seen
130     at=c or 0, -- column position
131     name=s or "", -- column name
132     _has={} -- kept data
133   } end
134
135 -- 'Num' ummarizes a stream of numbers.
136 function Num(c,s)
137   return {n=0,at=c or 0, name=s or "", _has={}, -- as per Sym
138     lo=math.huge, -- lowest seen
139     hi=-math.huge, -- highest seen
140     isSorted=true, -- no updates since last sort of data
141     w = ((s or ""):find"$" and -1 or 1)
142   } end
143
144 -- 'Row' holds one record
145 function Row(t) return {cells=t, -- one record
146   cooked=copy(t), -- used if we discretize data
147   isEvaluated=false -- true if y-values evaluated.
148 } end
149
150 -- ## Columns
151 -- Add one thing to 'col'. For Num, keep at most 'nums' items.
152 function Sym:add(v)
153   if v=="?" then
154     self.n=self.n+1; self._has[v] = 1 + (self._has[v] or 0) end end
155
156 -- Reservoir sampler. Keep at most 'the.nums' numbers
157 -- (and if we run out of room, delete something old, at random)..
158 function Num:add(col,v, pos)
159   if v=="?" then
160     self.n=self.n+1
161     self.lo = math.min(v, self.lo)
162     self.hi = math.max(v, self.hi)
163     if #self._has < the.nums then pos = 1 + (#self._has)
164     elseif math.random() < the.nums/col.n then pos = math.random(#self._has) end
165     if pos then self.isSorted = false
166     self._has[pos] = tonumber(v) end end end
167
168 -- Add a 'row' to 'data'. Calls 'add()' to updatie the 'cols' with new values.
169 function Data:add(xs)
170   local row= push(data.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
171   for _,todo in pairs(data.cols.x, data.cols.y) do
172     for _,col in pairs(todo) do
173       col:add(row.cells[col.at]) end end end
174
175 function Coils:new(names, col)
176   self.names = names
177   for c,s in pairs(names) do
178     local col = push(self.all, -- Numerics start with Uppercase.
179       (s:find"^[A-Z]" and Num or Sym)(c,s))
180     if not s:find"$" then -- some columns are skipped
181       push(s:find"[+-]" and self.y or self.x, col) -- some cols are goal cols
182       if s:find"$" then self.klass=col end end end
183
184 -- else read rows from a 'src' table. When reading, use row1 to define columns.
185 -- Generate rows from some 'src'. If 'src' is a string, read rows from file;
186 function records(src, data,head,body)
187
188   function body(t) -- treat first row differently (defines the columns)
189     if data.cols then record(data,t) else data.cols=head(t) end
190   end
191
192   data = Data()
193   if type(src)=="string" then csv(src, body) else
194     for _,t in pairs(src or {}) do body(t) end end
195   return data end
196
197 -- ## Query
198 -- Return kept numbers, sorted.
199 local function nums(num)
200   if not num.isSorted then table.sort(num._has); num.isSorted=true end
201   return num._has end
202
203 -- Diversity (standard deviation for Nums, entropy for Syms)
204 local function div(col)
205   if col.isNum then local a=nums(col); return (per(a,.9)-per(a,.1))/2.58 else
206     local function fun(p) return p*math.log(p,2) end
207     local e=0
208     for _,n in pairs(col._has) do if n>0 then e=e-fun(n/col.n) end end
209     return e end end
210
211 -- Central tendency (median for Nums, mode for Syms)
212 local function mid(col)
213   if col.isNum then return per(nums(col),.5) else
214     local most,mode = -1
215     for k,v in pairs(col._has) do if v>most then mode,most=k,v end end
216     return mode end end
217
218 -- Diversity (standard deviation for Nums, entropy for Syms)
219 local function div2(col)
220   if col.isNum then local a=nums(col); return (per(a,.9)-per(a,.1))/2.58 else
221     local function fun(p) return p*math.log(p,2) end
222     local e=0
223     for _,n in pairs(col._has) do if n>0 then e=e-fun(n/col.n) end end
224     return e end end
225
226 -- For 'showCols' (default='data.cols.x') in 'data', report 'fun' (default='mid').
227 showCols, fun = showCols or data.cols.x, fun or mid
228 showCols, fun = showCols or data.cols.y, fun or mid
229 t={}; for _,col in pairs(showCols) do t[col.name]=fun(col) end; return t end

```

```

230
231 -- -----
232 -- ## Test Engine
233 local eg, fails = {},0
234
235 -- [1] reset random number seed before running something.
236 -- [2] Cache the defaults settings, and [3] restore them after the test
237 -- [4] Print error messages or stack dumps as required.
238 -- Return true if this all went well.
239 local function runs(k, old, status, out, msg)
240     if not eg[k] then return end
241     math.randomseed(the.seed) -- reset seed [1]
242     old={}; for k,v in pairs(the) do old[k]=v end -- [2]
243     if the.dump then
244         status,out = true, eg[k]()
245     else
246         status,out = pcall(eg[k]) -- pcall means we do not crash and dump on error
247     end
248     for k,v in pairs(old) do the[k]=v end -- restore old settings [3]
249     msg = status and (out==true and "PASS") or "FAIL" or "CRASH" -- [4]
250     print("!!!!!!", msg, k, status)
251     return out or err end
252
253 -- -----
254 -- ## Tests
255 -- Test that the test happens when something crashes?
256 function eg.BAD() print(eg.dont.have.this.field) end
257
258 -- Sort all test names.
259 function eg.LIST( t)
260     t={}; for k,_ in pairs(eg) do t[1+#t]=k end; table.sort(t); return t end
261
262 -- List test names.
263 function eg.LS()
264     print("\nExamples lua csv -e...")
265     for _,k in pairs(eg.LIST()) do print(string.format("%s",k)) end
266     return true end
267
268 -- Run all tests
269 function eg.ALL()
270     for _,k in pairs(eg.LIST()) do
271         if k ~= "ALL" then
272             print("\n-----")
273             if not runs(k) then fails=fails+ 1 end end end
274     return true end
275
276 -- Settings come from big string top of "sam.lua"
277 -- (maybe updated from comand line)
278 function eg.the() oo(the); return true end
279
280 -- The middle and diversity of a set of symbols is called "mode"
281 -- and "entropy" (and the latter is zero when all the symbols
282 -- are the same).
283 function eg.sym( sym,entropy,mode)
284     sym= adds(Sym(), {"a","a","a","a","b","b","c"})
285     mode, entropy = mid(sym), div(sym)
286     entropy = (1000*entropy)//1/1000
287     oo({mid=mode, div=entropy})
288     return mode=="a" and 1.37 <= entropy and entropy <=1.38 end
289
290 -- The middle and diversity of a set of numbers is called "median"
291 -- and "standard deviation" (and the latter is zero when all the nums
292 -- are the same).
293 function eg.num( num)
294     num=Num()
295     for i=1,100 do add(num,i) end
296     local med,ent = mid(num), div(num)
297     print(mid(num) ,div(num))
298     return 50<= med and med<= 52 and 30.5 <ent and ent <32 end
299
300 -- Nums store only a sample of the numbers added to it (and that storage
301 -- is done such that the kept numbers span the range of inputs).
302 function eg.bignum( num)
303     num=Num()
304     the.nums = 32
305     for i=1,1000 do add(num,i) end
306     oo(nums(num))
307     return 32==#num._has; end
308
309 -- Show we can read csv files.
310 function eg.csv()
311     local n=0
312     csv("../data/aut093.csv",function(row)
313         n=n+1; if n> 10 then return else oo(row) end end); return true end
314
315 -- Print some stats on columns.
316 function eg.stats()
317     oo(stats(records("../data/aut093.csv"))); return true end
318
319 -- -----
320 the = cli(the)
321 runs(the.eg)
322 rogues()
323 os.exit(fails)

```