

Aug 27, 22 17:06

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 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 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 rogue 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,cmp
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[1+#t]=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 s1 in s:gmatch(sep) do t[1+#t] = coerce(s1) 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[1+#u] = 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(k,...)
109         x=ako({},k); return ako(x.new(t,...) or t,x) end}) end

```

Aug 27, 22 17:06

csv.lua

Page 2/3

```

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

```

Aug 27, 22 17:06

csv.lua

Page 3/3

```

215 -- -----
216 -- ## Test Engine
217 local eg, fails = {}, 0
218
219
220 -- [1] reset random number seed before running something.
221 -- [2] Cache the defaults settings, and [3] restore them after the test
222 -- [4] Print error messages or stack dumps as required.
223 -- Return true if this all went well.
224 local function runs(k, old, status, out, msg)
225     if not eg[k] then return end
226     math.randomseed(the.seed) -- reset seed [1]
227     old={}; for k,v in pairs(the) do old[k]=v end -- [2]
228     if the.dump then
229         status,out = true, eg[k]()
230     else
231         status,out = pcall(eg[k]) -- pcall means we do not crash and dump on error
232     end
233     for k,v in pairs(old) do the[k]=v end -- restore old settings [3]
234     msg = status and ((out==true and "PASS") or "FAIL") or "CRASH" -- [4]
235     print("!!!!!!", msg, k, status)
236     return out or err end
237
238 -- -----
239 -- ## Tests
240 -- Test that the test happens when something crashes?
241 function eg.BAD() print(eg.dont.have.this.field) end
242
243 -- Sort all test names.
244 function eg.LIST( t)
245     t={}; for k,_ in pairs(eg) do t[1+#t]=k end; table.sort(t); return t end
246
247 -- List test names.
248 function eg.LS()
249     print("\nExamples lua csv -e...")
250     for _,k in pairs(eg.LIST()) do print(string.format("%s",k)) end
251     return true end
252
253 -- Run all tests
254 function eg.ALL()
255     for _,k in pairs(eg.LIST()) do
256         if k ~= "ALL" then
257             print("\n-----")
258             if not runs(k) then fails=fails+ 1 end end end
259     return true end
260
261 -- Settings come from big string top of "sam.lua"
262 -- (maybe updated from comand line)
263 function eg.the() oo(the); return true end
264
265 -- The middle and diversity of a set of symbols is called "mode"
266 -- and "entropy" (and the latter is zero when all the symbols
267 -- are the same).
268 function eg.sym( sym,entropy,mode)
269     sym= Sym()
270     for _,x in pairs{"a","a","a","a","b","b","c"} do sym:add(x) end
271     mode, entropy = sym:mid(), sym:div()
272     entropy = (1000*entropy)//1/1000
273     oo{mid=mode, div=entropy}
274     return mode=="a" and 1.37 <= entropy and entropy <=1.38 end
275
276 -- The middle and diversity of a set of numbers is called "median"
277 -- and "standard deviation" (and the latter is zero when all the nums
278 -- are the same).
279 function eg.num( num)
280     num=Num()
281     for i=1,100 do num:add(i) end
282     local med,ent = mid(num), div(num)
283     print(mid(num),div(num))
284     return 50<= med and med<= 52 and 30.5 <ent and ent <32 end
285
286 -- Nums store only a sample of the numbers added to it (and that storage
287 -- is done such that the kept numbers span the range of inputs).
288 function eg.bignum( num)
289     num=Num()
290     the.nums = 32
291     for i=1,1000 do num:add(i) end
292     oo(nums(num))
293     return 32==#num._has; end
294
295 -- Show we can read csv files.
296 function eg.csv()
297     local n=0
298     csv("../data/auto93.csv",function(row)
299         n=n+1; if n> 10 then return else oo(row) end end); return true end
300
301 -- Print some stats on columns.
302 function eg.stats()
303     oo(stats(Data("../data/auto93.csv"))); return true end
304
305 -- -----
306 the = cli(the)
307 runs(the.eg)
308 rogues()
309 os.exit(fails)

```