

Aug 27, 22 19:06

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 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 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 -- ## 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[#t+1]=x; return x end
66
67 -- ## Call 'fun' on each row. Row cells are divided in 'the.separator'.
68 function csv(fname,fun, sep,src,s,t)
69   sep = "([" .. the.separator .. "])+"
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 sl in s:gmatch(sep) do t[#t+1] = coerce(sl) end
76       fun(t) end end end
77
78 -- ## Strings
79 local o,oo
80 -- 'o' is a telescope and 'oo' are some binoculars we use to exam stucls.
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[#u+1] = 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

```

Aug 27, 22 19:06

csv.lua

Page 2/3

```

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

```

```

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

```