

Aug 27, 22 19:08

csv.lua

Page 1/4

```

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/autog3.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 options; four blanks denote locals:
19  -- 2. prefix n,s,i,s,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(s)
29    if s=="true" then return true end
30    if s=="false" then return false end
31    return s
32  end
33  return math.tointeger(s) or tonumber(s) or fun(s:match"%s*(-)%s*$") end
34
35  -- Create a 'the' variables
36  the={}
37  help:gsub("n[-]|%S+{%s+}[-]|[-]|%S+|\\n|=|{|%S+}|",
38    function(k,x) the[k]=coerce(x) end)
39
40  -- Update settings from values on command-line flags. Booleans need no values
41  -- (we just flip the defaults).
42  function cli(t)
43    for slot,v in pairs(t) do
44      v = tostring(v)
45      for n,x in ipairs(arg) do
46        if x=="-." (slot:sub(1,1)) or x=="-."..slot then
47          v = v=="false" and "true" or v=="true" and "false" or arg[n+1] end end
48        t[slot] = coerce(v) end
49      if t.help then os.exit(print("n"..help.."n")) end
50      return t end
51
52  -- ## Linting code
53  -- ## Lists
54  local copy,per,push,sv
55  -- deepcopy
56  function copy(t, u)
57    if type(t) == "table" then return t end
58    u={}; for k,v in pairs(t) do u[k]=copy(v) end
59    return setmetatable(u, getmetatable(t)) end
60
61  -- Return the 'p'-th thing from the sorted list 't'.
62  function per(t,p)
63    p=math.floor(((p or .5)*#t)+.5); return t[math.max(1,math.min(#t,p))] end
64
65  -- Add to 't', return 'x'.
66  function push(t,x) t[#t+1]=x; return x end
67
68  -- ## Call 'fun' on each row. Row cells are divided in 'the.separator'.
69  function csv(fname,fun, sep,src,s,t)
70    sep = "([\\n] .. the.separator .. ")]"
71    src = io.input(fname)
72    while true do
73      s = io.read()
74      if not s then return io.close(src) else
75        t={}
76        for sl in s:gmatch(sep) do t[#t+1] = coerce(sl) end
77        fun(t) end end end
78
79  -- ## Strings
80  local o,oo
81  -- 'o' is a telescope and 'oo' are some binoculars we use to exam stucts.
82  -- 'o': generates a string from a nested table.
83  function o(t, show,u)
84    if type(t) == "table" then return tostring(t) end
85    function show(k,v)
86      if not tostring(k):find"^_" then
87        v = o(v)
88        return #t==0 and string.format("%s%s",k,v) or tostring(v) end end
89    u={}; for k,v in pairs(t) do u[#u+1] = show(k,v) end
90    if #t==0 then table.sort(u) end
91    return "{ "..table.concat(u, ", " ) .. "}" end
92
93  -- 'oo': prints the string from 'o'.
94  function oo(t) print(o(t)) return t end
95
96  -- ## Misc
97  local roques, rnd, obj
98  -- Find rogue locals.
99  function roques()
100    for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end end
101
102  -- ## Maths
103  function rnd(x, places)
104    local mult = 10^(places or 2)
105    return math.floor(x * mult + 0.5) / mult end
106
107  -- obj("Thing") enables a constructor Thing:new() ... and a pretty-printer
108  -- for Things.
109  function obj(s, t,i,new)
110    function new(k,...) i=setmetatable({},k);
111      return setmetatable(t,new(i,...) or i,k) end
112    t[["_tostring"]] = function(x) return s..o(x) end
113    t[["_index"]] = t;return setmetatable(t,{__call=new}) end

```

Aug 27, 22 19:08

csv.lua

Page 2/4

```

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, names=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={} -- dependent 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"^[a-z]" then -- some columns are skipped
146        push(s:find"[a-z]" and self.y or self.x, col) -- some cols are goal cols
147      if s:find"^[a-z]" 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    isEval=false -- true if y-values evaluated.
153  } end
154
155  -- 'Data' is a holder of 'rows' and their summaries (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
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 = fun(n/self.n) end end
175    return e end
176
177  -- ## Num
178  -- Return kept numbers, sorted.
179  function Num:nums()
180    if not self.isSorted then table.sort(self._has); self.isSorted=true end
181    return self._has end
182
183  -- Reservoir sampler. Keep at most 'the.nums' numbers
184  -- (and if we run out of room, delete something old, at random).,
185  function Num:add(v, pos)
186    if v=="?" then
187      self.n = self.n + 1
188      self.lo = math.min(v, self.lo)
189      self.hi = math.max(v, self.hi)
190      if #self._has < the.nums then pos = 1 + (#self._has)
191      elseif math.random() < the.nums/self.n then pos = math.random(#self._has) end
192      if pos then self.isSorted = false
193        self._has[pos] = tonumber(v) end end end
194
195  -- Diversity (standard deviation for Nums, entropy for Syms)
196  function Num:div(a) a=self:nums(); return (per(a,.9)-per(a,.1))/2.58 end
197
198  -- Central tendency (median for Nums, mode for Syms)
199  function Num:mid() return per(self:nums(),.5) end

```

Aug 27, 22 19:08

csv.lua

Page 3/4

```

201  -- ## Data
202  -- Add a 'row' to 'data'. Calls 'add()' to update the 'cols' with new values.
203  function Data:add(xs, row)
204    if not self.cols
205      then self.cols = Cols(xs)
206      else row=push(self.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
207      for _,col in pairs(self.cols.x, self.cols.y) 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      vtype(v)=="number" and rnd(v,places) or v
217      t[col.name]=v end; return t end
218
219
220

```

Aug 27, 22 19:08

csv.lua

Page 4/4

```

221 -----
222
223 -- ## Test Engine
224 local eg, fails = {},0
225
226 -- 1. reset random number seed before running something.
227 -- 2. Cache the defaults settings, and...
228 -- 3. ... restore them after the test
229 -- 4. Print error messages or stack dumps as required.
230 -- 5. Return true if this all went well
231 local function runs(k, old,status,out,msg)
232   if not eg[k] then return end
233   math.randomseed(the.seed) -- reset seed [1]
234   olde={}; for k,v in pairs(the) do old[k]=v end -- [2]
235   if the.dump then -- [4]
236     status,out = true, eg[k]()
237   else
238     status,out = pcall(eg[k]) -- pcall means we do not crash and dump on error
239   end
240   for k,v in pairs(old) do the[k]=v end -- restore old settings [3]
241   msg = status and ((out==true and "PASS") or "FAIL") or "CRASH" -- [4]
242   print("!!!!", msg, k, status)
243   return out or err end
244
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("%10s",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("-----")
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 comand 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","a","b","b","c"} do sym:add(x) end
279   mode, entropy = sym:mid(), sym:div()
280   entropy = (1000*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/auto093.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/auto093.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/auto093.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)

```