

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
10  -d --dump     on test failure, exit with stack dump = nothing
11  -f --file     file with csv data = ../data/aut93.csv
12  -h --help     show help
13  -n --nums     number of nums to keep = 512
14  -s --seed     random number seed = 10019
15  -S --separator field 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 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.."n")) 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
77
78  -- ## Strings
79  local o,oo
80  -- 'o' is a telescope and 'oo' are some binoculars we use to exam stucts.
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 roques, rnd, obj
97  -- Find rogue locals.
98  function roques()
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
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 summaries (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:do(row) end)
160   else for _,row in pairs(src or {}) do self:do(row) end end end
161
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 -- -----
178 -- ## Num
179 -- Return kept numbers, sorted.
180 function Num:nums()
181   if not self.isSorted then table.sort(self._has); self.isSorted=true end
182   return self._has end
183
184 -- Reservoir sampler. Keep at most 'the.nums' numbers
185 -- (and if we run out of room, delete something old, at random).,
186 function Num:add(v, pos)
187   if v=="?" then
188     self.n = self.n + 1
189     self.lo = math.min(v, self.lo)
190     self.hi = math.max(v, self.hi)
191     if #self._has < the.nums then pos = 1 + (#self._has)
192     elseif math.random() < the.nums/self.n then pos = math.random(#self._has) end
193     if pos then self.isSorted = false
194       self._has[pos] = tonumber(v) end end end
195
196 -- Diversity (standard deviation for Nums, entropy for Sym)
197 function Num:div(a) a=self:nums(); return (per(a,.9)-per(a,1))/2.58 end
198
199 -- Central tendency (median for Nums, mode for Sym)
200 function Num:mid() return per(self:nums(),.5) end
201
202 -- -----
203 -- ## Data
204 -- Add a 'row' to 'data'. Calls 'add()' to updatie the 'cols' with new values.
205 function Data:add(xs, row)
206   if not self.cols
207   then self.cols = Cols(xs)
208   else row = push(self.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
209     for _,todo in pairs(self.cols.x, self.cols.y) do
210       for _,col in pairs(todo) do
211         col:add(row.cells[col.at]) end end end end
212
213 -- For 'showCols' (default="data.cols.x") in 'data', report 'fun' (default="mid"),
214 -- rounding numbers to 'places',showCols,fun, t,v)
215 function Data:stats(places,showCols,fun, t,v)
216   showCols, fun = showCols or self.cols.y, fun or "mid"
217   t={}; for _,col in pairs(showCols) do
218     v=fun(col)
219     v=type(v)=="number" and rnd(v,places) or v
220     t[col.name]=v end; return t end

```

Aug 27, 22 19:06

csv.lua

Page 3/3

```

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   old={}; 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 -- ## Tests
247 -- Test that the test happens when something crashes?
248 function eg.BAD() print(eg.dont.have.this.field) end
249
250 -- Sort all test names.
251 function eg.LIST(t)
252   t={}; for k,_ in pairs(eg) do t[#t+1]=k end; table.sort(t); return t end
253
254 -- List test names.
255 function eg.LS()
256   print("nExamples lua csv-e..")
257   for _,k in pairs(eg.LIST()) do print(string.format("%s",k)) end
258   return true end
259
260 -- Run all tests
261 function eg.ALL()
262   for _,k in pairs(eg.LIST()) do
263     if k ~= "ALL" then
264       print("n-----")
265       if not runs(k) then fails=fails+ 1 end end
266     return true end
267
268 -- Settings come from big string top of "sam.lua"
269 -- (maybe updated from command line)
270 function eg.the() oo(the); return the end
271
272 -- The middle and diversity of a set of symbols is called "mode"
273 -- and "entropy" (and the latter is zero when all the symbols
274 -- are the same).
275 function eg.sym( sym,entropy,mode)
276   sym= Sym()
277   for _,x in pairs("a","a","a","a","b","b","c") do sym:add(x) end
278   mode, entropy = sym:mid(), sym:div()
279   entropy = (1000*entropy)/1/1000
280   oo((midmode, div=entropy))
281   return mode=="a" and 1.37 <= entropy and entropy <=1.38 end
282
283 -- The middle and diversity of a set of numbers is called "median"
284 -- and "standard deviation" (and the latter is zero when all the nums
285 -- are the same).
286 function eg.num( num,mid,div)
287   num=Num()
288   for i=1,100 do num:add(i) end
289   oo(num:nums())
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   the.nums = 32
298   for i=1,1000 do num:add(i) end
299   oo(num:nums())
300   return 32==#num._has; end
301
302 -- Show we can read csv files.
303 function eg.csv( n)
304   n=0
305   csv("../data/aut93.csv",function(row)
306     n=n+1; if n> 10 then return else oo(row) end end); return true end
307
308 -- Can I load a csv file into a Data?.
309 function eg.data( d)
310   d = Data("../data/aut93.csv")
311   for _,col in pairs(d.cols.y) do oo(col) end
312   return true
313 end
314
315 -- Print some stats on columns.
316 function eg.stats( data,mid,div)
317   data = Data("../data/aut93.csv")
318   div=function(col) return col:div() end
319   mid=function(col) return col:mid() end
320   print("xmid", o( data:stats(2,data.cols.x, mid)))
321   print("xdiv", o( data:stats(2,data.cols.x, div)))
322   print("ymid", o( data:stats(2,data.cols.y, mid)))
323   print("ydiv", o( data:stats(2,data.cols.y, div)))
324   return true
325 end
326
327 -- -----
328 the = cli(the)
329 runs(the.eg)
330 roques()
331 os.exit(fails)

```