

Aug 27, 22 8:11

csv.lua

Page 1/4

```

1 local help=[
2 SEEN : summarized csv file
3 (c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license
4
5 USAGE: lua seen.lua [OPTIONS]
6
7 OPTIONS:
8 -e --eg          start-up example          = nothing
9 -d --dump        on test failure, exit with stack dump = false
10 -f --file        file with csv data         = ../data/auto93.csv
11 -h --help        show help                  = false
12 -n --nums        number of nums to keep     = 512
13 -s --seed        random number seed        = 10019
14 -S --separator  feild seperator            = ,
15 ]]
16
17 -- ## Misc routines
18 -- ### Linting code
19 -- Find rogue locals.
20 local b4={}; for k,v in pairs(_ENV) do b4[k]=v end
21 local function rogues()
22   for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end end
23
24 -- ### Handle Settings
25 -- Parse 'the' config settings from 'help'.
26 local the={}
27 local function coerce(s)
28   local function coercel(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 coercel(s:match"^%s*(.)%s*$") end
33
34 help:gsub("\n[~]([%S]+[%s]+)[~]([~]([%S]+)[^~]+=([%S]+)",
35   function(k,x) the[k]=coerce(x) end)
36
37 -- Update settings from values on command-line flags. Booleans need no values
38 -- (we just flip the defaults).
39 local function cli(t)
40   for slot,v in pairs(t) do
41     v = tostring(v)
42     for n,x in ipairs(arg) do
43       if x=="-.."(slot:sub(1,1)) or x=="-.."slot then
44         v = v=="false" and "true" or v=="true" and "false" or arg[n+1] end end
45     t[slot] = coerce(v) end
46   if t.help then os.exit(print("\n"..help.."")) end
47   return t end
48
49 -- ### Strings
50 -- 'o' generates a string from a nested table.
51 local function o(t)
52   if type(t) ~= "table" then return tostring(t) end
53   local function show(k,v)
54     if not tostring(k):find"^_" then
55       v = o(v)
56       return #t==0 and string.format(":%s %s",k,v) or tostring(v) end end
57   local u={}; for k,v in pairs(t) do u[1+#u] = show(k,v) end
58   if #t==0 then table.sort(u) end
59   return (t._is or "").."["..table.concat(u, " ").."]" end
60
61 -- 'oo' prints the string from 'o'.
62 local function oo(t) print(o(t)) return t end
63
64 -- ### Lists
65 -- Deepcopy
66 local function copy(t)
67   if type(t) ~= "table" then return t end
68   local u={}; for k,v in pairs(t) do u[k] = copy(v) end
69   return setmetatable(u, getmetatable(t)) end
70
71 -- Return the 'p'-th thing from the sorted list 't'.
72 local function per(t,p)
73   p=math.floor(((p or .5)*#t)+.5); return t[math.max(1,math.min(#t,p))] end
74
75 -- Add to 't', return 'x'.
76 local function push(t,x) t[1+#t]=x; return x end
77
78 -- ## Call 'fun' on each row. Row cells are divided in 'the.seperator'.
79 local function csv(fname,fun)
80   local sep = "([^\n] .. the.seperator .. ")+)"
81   local src = io.input(fname)
82   while true do
83     local s = io.read()
84     if not s then return io.close(src) else
85       local t={}
86       for s1 in s:gmatch(sep) do t[1+#t] = coerce(s1) end
87       fun(t) end end end

```

Aug 27, 22 8:11

csv.lua

Page 2/4

```

88 -- -----
89 -- ## Objects
90 local Data, Cols, Sym, Num, Row
91
92 -- 'Data' is a holder of 'rows' and their summaries (in 'cols').
93 function Data() return {_is = "Data",
94   cols= nil, -- summaries of data
95   rows= {} -- kept data
96 } end
97
98 -- 'Columns' Holds of summaries of columns.
99 -- Columns are created once, then may appear in multiple slots.
100 function Cols() return {
101   _is = "Cols",
102   names={}, -- all column names
103   all={}, -- all the columns (including the skipped ones)
104   klass=nil, -- the single dependent klass column (if it exists)
105   x={}, -- independent columns (that are not skipped)
106   y={} -- dependet columns (that are not skipped)
107 } end
108
109 -- 'Sym' summarizes a stream of symbols.
110 function Sym(c,s)
111   return {_is="Sym",
112     n=0, -- items seen
113     at=c or 0, -- column position
114     name=s or "", -- column name
115     _has={} -- kept data
116 } end
117
118 -- 'Num' ummarizes a stream of numbers.
119 function Num(c,s)
120   return {_is="Nums",
121     n=0, at=c or 0, name=s or "", _has={}, -- as per Sym
122     isNum=true, -- mark that this is a number
123     lo= math.huge, -- lowest seen
124     hi= -math.huge, -- highest seen
125     isSorted=true, -- no updates since last sort of data
126     w = ((s or ""):find"$" and -1 or 1)
127 } end
128
129 -- 'Row' holds one record
130 function Row(t) return {_is="Row",
131   cells=t, -- one record
132   cooked=copy(t), -- used if we discretize data
133   isEvald=false -- true if y-values evaluated.
134 } end
135
136 -- ## Data
137 -- Add one thing to 'col'. For Num, keep at most 'nums' items.
138 local function add(col,v)
139   if v=="?" then
140     col.n = col.n + 1
141     if not col.isNum then col._has[v] = 1 + (col._has[v] or 0) else
142       col.lo = math.min(v, col.lo)
143       col.hi = math.max(v, col.hi)
144     local pos
145     if #col._has < the.nums then pos = 1 + (#col._has)
146     elseif math.random() < the.nums/col.n then pos = math.random(#col._has) end
147     if pos then col.isSorted = false
148       col._has[pos] = tonumber(v) end end end end
149
150 local function adds(col,t) for _,x in pairs(t) do add(col,x) end; return col end
151
152 -- Add a 'row' to 'data'. Calls 'add()' to update the 'cols' with new values.
153 local function record(data,xs)
154   local row= push(data.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
155   for _,todo in pairs(data.cols.x, data.cols.y) do
156     for _,col in pairs(todo) do
157       add(col, row.cells[col.at]) end end end
158
159 -- Generate rows from some 'src'. If 'src' is a string, read rows from file;
160 -- else read rows from a 'src' table. When reading, use row1 to define columns.
161 local function records(src, data, head, body)
162   function head(sNames)
163     local cols = Cols()
164     cols.names = sNames
165     for c,s in pairs(sNames) do
166       local col = push(cols.all, -- Numerics start with Uppercase.
167         (s:find"^[A-Z]" and Num or Sym)(c,s))
168       if not s:find"$" then -- some columns are skipped
169         push(s:find"[+-]" and cols.y or cols.x, col) -- some cols are goal cols
170       if s:find"$" then cols.klass=col end end end
171     return cols
172   end
173   function body(t) -- treat first row differently (defines the columns)
174     if data.cols then record(data,t) else data.cols=head(t) end
175   end
176   data = Data()
177   if type(src)=="string" then csv(src, body) else

```

Aug 27, 22 8:11

csv.lua

Page 3/4

```

178   for _,t in pairs(src or {}) do body(t) end end
179   return data end
180
181 -- ### Query
182 -- Return kept numbers, sorted.
183 local function nums(num)
184   if not num.isSorted then table.sort(num._has); num.isSorted=true end
185   return num._has end
186
187 -- Diversity (standard deviation for Nums, entropy for Syms)
188 local function div(col)
189   if col.isNum then local a=nums(col); return (per(a,.9)-per(a,.1))/2.58 else
190     local function fun(p) return p*math.log(p,2) end
191     local e=0
192     for _,n in pairs(col._has) do if n>0 then e=e-fun(n/col.n) end end
193     return e end end
194
195 -- Central tendency (median for Nums, mode for Syms)
196 local function mid(col)
197   if col.isNum then return per(nums(col),.5) else
198     local most,mode = -1
199     for k,v in pairs(col._has) do if v>most then mode,most=k,v end end
200     return mode end end
201
202 -- Diversity (standard deviation for Nums, entropy for Syms)
203 local function div(col)
204   if col.isNum then local a=nums(col); return (per(a,.9)-per(a,.1))/2.58 else
205     local function fun(p) return p*math.log(p,2) end
206     local e=0
207     for _,n in pairs(col._has) do if n>0 then e=e-fun(n/col.n) end end
208     return e end end
209
210 -- For 'showCols' (default='data.cols.x') in 'data', report 'fun' (default='mid').
211 local function stats(data, showCols,fun, t)
212   showCols, fun = showCols or data.cols.y, fun or mid
213   t={}; for _,col in pairs(showCols) do t[col.name]=fun(col) end; return t end

```

Aug 27, 22 8:11

csv.lua

Page 4/4

```

214 -----
215 local eg, fails = {},0
216
217 local function runs(k, old,status,out,msg)
218   if not eg[k] then return end
219   math.randomseed(the.seed) -- reset seed
220   old={}; for k,v in pairs(the) do old[k]=v end
221   if the.dump then
222     status,out = true, eg[k]()
223   else
224     status,out = pcall(eg[k]) -- pcall sets status=false if crash
225   end
226   for k,v in pairs(old) do the[k]=v end -- restore old settings
227   msg = status and ((out==true and "PASS") or "FAIL") or "CRASH"
228   print("!!!!!!", msg, k, status)
229   return out or err end
230
231 function eg.BAD() print(eg.ab.sent) end
232
233 function eg.LIST(t)
234   t={}; for k,_ in pairs(eg) do t[l+#t]=k end; table.sort(t); return t end
235
236 function eg.LS()
237   print("\nExamples lua csv -e...")
238   for _,k in pairs(eg.LIST()) do print(string.format("%t%s",k)) end
239   return true end
240
241 function eg.ALL()
242   for _,k in pairs(eg.LIST()) do
243     if k ~= "ALL" then
244       print("\n-----")
245       if not runs(k) then fails=fails+ 1 end end end
246   return true end
247
248 -- Settings come from big string top of "sam.lua"
249 -- (maybe updated from comandnd line)
250 function eg.the() oo(the); return true end
251
252 -- The middle and diversity of a set of symbols is called "mode"
253 -- and "entropy" (and the latter is zero when all the symbols
254 -- are the same).
255 function eg.sym( sym,entropy,mode)
256   sym= adds(Sym(), {"a","a","a","a","b","b","c"})
257   mode, entropy = mid(sym), div(sym)
258   entropy = (1000*entropy)//1/1000
259   oo({mid=mode, div=entropy})
260   return mode=="a" and 1.37 <= entropy and entropy <=1.38 end
261
262 -- The middle and diversity of a set of numbers is called "median"
263 -- and "standard deviation" (and the latter is zero when all the nums
264 -- are the same).
265 function eg.num( num)
266   num=Num()
267   for i=1,100 do add(num,i) end
268   local med,ent = mid(num), div(num)
269   print(mid(num),div(num))
270   return 50<= med and med<= 52 and 30.5 <ent and ent <32 end
271
272 -- Nums store only a sample of the numbers added to it (and that storage
273 -- is done such that the kept numbers span the range of inputs).
274 function eg.bignum( num)
275   num=Num()
276   the.nums = 32
277   for i=1,1000 do add(num,i) end
278   oo(nums(num))
279   return 32==#num._has; end
280
281 -- Show we can read csv files.
282 function eg.csv()
283   csv("../data/auto93.csv",oo); return true end
284
285 -- Print some stats on columns.
286 function eg.stats()
287   oo(stats(records("../data/auto93.csv"))); return true end
288
289 -----
290 the = cli(the)
291 runs(the.eg)
292 rogues()
293 os.exit(fails)

```