

Aug 27, 22 7:22

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 -f --file    file with csv data    = ../data/auto93.csv
10 -h --help    show help              = false
11 -n --nums    number of nums to keep = 512
12 -s --seed    random number seed    = 10019
13 -S --separator feild separator      = ,]
14
15 -- ## Misc routines
16 -- ### Linting code
17 -- Find rogue locals.
18 local b4={}; for k,v in pairs(_ENV) do b4[k]=v end
19 local function rogues()
20   for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end end
21
22 -- ### Handle Settings
23 -- Parse 'the' config settings from 'help'.
24 local the={}
25 local function coerce(s)
26   local function coercel(s1)
27     if s1=="true" then return true end
28     if s1=="false" then return false end
29     return s1 end
30   return math.tointeger(s) or tonumber(s) or coercel(s:match("^%s*(-)%s*$") end
31
32 help:gsub("ln[-][%S]+[%s]+[-][-][%S]+[^\n]+([%S]+)",
33   function(k,x) the[k]=coerce(x) end)
34
35 -- Update settings from values on command-line flags. Booleans need no values
36 -- (we just flip the defaults).
37 local function cli(t)
38   for slot,v in pairs(t) do
39     v = tostring(v)
40     for n,x in pairs(arg) do
41       if x=="-.."(slot:sub(1,1)) or x=="-.."slot then
42         v = v=="false" and "true" or v=="true" and "false" or arg[n+1] end end
43     t[slot] = coerce(v) end
44   if t.help then os.exit(print("ln"..help.."ln")) end
45   return t end
46
47 -- ## Strings
48 -- 'o' generates a string from a nested table.
49 local function o(t)
50   if type(t) ~= "table" then return tostring(t) end
51   local function show(k,v)
52     if not tostring(k):find"^_" then
53       v = o(v)
54       return #t==0 and string.format(":%s %s",k,v) or tostring(v) end end
55   local u={}; for k,v in pairs(t) do u[1+#u] = show(k,v) end
56   if #t==0 then table.sort(u) end
57   return (t._is or "").."["..table.concat(u, " ").."]" end
58
59 -- 'oo' prints the string from 'o'.
60 local function oo(t) print(o(t)) return t end
61
62 -- ## Lists
63 -- Deepcopy
64 local function copy(t)
65   if type(t) ~= "table" then return t end
66   local u={}; for k,v in pairs(t) do u[k] = copy(v) end
67   return setmetatable(u, getmetatable(t)) end
68
69 -- Return the 'p'-th thing from the sorted list 't'.
70 local function per(t,p)
71   p=math.floor(((p or .5)*#t)+.5); return t[math.max(1,math.min(#t,p))] end
72
73 -- Add to 't', return 'x'.
74 local function push(t,x) t[1+#t]=x; return x end
75
76 -- ## Call 'fun' on each row. Row cells are divided in 'the.sep'.
77 local function csv(fname,fun)
78   local sep = "[^" .. the.sep .. "]"
79   local src = io.input(fname)
80   while true do
81     local s = io.read()
82     if not s then return io.close(src) else
83       local t={}
84       for s1 in s:gmatch(sep) do t[1+#t] = coerce(s1) end
85       fun(t) end end end

```

Aug 27, 22 7:22

csv.lua

Page 2/4

```

86 -- -----
87 -- ## Objects
88 local Data, Cols, Sym, Num, Row
89
90 -- 'Data' is a holder of 'rows' and their summaries (in 'cols').
91 function Data() return {_is = "Data",
92   cols=nil, -- summaries of data
93   rows={}, -- kept data
94 } end
95
96 -- 'Columns' Holds of summaries of columns.
97 -- Columns are created once, then may appear in multiple slots.
98 function Cols() return {
99   _is = "Cols",
100   names={}, -- all column names
101   all={}, -- all the columns (including the skipped ones)
102   klass=nil, -- the single dependent klass column (if it exists)
103   x={}, -- independent columns (that are not skipped)
104   y={} -- dependent columns (that are not skipped)
105 } end
106
107 -- 'Sym' summarizes a stream of symbols.
108 function Sym(c,s)
109   return {_is="Sym",
110     n=0, -- items seen
111     at=c or 0, -- column position
112     name=s or "", -- column name
113     _has={} -- kept data
114 } end
115
116 -- 'Num' ummarizes a stream of numbers.
117 function Num(c,s)
118   return {_is="Nums",
119     n=0, at=c or 0, name=s or "", _has={}, -- as per Sym
120     isNum=true, -- mark that this is a number
121     lo= math.huge, -- lowest seen
122     hi= -math.huge, -- highest seen
123     isSorted=true, -- no updates since last sort of data
124     w = ((s or ""):find"$" and -1 or 1)
125 } end
126
127 -- 'Row' holds one record
128 function Row(t) return {_is="Row",
129   cells=t, -- one record
130   cooked=copy(t), -- used if we discretize data
131   isEval=false -- true if y-values evaluated.
132 } end
133
134 -- ## Data
135 -- Add one thing to 'col'. For Num, keep at most 'nums' items.
136 local function add(col,v)
137   if v=="?" then
138     col.n = col.n + 1
139     if not col.isNum then col._has[v] = 1 + (col._has[v] or 0) else
140       col.lo = math.min(v, col.lo)
141       col.hi = math.max(v, col.hi)
142       local pos
143       if #col._has < the.nums then pos = 1 + (#col._has)
144       elseif math.random() < the.nums/col.n then pos = math.random(#col._has) end
145       if pos then col.isSorted = false
146       if pos then col._has[pos] = tonumber(v) end end end end
147
148 local function adds(col,t) for _,x in pairs(t) do add(col,x) end; return col end
149
150 -- Add a 'row' to 'data'. Calls 'add()' to update the 'cols' with new values.
151 local function record(data,xs)
152   local row= push(data.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
153   for _,todo in pairs(data.cols.x, data.cols.y) do
154     for _,col in pairs(todo) do
155       add(col, row.cells[col.at]) end end end
156
157 -- Generate rows from some 'src'. If 'src' is a string, read rows from file;
158 -- else read rows from a 'src' table. When reading, use row1 to define columns.
159 local function records(src, data, head, body)
160   function head(sNames)
161     local cols = Cols()
162     cols.names = sNames
163     for c,s in pairs(sNames) do
164       local col = push(cols.all, -- Numerics start with Uppercase.
165         (s:find"^[A-Z]" and Num or Sym)(c,s))
166       if not s:find"$" then -- some columns are skipped
167         push(s:find"[+-]" and cols.y or cols.x, col) -- some cols are goal cols
168         if s:find"$" then cols.klass=col end end end
169     return cols
170   end
171   function body(t) -- treat first row differently (defines the columns)
172     if data.cols then record(data,t) else data.cols=head(t) end
173   end
174   data = Data()
175   if type(src)=="string" then csv(src, body) else

```

Aug 27, 22 7:22

csv.lua

Page 3/4

```

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

```

Aug 27, 22 7:22

csv.lua

Page 4/4

```

213 -----
214 local eg, fails = {},0
215
216 local function runs(k)
217   if not eg[k] then return end
218   math.randomseed(the.seed) -- reset seed
219   local old={}; for k,v in pairs(the) do old[k]=v end
220   local out=eg[k]()
221   for k,v in pairs(old) do the[k]=v end -- restore old settings
222   print("!!!!!!", k, out and "PASS" or "FAIL")
223   return out end
224
225 function eg.FAIL() print(eg.ab.sent) end
226
227 function eg.LIST( t)
228   t={}; for k,_ in pairs(eg) do t[1+#t]=k end; table.sort(t); return t end
229
230 function eg.LS()
231   print("\nExamples lua csv -e ...")
232   for _,k in pairs(eg.LIST()) do print(string.format("%t%s",k)) end
233   return true end
234
235 function eg.ALL()
236   for _,k in pairs(eg.LIST()) do
237     if k ~= "ALL" then
238       print("\n-----")
239       local status, err = pcall(function () runs(k) end)
240       if not status or err then fails=fails+1 end end end
241   return true end
242
243 -- Settings come from big string top of "sam.lua"
244 -- (maybe updated from comand line)
245 function eg.the() oo(the); return true end
246
247 -- The middle and diversity of a set of symbols is called "mode"
248 -- and "entropy" (and the latter is zero when all the symbols
249 -- are the same).
250 function eg.ent( sym,ent)
251   sym= adds(Sym(), {"a","a","a","a","b","b","c"})
252   ent= div(sym)
253   print(ent,mid(sym))
254   return 1.37 <= ent and ent <=1.38 end
255
256 -- The middle and diversity of a set of numbers is called "median"
257 -- and "standard deviation" (and the latter is zero when all the nums
258 -- are the same).
259 function eg.num( num)
260   num=Num()
261   for i=1,100 do add(num,i) end
262   local med,ent = mid(num), div(num)
263   print(mid(num),div(num))
264   return 50<= med and med<= 52 and 30.5 <ent and ent <32 end
265
266 -- Nums store only a sample of the numbers added to it (and that storage
267 -- is done such that the kept numbers span the range of inputs).
268 function eg.bignum( num)
269   num=Num()
270   the.nums = 32
271   for i=1,1000 do add(num,i) end
272   oo(nums(num))
273   return 32==#num._has; end
274
275 -- Show we can read csv files.
276 function eg.csv()
277   csv("../data/auto93.csv",oo); return true end
278
279 -- Print some stats on columns.
280 function eg.stats()
281   oo(stats(records("../data/auto93.csv"))); return true end
282
283 -----
284 the = cli(the)
285 runs(the.eg)
286 rogues()
287 os.exit(fails)

```