

```

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225

```

SAM : Semi-supervised And Multi-objective explanations  
(c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license

eg.lua

```

-- In this code:
-- Line strive to be 80 chars (or less)
-- Two spaces before function arguments denote optionals.
-- Four spaces before function arguments denote local variables.
-- Private functions start with '_'
-- Arguments of private functions do anything at all
-- Local variables inside functions do anything at all
-- Arguments of public functions use type hints
-- Variable 'x' is anything
-- Prefix 'is' is a boolean
-- Prefix 'fun' is a function
-- Prefix 'f' is a filename
-- Prefix 'n' is a string
-- Prefix 's' is a string
-- Prefix 'c' is a column index
-- 'col' denotes 'num' or 'sym'
-- 'x' is anything (table or number of boolean or string)
-- 'v' is a simple value (number or boolean or string)
-- Suffix 'a' is a list of things
-- Tables are 't' or, using the above, a table of numbers would be 'ns'
-- Type names are lower case versions of constructors; e.g 'col' isa 'Cols'.

-- All demo functions 'eg.fun1' can be called via 'lua eg.lua -e fun1'.
local eg = {}

local l=require"lib"
local _require"sam"
local o,oo,per,push,rnd = 1,0.1,0.0,1,per,1,push,1,rnd
local add,adds,dist,div = _add,_adds,_dist,_div
local mid, records, the = _mid,_records,_the
local Num, Sym = _Num, _Sym

-- Settings come from big string top of "sam.lua"
-- (maybe updated from command line)
function eg.the(i) oo(the); return true end

-- The middle and diversity of a set of symbols is called "mode"
-- and "entropy" (and the latter is zero when all the symbols
-- are the same).
function eg.ent( sym,ent)
  sym= adds(Sym(), { "a","A","a","a","b","b","b","c","c" })
  ent= div(sym)
  print(ent,mid(sym))
  return 1.37 <= ent and ent <= 1.38 end

-- The middle and diversity of a set of numbers is called "median"
-- and "standard deviation" (and the latter is zero when all the nums
-- are the same).
function eg.num( num)
  num=Num()
  for i=1,100 do add(num,i) end
  local med,ent = mid(num), rnd(div(num),2)
  print(mid(num),rnd(div(num),2))
  return 50<= med and med<= 52 and 30.5 <=ent and ent <32 end

-- Nums store only a sample of the numbers added to it (and that storage
-- is done such that the kept numbers span the range of inputs).
function eg.bignum( num)
  num=Num()
  the.nums = 32
  for i=1,1000 do add(num,i) end
  oo(_nums(num))
  return 32==#num._has end

-- We can read data from disk-based csv files, where row1 lists a
-- set of columns names. These names are used to work out what are Nums, or
-- ro Syms, or goals to minimize/maximize, or (indeed) what columns to ignore.
function eg.records()
  oo(records("../data/auto93.csv").cols.y); return true end

-- Any two rows have a distance 0..1 that satisfies equality, symmetry
-- and the triangle inequality.
function eg.dist( data,t)
  data=records("../data/auto93.csv")
  t={}
  for i=1,100 do
    local A,B,C = 1.any(data.rows), 1.any(data.rows), 1.any(data.rows)
    local a,b,c = dist(data,B,C), dist(data,A,C), dist(data,A,B)
    assert(a<=1 and b<=1 and c<=1)
    assert(a>=0 and b>=0 and c>=0)
    assert( dist(data,A,A) = 0) -- equality
    assert( dist(data,A,B) = dist(data,B,A)) -- symmetry
    assert( a+b<=c) -- triangle inequality
    for _,x in pairs(a) do push(t,rnd(x,2)) end end
  table.sort(t)
  oo(t)
  return true end

the = l.cli(the)
os.exit(1,runs(the.eg, eg, the))

```

```

107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225

```

eg.lua (https://github.com/timm/lua/blob/main/src/sam/eg.lua).

SAM : Semi-supervised And Multi-objective explanations  
(c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license

USAGE: lua eg.lua [OPTIONS]

```

OPTIONS:
  -b --bins      number of bins      = 8
  -c --cohen     small effect        = .35
  -e --eg       start-up example     = nothing
  -f --far       far away            = .95
  -f --file      file with csv data   = ../docs/auto93.csv
  -h --help      show help           = false
  -m --min       min size = n(the.min) = 5
  -n --nums      how many numbers to keep = 256
  -p --p         distance coefficient = 2
  -s --seed      random number seed   = 10019
  -S --sample    how many rows to search = 512)))
  -- Commonly used lib functions.
local lt,o,oo,per,push,sort = 1,lt,1.0,1.0,0.1,per,1,push,1,sort

----- Classes
local Data,Cols,Sym,Num,Row
-- Holder of 'rows' and their summaries (in 'cols').
function Data(i) return { _is = "Data", -- summaries of data
  rows = {}, -- kept data
} end

-- Holds of summaries of columns.
-- Columns are created once, then may appear in multiple slots.
function Cols() return {
  _is = "Cols",
  names={}, -- all column names
  all={} -- the columns (including the skipped ones)
  klass=nil, -- the single dependent klass column (if it exists)
  x={}, -- independent columns (that are not skipped)
  y={} -- dependent columns (that are not skipped)
} end

-- Summarizers a stream of symbols.
function Sym(c,s)
  return { _is = "Sym",
    n=0, -- items seen
    at=c or 0, -- column position
    names=s or "", -- column name
    _has={} -- kept data
  } end

-- Summarizes a stream of numbers.
function Num(c,s)
  return { _is="Nums",
    n=0,at=c or 0, names=s or "", _has={}, -- as per Sym
    isNum=true, -- mark that this is a number
    lo= math.huge, -- lowest seen
    hi= -math.huge, -- highest seen
    isSorted=true, -- no updates since last sort of data
    w=(s or ""):find"$" and -1 or 1 -- minimizing if w=-1
  } end

-- Holds one record
function Row(t) return { _is="Row",
  cells=t,
  cooked=i.copy(t) -- used if we discretize data
} end

----- Data Functions
local add,adds,clone,div,mid,norm,nums,record,records,stats
-- Create
-- Generate rows from some 'src'. If 'src' is a string, read rows from file;
-- else read rows from a 'src' table. When reading, use row1 to define columns.
function records(src, data,head,body)
  local function head(sNames)
    local cols = Cols()
    cols.names = names
    for c,s in pairs(sNames) do
      local col = push(cols.all, -- Numerics start with Uppercase.
        (s:find"^[A-Z]" and Num or Sym)(c,s))
      if not s:find"$" then -- some columns are skipped
        push(s:find"[|+]" and cols.y or cols.x, col) -- some cols are goal cols
      if s:find"$" then cols.klass=col end end end
      return cols
    end
    function body(t) -- treat first row differently (defines the columns)
      if data.cols then record(data,t) else data.cols=head(t) end
    end
    data = Data()
    if type(src)=="string" then l.csv(src, body) else
      for _,t in pairs(src or {}) do body(t) end end
    return data end

-- Return a new data with same structure as 'data1'. Optionally, oad in 'rows'.
function clone(data1, rows)
  data2=Data()
  data2.cols = _head(data1.cols.names)
  for _,row in pairs(rows or {}) do record(data2,row) end
  return data2 end

----- Update
-- Add one thing to 'col'. For Num, keep at most 'nums' items.
function add(col,v)
  if v=="?" then
    col.n = col.n + 1
    if not col.isNum then col._has[v] = 1 + (col._has[v] or 0) else
      col.lo = math.min(v, col.lo)
      col.hi = math.max(v, col.hi)
    local pos
    if #col._has < the.nums then pos = 1 + (#col._has)
    elseif math.random() < the.nums/col.n then pos = math.random(#col._has) end
    if pos then col.isSorted = false
      col._has[pos] = tonumber(v) end end end end

```

```

226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344

```

-- Add many things to col

function adds(col,t) for \_,v in pairs(t) do add(col,v) end; return col end

-- Add a 'row' to 'data'. Calls 'add()' to update the 'cols' with new values.

function record(data,xs)
 local row= push(data.rows, xs.ccells and xs or Row(xs)) -- ensure xs is a Row
 for \_,todo in pairs(data.cols.x, data.cols.y) do
 for \_,col in pairs(todo) do
 add(col, row.ccells[col.at]) end end end

----- Query
-- Return kept numbers, sorted.
function nums(num)
 if not num.isSorted then num.\_has = sort(num.\_has); num.isSorted=true end
 return num.\_has end

-- Normalized numbers 0..1. Everything else normalizes to itself.
function norm(col,n)
 return xs=="?" or not col.isNum and x or (n-col.lo)/(col.hi-col.lo + 1E-32) end

-- Diversity (standard deviation for Nums, entropy for Syms)
function div(col)
 if col.isNum then local a=nums(col); return (per(a,.9)-per(a,.1))/2.58 else
 local function fun(p) return p\*math.log(p,2) end
 local e=0
 for \_,n in pairs(col.\_has) do if n>0 then e=-fun(n/col.n) end end
 return e end end

-- Central tendency (median for Nums, mode for Syms)
function mid(col)
 if col.isNum then return per(nums(col),.5) else
 local most,mode = -1
 for k,v in pairs(col.\_has) do if v>most then mode,most=k,v end end
 return mode end end

-- For 'showCols' (default='data.cols.x') in 'data', report 'fun' (default='mid').
function stats(data, showCols, fun, t)
 showCols, fun = showCols or data.cols.y, fun or mid
 t={}; for \_,col in pairs(showCols) do t[col.name]=fun(col) end; return t end

----- Discretization
-- Find ranges within a num (unsupervised).
function bins(num)
 local a, epsilon = nums(num), the.cohen\*div(num)
 local enough = #a>the.min
 local one = {lo=a[1], hi=a[1], n=0}
 local t = {one}
 for i,x in pairs(a) do
 if i < #a-enough and x ~= a[i+1] and n > enough and hi-lo > epsilon then
 one = push(t, {lo=one.hi, hi=a[i], n=0}) end
 one.hi = a[i]
 one.n = 1 + one.n end
 t[1].lo = -math.huge
 t[#t].ho = math.huge
 return t end

-- Fill in discretized values (in 'cooked').
function cook(data)
 for \_,num in pairs(data.cols.x) do
 if num.isNum then local t = bins(num)
 for \_,row in pairs(data.rows) do
 local v = row.ccells[num.at]
 if v ~= "?" then
 for \_,bin in pairs(t) do
 if v > bin.lo and v <= bin.hi then
 row.cooked[col.at] = bin.lo
 break end end end end end end

```

345     push(i < #rows/2 and lefts or rights, pair.row) end
346     return lefts,rights,left,right,c end
347
348 -- Recursively split 'rows' (default='data.rows') in half.
349 function halves(data,rows, rowAbove,stop)
350     rows = rows or data.rows
351     stop = stop or (#rows)^the.min
352     local here = {here=rows}
353     if #rows>stop then
354         local lefts,rights,left,right = half(data,rows,rowAbove)
355         here.lefts = halves(data,lefts,left,stop)
356         here.rights = halves(data,rights,right,stop) end
357     return here end
358
359 -----
360 -- That's all folks.
361 return {
362     the=the, Data=Data, Cols=Cols, Sym=Sym, Num=Num, Row=Row, add=add,
363     adds=adds, around=around, clone=clone, dist=dist, div=div,
364     far=far, half=half, halves=halves, mid=mid, nums=nums, records=records,
365     record=record, stats=stats}
366
367
368
369
370
371 -- lib.lua: misc LUA functions
372 -- (c)2022 Tim Menzies <timm@ieee.org> BSD-2 licence
373 local l={}
374
375 ----- Meta
376 -- Find rogue locals.
377 l.b4={}; for k,v in pairs(_ENV) do l.b4[k]=v end
378 function l.rogues()
379     for k,v in pairs(_ENV) do if not l.b4[k] then print("?",k,type(v)) end end end
380
381 ----- Lists
382 -- Add 'x' to a list. Return 'x'.
383 function l.push(t,x) t[1+#t]=x; return x end
384
385 -- Sample one item
386 function l.any(t) return t[math.random(#t)] end
387
388 -- Sample many items
389 function l.many(t,n, u) u={}; for i=1,n do u[1+#u]=l.any(t) end; return u end
390
391 -- Deepcopy
392 function l.copy(t)
393     if type(t) ~= "table" then return t end
394     local u={}; for k,v in pairs(t) do u[k] = l.copy(v) end
395     return setmetatable(u,getmetatable(t)) end
396
397 -- Round
398 function l.rnd(n, nPlaces)
399     local mult = 10^(nPlaces or 3)
400     return math.floor(n * mult + 0.5) / mult end
401
402 -- Deepcopy
403 function l.copy(t)
404     if type(t) ~= "table" then return t end
405     local u={}; for k,v in pairs(t) do u[k] = l.copy(v) end
406     return u end
407
408 -- Return the 'p'-th thing from the sorted list 't'.
409 function l.per(t,p)
410     p=math.floor(((p or .5)*#t)+.5); return t[math.max(1,math.min(#t,p))] end
411
412 ----- Strings
413 -- 'o' generates a string from a nested table.
414 function l.o(t)
415     if type(t) ~= "table" then return tostring(t) end
416     local function show(k,v)
417         if not tostring(k):find"_" then
418             v = l.o(v)
419             return #t==0 and string.format("%s%s",k,v) or tostring(v) end end
420     local u={}; for k,v in pairs(t) do u[1+#u] = show(k,v) end
421     return (t._is or "").."["..table.concat(#t==0 and sort(u) or u, ",").."]" end
422
423 -- 'oo' prints the string from 'o'.
424 function l.oo(t) print(l.o(t)) return t end
425
426 -- Convert string to something else.
427 function l.coerce(s)
428     local function coerce1(s1)
429         if s1=="true" then return true end
430         if s1=="false" then return false end
431         return s1 end
432     return math.tointeger(s) or tonumber(s) or coerce1(s:match"^(%s*)(-)%s*$") end
433
434 -- Iterator over csv files. Call 'fun' for each record in 'fname'.
435 function l.csv(fname,fun)
436     local src = io.input(fname)
437     while true do
438         local s = io.read()
439         if not s then return io.close(src) else
440             local t={}
441             for s1 in s:gmatch("[^,]+") do t[1+#t] = l.coerce(s1) end
442             fun(t) end end end
443
444 ----- Settings
445 -- Parse help string looking for slot names and default values
446 function l.settings(s)
447     local t={}
448     s:gsub("%n[-][%S]+[%s]+[-][-]([%S+)%n")+="( [%S+)%",
449         function(k,x) t[k]=l.coerce(x) end)
450     t._help = s
451     return t end
452
453 -- Update 't' from values after command-line flags. Booleans need no values
454 -- (we just flip the defaults).
455 function l.cli(t)
456     for slot,v in pairs(t) do
457         v = tostring(v)
458         for n,x in ipairs(arg) do
459             if x=="-."(slot:sub(1,1)) or x=="-."slot then
460                 v = v=="false" and "true" or v=="true" and "false" or arg[n+1] end end
461             t[slot] = l.coerce(v) end
462     if t.help then os.exit(print("un"..t._help.."un")) end
463     return t end

```

0015 0 0000