



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



```
1  -- In this code:
2  -- Line strive to be 80 chars (or less)
3  -- Two spaces before function arguments denote optionals.
4  -- Four spaces before function arguments denote local variables.
5  -- Private functions start with '_'
6  -- Arguments of private functions do anything at all
7  -- Local variables inside functions do anything at all
8  -- Arguments of public functions use type hints
9  -- Variable 'x' is anything
10 -- Prefix 'is' is a boolean
11 -- Prefix 'fun' is a function
12 -- Prefix 'f' is a filename
13 -- Prefix 'n' is a string
14 -- Prefix 's' is a string
15 -- Prefix 'c' is a column index
16 -- 'col' denotes 'num' or 'sym'
17 -- 'x' is anything (table or number of boolean or string)
18 -- 'v' is a simple value (number or boolean or string)
19 -- Suffix 'a' is a list of things
20 -- Tables are 't' or, using the above, a table of numbers would be 'ns'
21 -- Type names are lower case versions of constructors; e.g 'col' isa 'Cols'.
22
23 -- All the demo functions 'eg.fun1', 'eg.fun2', etc can be called via
24 -- e.g. 'lua eg.lua -e fun1'.
25 local eg = {}
26
27 -----
28 local l=require"lib"
29 local _require="sam"
30
31 local o,oo,per,push,rnd = l.o,l.oo,l.per,l.push,l.rnd
32 local add,adds,dist,div = _add,_adds,_dist,_div
33 local mid,records,the = _mid,_records,_the
34 local Num,Sym = _Num,_Sym
35
36 -- Settings come from big string top of 'sam.lua'
37 -- (maybe updated from command line)
38 function eg.the(i) oo(the); return true end
39
40 -- The middle and diversity of a set of symbols is called "mode"
41 -- and "entropy" (and the latter is zero when all the symbols
42 -- are the same).
43 function eg.ent( sym,ent)
44 sym= adds(Sym(), { "a","a","a","a","b","b","c","c" })
45 ent= div(sym)
46 print(ent,mid(sym))
47 return 1.37 <= ent and ent <=1.38 end
48
49 -- The middle and diversity of a set of numbers is called "median"
50 -- and "standard deviation" (and the latter is zero when all the nums
51 -- are the same).
52 function eg.num( num)
53 num=Num()
54 for i=1,100 do add(num,i) end
55 local med,ent = mid(num), rnd(div(num),2)
56 print( med(num),rnd(div(num),2))
57 return 50<= med and med<= 52 and 30.5 <ent and ent <32 end
58
59 -- Nums store only a sample of the numbers added to it (and that storage
60 -- is done such that the kept numbers span the range of inputs).
61 function eg.bignum( num)
62 the.num= 32
63 for i=1,1000 do add(num,i) end
64 oo(_nums(num))
65 return 32==#num._has end
66
67 -- We can read data from disk-based csv files, where row1 lists a
68 -- set of columns names. These names are used to work out what are Nums, or
69 -- ro Syms, or goals to minimize/maximize, or (indeed) what columns to ignore.
70 function eg.records(i)
71 oo(records("../data/auto93.csv").cols.y); return true end
72
73 -- Any two rows have a distance 0.1 that satisfies equality, symmetry
74 -- and the triangle inequality.
75 function eg.dist( data,t)
76 data=records("../data/auto93.csv")
77 t={}
78 for i=1,100 do
79 local A,B,C = l.any(data.rows), l.any(data.rows), l.any(data.rows)
80 local a,b,c = dist(data,B,C), dist(data,A,C), dist(data,A,B)
81 assert(a<=1 and b<=1 and c<=1)
82 assert(a==0 and b==0 and c==0)
83 assert( dist(data,A,A) == 0)
84 assert( dist(data,A,B) == dist(data,B,A)) -- symmetry
85 assert(a+b>=c) -- triangle inequality
86 for _,x in pairs(a) do push(t,rnd(x,2)) end end
87 table.sort(t)
88 oo(t)
89 return true end
90
91 -----
92 the = l.cli(the)
93 os.exit( l.runs(the.eg, eg, the))
```



```
101 -- For a list of coding conventions in this file, see
102 -- [eg.lua](https://github.com/timm/lua/blob/main/src/sam/eg.lua).
103 local l=require"lib"
104 local the,settings={}
105 SAM : Semi-supervised And Multi-objective explanations
106 (c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license
107
108 USAGE: lua eg.lua [OPTIONS]
109
110 OPTIONS:
111 -e --eg start-up example = nothing
112 -h --help show help = false
113 -n --nums how many numbers to keep = 256
114 -p --p distance coefficient = 2
115 -s --seed random number seed = 10019]]
116 -- Commonly used lib functions.
117 local o,oo,per,push = l.o,l.oo,l.per,l.push
118
119 -----
120 local Data,Cols,Sym,Num,Row
121 -- Holder of 'rows' and their summaries (in 'cols').
122 function Data() return {cols=nil, -- summaries of data
123 rows={}, -- kept data
124 } end
125
126 -- Holder of summaries of columns.
127 -- Columns are created once, then shared across the following slots.
128 function Cols() return {
129 names={}, -- all column names
130 all={}, -- holds all the columns (including the skipped ones)
131 klass=nil, -- shares the symbolic klass column (if it exists)
132 x={}, -- shares the independent columns (that are not skipped)
133 y={} -- shared the dependet columns (that are not skipped)
134 } end
135
136 -- Summary of a stream of symbols.
137 function Sym(c,s)
138 return {n=0,at=c or 0, -- items seen
139 at=c or 0, -- column position
140 names=or "", -- column name
141 _has={} -- kept data
142 } end
143
144 -- Summary of a stream of numbers.
145 function Num(c,s)
146 return {n=0,at=c or 0, names=or "", _has={}, -- as per Sym
147 isNum=true, -- mark that this is a number
148 lo= math.huge, -- lowest seen
149 hi= -math.huge, -- highest seen
150 sorted=true, -- no updates since last sort of data
151 w={s or ""}.find"$" and -1 or 1 -- minimizing if w=-1
152 } end
153
154 -- Hold one record
155 function Row(t) return {cells=t, -- one record
156 cooked=nil -- used if we discretize data
157 } end
158
159 -----
160 local add,adds,clone,div,mid,norm,nums,record,records,stats
161 ----- Create
162 -- Generate rows from some 'src'. If 'src' is a string, read rows from file;
163 -- else read rows from a 'src' table. When reading, use row1 to define columns.
164 function records(src, data,head,body)
165 local cols = Cols()
166 cols.names = namesa
167 for c,s in pairs(sNames) do
168 local col = push(cols.all, -- Numerics start with Uppercase.
169 (s:find"^[A-Z]" and Num or Sym)(c,s))
170 if not s:find"$" then -- some columns are skipped
171 s:find"^[a-z]" and cols.x, col) -- some cols are goal cols
172 if s:find"$" then cols.klass=col end end end
173 return cols
174 end
175 function body(t) -- treat first row differently (defines the columns)
176 if data.cols then record(data,t) else data.cols=head(t) end
177 end
178 data = Data()
179 if type(src)=="string" then l.csv(src, body) else
180 for _,t in pairs(src or {}) do body(t) end end
181 return data end
182
183 -- Return a new data with same structure as 'data'. Optionally, oad in 'rows'.
184 function clone(data1, rows)
185 data2=Data()
186 data2.cols = _head(data1.cols.names)
187 for _,row in pairs(rows or {}) do record(data2,row) end
188 return data2 end
189
190 -----
191 -- Add one thing to 'col'. For Num, keep at most 'nums' items.
192 function add(col,v)
193 if v=="$" then
194 col.n = col.n + 1
195 if not col.isNum then col._has[v] = 1 + (col._has[v] or 0) else
196 col.lo = math.min(v, col.lo)
197 col.hi = math.max(v, col.hi)
198 local pos
199 if #col._has < the.nums then pos = 1 + (#col._has)
200 elseif math.random() < the.nums/col.n then pos = math.random(#col._has) end
201 if pos then col._has[pos] = false
202 col._has[pos] = tonumber(v) end end end end
203
204 -- Add many things to col
205 function adds(col,t) for _,v in pairs(t) do add(col,v) end; return col end
206
207 -- Add a 'row' to 'data'. Calls 'add()' to update the 'cols' with new values.
208 function record(data,xs)
209 local row= push(data.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
210 for _,todo in pairs(data.cols.x, data.cols.y) do
211 for _,col in pairs(todo) do
212 add(col, row.cells[col.at]) end end end
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
```

```
229 -----
230 -- Return kept numbers, sorted.
231 function nums(num)
232 if not num.sorted then table.sort(num._has); num.sorted=true end
233 return num._has end
234
235 -- Normalized numbers 0..1. Everything else normalizes to itself.
236 function norm(col,n)
237 return x=="$" or not col.isNum and x or (n-col.lo)/(col.hi-col.lo + 1E-32) end
238
239 -- Diversity (standard deviation for Nums, entropy for Syms)
240 function div(col)
241 if col.isNum then local a=nums(col); return (per(a,.9)-per(a,.1))/2.58 else
242 local function fun(p) return p*math.log(p,2) end
243 local e=0
244 for _,n in pairs(col._has) do if n>0 then e=fun(n/col.n) end end
245 return e end end
246
247 -- Central tendency (median for Nums, mode for Syms)
248 function mid(col)
249 if col.isNum then return per(nums(col),.5) else
250 local most,mode = -1
251 for k,v in pairs(col._has) do if v>most then mode,most=k,v end end
252 return mode end end
253
254 -- For 'showCols' (default='data.cols.x') in 'data', report 'fun' (default='mid').
255 function stats(data, showCols,fun, t)
256 showCols, fun = showCols or data.cols.y, fun or mid
257 t={}; for _,col in pairs(showCols) do t[col.name]=fun(col) end; return t end
258
259 -----
260 -- Distance functions
261 local dist
262 -- Distance between rows (returns 0..1). For unknown values, assume max distance.
263 function dist(data,t1,t2)
264 local function fun(col, v1,v2)
265 if v1=="$" and v2=="$" then return 1 end
266 if not col.isNum then return v1==v2 and 0 or 1 end
267 v1,v2 = norm(col,v1), norm(col,v2)
268 if v1=="$" then v1 = v2<.5 and 1 or 0 end
269 if v2=="$" then v2 = v1<.5 and 1 or 0 end
270 return math.abs(v1-v2)
271 end
272 local d = 0
273 for _,col in pairs(data.cols.x) do
274 d = d + fun(col, t1.cells[col.at], t2.cells[col.at])^the.p end
275 return (d/#data.cols.x)^(1/the.p) end
276
277 -----
278 -- That's all folks.
279 return {the=the,
280 Data=Data, Cols=Cols, Sym=Sym, Num=Num, Row=Row,
281 add=adds, adds=adds, clone=clone, dist=dist, div=div,
282 mid=mid, nums=nums, records=records, record=record, stats=stats}
```

[illegible]

```

402 else if [k==1] then -- run all
403   for _k in pairs(_egs(1)) do
404     fails=fails + (l.run(k,funs,settings) and 0 or 1) end
405   else funs[k] then -- run one
406     math.randomseed(settings.seed) -- reset seed
407     local b4={}; for k,v in pairs(settings) do b4[k]=v end
408     local out=funs[k]()
409     for k,v in pairs(b4) do settings[k]=v end -- restore old settings
410     print("****", k, out and "PASS" or "FAIL") end
411   l.roques(1)
412   return fails end
413
414 -----
415 -- That's all folks.
416 return l

```