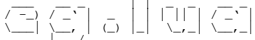




SAM : Semi-supervised And Multi-objective explanations
(c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license



```
1
2
3
4
5
6
7
8
9
10
11
12 SAM : Semi-supervised And Multi-objective explanations
13 (c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license
14
15
16
17
18 -- In this code:
19 -- Line strive to be 80 chars (or less)
20 -- Two spaces before function arguments denote optionals.
21 -- Four spaces before function arguments denote local variables.
22 -- Private functions start with '_'
23 -- Arguments of private functions do anything at all
24 -- Local variables inside functions do anything at all
25 -- Arguments of public functions use type hints
26 -- Variable 'x' is anything
27 -- Prefix 'is' is a boolean
28 -- Prefix 'fun' is a function
29 -- Prefix 'f' is a filename
30 -- Prefix 'n' is a string
31 -- Prefix 's' is a string
32 -- Prefix 'c' is a column index
33 -- 'col' denotes 'num' or 'sym'
34 -- 'x' is anything (table or number of boolean or string)
35 -- 'v' is a simple value (number or boolean or string)
36 -- Suffix 'a' is a list of things
37 -- Tables are 't' or, using the above, a table of numbers would be 'ns'
38 -- Type names are lower case versions of constructors; e.g 'col' isa 'Cols'.
39
40 -- All demo functions 'eg.fun1' can be called via 'lua eg.lua -e fun1'.
41 local eg = {}
42
43 local l=require"lib"
44 local _require="sam"
45 local O,oo,per,push,rnd = l.o,l.oo,l.per,l.push,l.rnd
46 local add,adds,dist,div = _,add,_,adds,_,dist,_,div
47 local mid, records, the = _,mid,_,records,_,the
48 local Num,Sym = _,Num, _,Sym
49
50 -- Settings come from big string top of "sam.lua"
51 -- (maybe updated from command line)
52 function eg.the(i) oo(the); return true end
53
54 -- The middle and diversity of a set of symbols is called "mode"
55 -- and "entropy" (and the latter is zero when all the symbols
56 -- are the same).
57 function eg.ent(i, sym,ent)
58   sym= adds(Sym(i), { "a","a","a","a","b","b","b","c","c" })
59   ent= div(sym)
60   print(ent,mid(sym))
61   return 1.37 <= ent and ent <= 1.38 end
62
63 -- The middle and diversity of a set of numbers is called "median"
64 -- and "standard deviation" (and the latter is zero when all the nums
65 -- are the same).
66 function eg.num(i num)
67   num=Num(i)
68   for i=1,100 do add(num,i) end
69   local med,ent = mid(num), rnd(div(num),2)
70   print(mid(num),rnd(div(num),2))
71   return 50<= med and med<= 52 and 30.5 <=ent and ent <= 32 end
72
73 -- Nums store only a sample of the numbers added to it (and that storage
74 -- is done such that the kept numbers span the range of inputs).
75 function eg.bignum(i num)
76   num=Num(i)
77   the.nums = 32
78   for i=1,1000 do add(num,i) end
79   oo(_nums(num))
80   return 32==#num._has end
81
82 -- We can read data from disk-based csv files, where row1 lists a
83 -- set of column names. These names are used to work out what are Nums, or
84 -- ro Syms, or goals to minimize/maximize, or (indeed) what columns to ignore.
85 function eg.records(i)
86   oo(records("./data/aut93.csv").cols.y); return true end
87
88 -- Any two rows have a distance 0..1 that satisfies equality, symmetry
89 -- and the triangle inequality.
90 function eg.dist(i data,t)
91   data=records("./data/aut93.csv")
92   t={}
93   for i=1,100 do
94     local A,B,C = l.any(data.rows), l.any(data.rows), l.any(data.rows)
95     local a,b,c = dist(data,B,C), dist(data,A,C), dist(data,A,B)
96     assert(a<=1 and b<=1 and c<=1)
97     assert(a>=0 and b>=0 and c>=0)
98     assert( dist(data,A,A) == 0) -- equality
99     assert( dist(data,A,B) == dist(data,B,A)) -- symmetry
100    assert(a+b>=c) -- triangle inequality
101    for _,x in pairs(a) do push(t,rnd(x,2)) end end
102    table.sort(t)
103    oo(t)
104    return true end
105
106 the = l.cli(the)
107 os.exit( l.runs(the.eg, eg, the))
```



```
108
109
110
111
112 -- For a list of coding conventions in this file, see
113 -- [eg.lua](https://github.com/timm/lua/blob/main/src/sam/eg.lua).
114 local l=require"lib"
115 local the,l.settings={
116   SAM : Semi-supervised And Multi-objective explanations
117   (c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license
118
119   USAGE: lua eg.lua [OPTIONS]
120
121   OPTIONS:
122   -e --eg start-up example = nothing
123   -c --cohen small effect = .35
124   -h --help show help = false
125   -m --min min size = n^(the.min) = .5
126   -n --nums how many numbers to keep = 256
127   -p --p distance coefficient = 2
128   -s --seed random number seed = 10019]]
129 -- Commonly used lib functions.
130 local O,oo,per,push = l.o,l.oo,l.per, l.push
131
132 ----- Classes
133 local Data,Cols,Sym,Num,Row
134 -- Holder of 'rows' and their summaries (in 'cols').
135 function Data(i) return {cols=nil, -- summaries of data
136   rows={}, -- kept data
137   } end
138
139 -- Holds of summaries of columns.
140 -- Columns are created once, then may appear in multiple slots.
141 function Cols(i) return {
142   names={}, -- all column names
143   all={}, -- all the columns (including the skipped ones)
144   klass=nil, -- the single dependent klass column (if it exists)
145   x={}, -- independent columns (that are not skipped)
146   y={} -- depedent columns (that are not skipped)
147   } end
148
149 -- Summarizes a stream of symbols.
150 function Sym(C,s)
151   return {n=0, -- items seen
152     at=c or 0, -- column position
153     names= C or "", -- column name
154     _has={} -- kept data
155   } end
156
157 -- Summarizes a stream of numbers.
158 function Num(C,s)
159   return {n=0,at=c or 0, names= C or "", _has={}, -- as per Sym
160     isNum=true, -- mark that this is a number
161     low= math.huge, -- lowest seen
162     his= -math.huge, -- highest seen
163     sorted=true, -- no updates since last sort of data
164     w=(s or ""):find"$" and -1 or 1 -- minimizing if w=-1
165   } end
166
167 -- Holds one record
168 function Row(t) return {cells=t, -- one record
169   cooked=i.copy(t) -- used if we discretize data
170   } end
171
172 ----- Data Functions
173 local add,adds,clone,div,mid,norm,nums,record,records,stats
174 ----- Create
175 -- Generate rows from some 'src'. If 'src' is a string, read rows from file;
176 -- else read rows from a 'src' table. When reading, use row1 to define columns.
177 function records(src, data,head,body)
178   function head(sNames)
179     local cols = Cols(i)
180     cols.names = sNames
181     for c,s in pairs(sNames) do
182       local col = push(cols.all, -- Numerics start with Uppercase.
183         s:find"^[A-Z]" and Num or Sym(c,s))
184       if not s:find"$" then -- some columns are skipped
185         push(cols.find["I-"] and cols.y or cols.x, col) -- some cols are goal cols
186         if s:find"$" then cols.klass=col end end end
187     return cols
188   end
189   function body(t) -- treat first row differently (defines the columns)
190     if data.cols then record(data,t) else data.cols=head(t) end
191   end
192   data = Data(i)
193   if type(src)=="string" then l.csv(src, body) else
194     for _,t in pairs(src or {}) do body(t) end end
195   return data end
196
197 -- Return a new data with same structure as 'data1'. Optionally, oad in 'rows'.
198 function clone(data1, rows)
199   data2=Data(i)
200   data2.cols = _head(data1.cols.names)
201   for _,row in pairs(rows or {}) do record(data2,row) end
202   return data2 end
203
204 ----- Update
205 -- Add one thing to 'col'. For Num, keep at most 'nums' items.
206 function add(col,v)
207   if v=="?" then
208     col.n = col.n + 1
209     if not col.isNum then col._has[v] = 1 + (col._has[v] or 0) else
210       col.lo = math.min(v, col.lo)
211       col.hi = math.max(v, col.hi)
212     local pos
213     if #col._has < the.nums then pos = 1 + (#col._has)
214     elseif math.random(i) < the.nums/col.n then pos = math.random(#col._has) end
215     if pos then col.sorted = false
216       col._has[pos] = tonumber(v) end end end end
217
218 -- Add many things to col
219 function adds(col,t) for _,v in pairs(t) do add(col,v) end; return col end
220
221 -- Add a 'row' to 'data'. Calls 'add()' to update the 'cols' with new values.
222 function record(data,xs)
223   local row= push(data.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
224   for _,todo in pairs(data.cols.x, data.cols.y) do
225     for _,col in pairs(todo) do
```

```
226     add(col, row.cells[col.at]) end end end
227
228 -- Unsupervised discretization.
229 -- ??? have i go the unknowns handles rite?
230 function unsuper(data)
231   local function cell(row,col) return row.cells[col.at] end
232   local function sorter(col) return
233     function (row1,row2)
234       local x,y = row1.cells[col.at], row2.cells[cols.at]
235       return (x=="?" and math.huge or x) < (y=="?" and math.huge or y) end end
236   for _,col in pairs(data.cols.x) do
237     if col.isNum then
238       local enough = (#data.rows)*the.min
239       local epsilon = div(col) *the.cohen
240       table.sort(data.rows, sorter(col))
241       for i,row in pairs(data.rows) do
242         v = cell(data.rows[i+1],col)
243         if v == "?" then
244           if not n then
245             local n,lo,hi = 0,cell(row,col), cells(row,col) end
246             if i < #data.rows - enough then
247               if v == w and (hi - lo) > epsilon and n > enough then
248                 n,lo,hi = 0,v,v end end end
249                 n = n+1
250                 hi = v end
251               row.cooked[col.at] = lo end end end
252
253 ----- Query
254 -- Return kept numbers, sorted.
255 function nums(num)
256   if not num.sorted then table.sort(num._has); num.sorted=true end
257   return num._has end
258
259 -- Normalized numbers 0..1. Everything else normalizes to itself.
260 function norm(col,n)
261   return x=="?" or not col.isNum and x or (n-col.lo)/(col.hi-col.lo + 1E-32) end
262
263 -- Diversity (standard deviation for Nums, entropy for Syms)
264 function div(col)
265   if col.isNum then local a=nums(col); return (per(a,.9)-per(a,.1))/2.58 else
266     local function fun(p) return p*math.log(p,2) end
267     local e=0
268     for _,n in pairs(col._has) do if n>0 then e=fun(n/col.n) end end
269     return e end end
270
271 -- Central tendency (median for Nums, mode for Syms)
272 function mid(col)
273   if col.isNum then return per(nums(col),.5) else
274     local most,mode = -1
275     for k,v in pairs(col._has) do if v>most then mode,most=k,v end end
276     return mode end end
277
278 -- For 'showCols' (default='data.cols.x') in 'data', report 'fun' (default='mid').
279 function stats(data, showCols,fun, t)
280   showCols, fun = showCols or data.cols.y, fun or mid
281   t={}; for _,col in pairs(showCols) do t[col.name]=fun(col) end; return t end
282
283 ----- Distance functions
284 local dist
285 -- Distance between rows (returns 0..1). For unknown values, assume max distance.
286 function dist(data,t1,t2)
287   local function fun(col, v1,v2)
288     if v1=="?" and v2=="?" then return 1 end
289     if not col.isNum then return v1==v2 and 0 or 1 end
290     v1,v2 = norm(col,v1), norm(col,v2)
291     if v1=="?" then v1 = v2<.5 and 1 or 0 end
292     if v2=="?" then v2 = v1<.5 and 1 or 0 end
293     return math.abs(v1-v2)
294   end
295   local d = 0
296   for _,col in pairs(data.cols.x) do
297     d = d + fun(col, t1.cells[col.at], t2.cells[col.at])*the.p end
298   return (d/#data.cols.x)^(1/the.p) end
299
300 ----- That's all folks.
301 return {the=the,
302   Data=Data, Cols=Cols, Sym=Sym, Num=Num, Row=Row,
303   add=add, adds=adds, clone=clone, dist=dist, div=div,
304   mid=mid, nums=nums, records=records, record=record, stats=stats}
```

[illegible]

```

self.funs["all"] = then -- run all
for _k in pairs(_egs) do
  fails=fails + (1,run(k,funs,settings) and 0 or 1) end
return funs(k) then -- run one
math.randomseed(settings.seed) -- reset seed
local b4={}; for k,v in pairs(settings) do b4[k]=v end
local oldfuns={k()
for k,v in pairs(b4) do settings[k]=v end -- restore old settings
print("*****", k, out and "PASS" or "FAIL") end
--roques()
return fails end
-----
That's all folks.
return 1

```