

```

1
2
3
4
5
6
7
8
9
10 SAM : Semi-supervised And Multi-objective explanations
11 (c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license
12
13
14
15
16 -- In this code:
17 -- Line strive to be 80 chars (or less)
18 -- Two spaces before function arguments denote optionals.
19 -- Four spaces before function arguments denote local variables.
20 -- Private functions start with '_'
21 -- Arguments of private functions 'do anything at all'
22 -- Local variables inside functions do anything at all
23 -- Arguments of public functions use type hints
24 -- Variable 'x' is anything
25 -- Prefix 'is' is a boolean
26 -- Prefix 'fun' is a function
27 -- Prefix 'f' is a filename
28 -- Prefix 'n' is a string
29 -- Prefix 's' is a string
30 -- Prefix 'c' is a column index
31 -- 'col' denotes 'num' or 'sym'
32 -- 'x' is anything (table or number of boolean or string
33 -- 'v' is a simple value (number or boolean or string)
34 -- Suffix 's' is a list of things
35 -- Tables are 't' or, using the above, a table of numbers would be 'ns'
36 -- Type names are lower case versions of constructors; e.g 'col' is a 'Cols'.
37
38 -- All demo functions 'eg.fun1' can be called via 'lua eg.lua -e fun1'.
39
40 local eg = {}
41
42 local l=require"lib"
43 local _require="sam"
44 local o,oo,per,push,rnd = 1.0,1.00,1.per,1.push,1.rnd
45 local add,adds,dist,div = _add,_adds,_dist,_div
46 local mid, records, the = _mid,_records,_the
47 local Num,Sym = _Num, _Sym
48
49 -- Settings come from big string top of "sam.lua"
50 -- (maybe updated from command line)
51 function eg.the() oo(the); return true end
52
53 -- The middle and diversity of a set of symbols is called "mode"
54 -- and "entropy" (and the latter is zero when all the symbols
55 -- are the same).
56 function eg.ent( sym,ent)
57   sym= adds(Sym(), { "a","a","a","a","b","b","b","c" })
58   ent= div(sym)
59   print(ent,mid(sym))
60   return 1.37 <= ent and ent <= 1.38 end
61
62 -- The middle and diversity of a set of numbers is called "median"
63 -- and "standard deviation" (and the latter is zero when all the nums
64 -- are the same).
65 function eg.num( num)
66   num=Num()
67   for i=1,100 do add(num,i) end
68   local med,ent = mid(num), rnd(div(num),2)
69   print(mid(num),rnd(div(num),2))
70   return 50<= med and med<= 52 and 30.5 <ent and ent <32 end
71
72 -- Nums store only a sample of the numbers added to it (and that storage
73 -- is done such that the kept numbers span the range of inputs).
74 function eg.bignum( num)
75   num=Num()
76   the.nums = 32
77   for i=1,1000 do add(num,i) end
78   oo(_nums(num))
79   return 32==#num._has end
80
81 -- We can read data from disk-based csv files, where row1 lists a
82 -- set of columns names. These names are used to work out what are Nums, or
83 -- ro Syms, or goals to minimize/maximize, or (indeed) what columns to ignore.
84 function eg.records()
85   oo(records("../data/auto93.csv").cols.y); return true end
86
87 -- Any two rows have a distance 0..1 that satisfies equality, symmetry
88 -- and the triangle inequality.
89 function eg.dist( data,t)
90   data=records("../data/auto93.csv")
91   t={}
92   for i=1,100 do
93     local A,B,C = 1.any(data.rows), 1.any(data.rows), 1.any(data.rows)
94     local a,b,c = dist(data,A,B,C), dist(data,A,C), dist(data,A,B)
95     assert(a<=1 and b<=1 and c<=1)
96     assert(a>=0 and b>=0 and c>=0)
97     assert( dist(data,A,A) == 0) -- equality
98     assert( dist(data,A,B) == dist(data,B,A)) -- symmetry
99     assert(a+b>=c) -- triangle inequality
100    for _,x in pairs(a) do push(t,rnd(x,2)) end end
101    table.sort(t)
102    oo(t)
103    return true end
104
105 function eg.far( data)
106   data = records("../data/auto93.csv")
107   oo(data.rows[1].cells)
108   for i,t in pairs(_around(data,data.rows[1])) do
109     if i>390 or i< 10 then print(o(t.row.cells),t.dist) end end
110   oo(_far(data, data.rows[1]).cells)
111   return true end
112
113 function eg.half( data)
114   data = records("../data/auto93.csv")
115   _halves(data)
116   _tree(_halves(data),function(t) tostring(10) end)
117   return true end
118
119 -----

```

08/25/22

```

122
123
124
125
126
127 -- For a list of coding conventions in this file, see
128 -- [eg.lua] (https://github.com/timm/lua/blob/main/src/sam/eg.lua).
129 local l=require"lib"
130 local the=_settings{[["
131 SAM : Semi-supervised And Multi-objective explanations
132 (c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license
133
134 USAGE: lua eg.lua [OPTIONS]
135
136 OPTIONS:
137 -b --bins      number of bins      = 8
138 -c --cohen     small effect        = .35
139 -e --eg        start-up example    = nothing
140 -F --far       far away            = .95
141 -f --file      file with csv data  = ../docs/auto93.csv
142 -h --help      show help           = false
143 -m --min       min size = n^(the.min) = 5
144 -n --nums      how many numbers to keep = 256
145 -p --p         distance coefficient = 2
146 -s --seed      random number seed  = 10019
147 -S --sample    how many rows to search = 512]])
148 -- Commonly used lib functions.
149 local lt,o,oo,map = 1.lt,1.0,1.oo,1.map
150 local per,push,sort = 1.per, 1.push,1.sort
151
152 ----- Classes
153 local Data,Cols,Sym,Num,Row
154 -- Holder of 'rows' and their summaries (in 'cols').
155 function Data() return { _is = "Data",
156   cols= nil, -- summaries of data
157   rows= {}, -- kept data
158 } end
159
160 -- Holds of summaries of columns.
161 -- Columns are created once, then may appear in multiple slots.
162 function Cols() return {
163   _is = "Cols",
164   names={}, -- all column names
165   all={}, -- all the columns (including the skipped ones)
166   klass=nil, -- the single dependent klass column (if it exists)
167   x={}, -- independent columns (that are not skipped)
168   y={} -- dependent columns (that are not skipped)
169 } end
170
171 -- Summarizes a stream of symbols.
172 function Sym(c,s)
173   return { _is="Sym",
174     n=0, -- items seen
175     at=c or 0, -- column position
176     name=s or "", -- column name
177     _has={} -- kept data
178   } end
179
180 -- Summarizes a stream of numbers.
181 function Num(c,s)
182   return { _is="Nums",
183     n=0,at=c or 0, name=s or "", _has={}, -- as per Sym
184     isNum=true, -- mark that this is a number
185     lo= math.huge, -- lowest seen
186     hi= -math.huge, -- highest seen
187     isSorted=true, -- no updates since last sort of data
188     w=(s or ""):find"$" and -1 or 1 -- minimizing if w=-1
189   } end
190
191 -- Holds one record
192 function Row(t) return { _is="Row",
193   cells=t, -- one record
194   cooked=l.copy(t), -- used if we discretize data
195   isEval=false -- true if y-values evaluated.
196 } end
197
198 ----- Data Functions
199 local add,adds,clone,div,mid,norm,nums,record,records,stats
200 ----- Create
201 -- Generate rows from some 'src'. If 'src' is a string, read rows from file;
202 -- else read rows from a 'src' table. When reading, use row1 to define columns.
203 function records(src, data,head,body)
204   function head(sNames)
205     local cols = Cols()
206     cols.names = sNames
207     for c,s in pairs(sNames) do
208       local col = push(cols.all, -- Numerics start with Uppercase.
209         (s:find"^[A-Z]" and Num or Sym)(c,s))
210       if not s:find"$" then -- some columns are skipped
211         push(c:find"^[!-]" and cols.y or cols.x,col) -- some cols are goal cols
212       if s:find"$" then cols.klass=col end end end
213     return cols
214   end
215   function body(t) -- treat first row differently (defines the columns)
216     if data.cols then record(data,t) else data.cols=head(t) end
217   end
218   data = Data()
219   if type(src)=="string" then l.csv(src, body) else
220     for _,t in pairs(src or {}) do body(t) end end
221   return data end
222
223 -- Return a new data with same structure as 'data1'. Optionally, oad in 'rows'.
224 function clone(data1, rows)
225   data2=Data()
226   data2.cols = _head(data1.cols.names)
227   for _,row in pairs(rows or {}) do record(data2,row) end
228   return data2 end
229
230 ----- Update
231 -- Add one thing to 'col'. For Num, keep at most 'nums' items.
232 function add(col,v)
233   if v==" " then
234     col.n = col.n + 1
235     if not col.isNum then col._has[v] = 1 + (col._has[v] or 0) else
236       col.lo = math.min(v, col.lo)
237       col.hi = math.max(v, col.hi)
238       local pos
239       if #col._has < the.nums then pos = 1 + (#col._has)
240       elseif math.random() < the.nums/col.n then pos = math.random(#col._has) end

```

Page 3/6

