



```
-- For a list of coding conventions in this file, see
-- [eg.lua] (https://github.com/timm/lua/blob/main/src/sam/eg.lua).
local require = lib"
local the = l.settings({{
SAM : Semi-supervised And Multi-objective explanations
(c) 2022 Tim Menzies <timmm@ieee.org> BSD-2 license
```

```

OPTIONS:
-b --bins      number of bins          = 8
-c --cohen     small effect            = .35
-e --eg        start-up example        = nothing
-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      = 100191))
-- Commonly used lib functions

local 0.00,per,push = 1.0,1.00,1.per, 1.push

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

-- Holds of summaries of columns.
-- Columns are created here, then may appear in multiple slots.
function Cols() return {
  names={}, -- all column names
  skip={}, -- 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)
}

```

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

```

-- Summarizes a stream of numbers.
function Num(c,s)
  return (n=0,at=c or 0, name=s or "", _has={}, -- as per Sym
    isNum=true, -- mark that this is a number
    lo= math.huge, -- lowest seen
    hi= math.huge, -- highest seen
    sorted=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 {cellar,t} -- one record
                                -- 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)
  function head (aNames)

```

```
cols.names = names
for c,i in pairs(nNames) do
    local col = push(cols.all, -- Numerics start with Uppercase.
                    (s:find("[A-Z]" == Num or Symb)(c,i))
    if not s:find("s") then -- some columns are skipped
        push(s:find("[+]" == col.s or col.x, col) -- some cols are goal cols
            if s:find("s") then cols.klass[col.i] end end
        return cols
    end
end
-----
function body(t) -- treat first row differently (defines the columns)
    if data.cols == 1 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 - pos) then pos = math.random(#col._has) + pos
    else pos = #col._has
    col._has[pos] = tonumber(v) end end end end
```

```
-- 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.cells and xs or Row(xs)) -- ensure xs is a Row
  for _,todo in pairs(data.cols.x, data.cols.y) do
```

```

219 for _,col in pairs(todo) do
220     add(col, row.cells[col.at]) end end end
221
222 ----- Query
223 -- Return kept numbers, sorted.
224 function nums(num)
225     if not num.sorted then table.sort(num._has); num.sorted=true end
226     return num._has end
227
228 -- Normalized numbers 0..1. Everything else normalizes to itself.
229 function norm(col,n)
230     return x=="?" or not col.isNum and x or (n-col.lo)/(col.hi-col.lo + 1E-32) end
231
232 ----- Diversity (standard deviation for Nums, entropy for Syms)
233 function div(col)
234     if col.isNum then local a=nums(col); return (per(a,.1)-per(a,.1))/2.58 else
235         local function fun(p) return p*math.log(p,2) end
236         local e=0
237         for _,n in pairs(col._has) do if n>0 then e=e+fun(n/(col.n)) end end
238         return e end
239
240 -- Central tendency (median for Nums, mode for Syms)
241 function mid(col)
242     if col.isNum then return per(nums(col),.5) else
243         local most,mode = -1
244         for k,v in pairs(col._has) do if v>most then mode,most=k,v end end
245         return mode end end
246
247 -- For 'showCols' (default='data.cols.x') in 'data', report 'fun' (default='mid')
248 function stats(data, showCols,fun, t)
249     showCols, fun = showCols or data.cols.y, fun or mid
250     t={}; for _,col in pairs(showCols) do t[col.name]=fun(col) end; return t end
251
252 ----- Discretization
253 -- Find ranges within a num (unsupervised).
254 function bins(num)
255     local a, epsilon = nums(num), the.cohen*div(num)
256     local enought = #a^the.min
257     local one = {lo=a[1], hi=a[1], n=0}
258     local t = {one}
259     for i,x in pairs(a) do
260         if i < a-enought and x ~= a[i+1] and n > enought and hi-lo > epsilon then
261             one = push(t, {lo=one.hi, hi=a[i], n=0}) end
262         one.hi = a[i]
263         one.n = 1 + one.n end
264     t[1].lo = -math.huge
265     t[#t].ho = math.huge
266     return t end
267
268 -- Fill in discretized values (in 'cooked').
269 function cook(data)
270     for _,num in pairs(data.cols.x) do
271         if num.isNum then local t = bins(num)
272             for _,row in pairs(data.rows) do
273                 local v = row.cells[num.at]
274                 if v ~= "?" then
275                     for _,bin in pairs(t) do
276                         if v > bin.lo and v <= bin.hi then
277                             row.cooked[col.at] = bin.lo
278                             break end end end end end end
279
280 -- Sum the entropy of the cooked independent columns.
281 function divs(data,rows)
282     local n = 0
283     for _,col in pairs(data.cols.x) do
284         local sym = Sym()
285         for _,row in pairs(rows or data.rows) do
286             v = row.cooked[col.at]
287             if v ~= "?" then add(s, v) end
288         n = n + div(sym) end
289     return n end
290
291 ----- Distance functions
292 local dist
293 -- Distance between rows (returns 0..1). For unknown values, assume max distance.
294 function dist(data,t1,t2)
295     local function fun(col, v1,v2)
296         if v1=="?" and v2=="?" then return 1 end
297         if not col.isNum then return v1==v2 and 0 or 1 end
298         v1,v2 = norm(col,v1), norm(col,v2)
299         if v1=="?" then v1 = v2<.5 and 1 or 0 end
300         if v2=="?" then v2 = v1<.5 and 1 or 0 end
301         return math.abs(v1-v2)
302     end
303     local d = 0
304     for _,col in pairs(data.cols.x) do
305         d = d + fun(col, t1.cells[col.at], t2.cells[col.at])^the.p end
306     return (d/#data.cols.x)^(1/the.p) end
307
308 -----
309 -- That's all folks.
310 return {thestats,
311         Data=Data, Cols=Cols, Sym=Sym, Num=Num, Row=Row,
312         add=add, addc=addc, clone=clone, dist=dist, div=div,
313         mid=mid, nums=nums, records=records, record=record, stats=stats}

```

[illegible]