

```

1  -- vim : ft=lua et sts=2 sw=2 ts=2 :
2  local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end --used later (to find rogues)
3  local help = {}
4
5  s1 == S.U.B.L.I.M.E. == Sublime's unsupervised
6  bifurcation: let's infer minimal explanations.
7  (c) 2022, Tim Menzies, BSD 2-clause license.
8
9  USAGE:
10     lua sl.lua [OPTIONS]
11
12  OPTIONS:
13     -Dump          stack dump on assert fails = false
14     -data N        data file = etc/data/auto93.csv
15     -enough F      recurse until rows'enough = .5
16     -far F         far = .9
17     -keep P        max kept items = 512
18     -p P          distance coefficient = 2
19     -seed P        set seed = 10019
20     -todo S        start up action (or 'all') = nothing
21     -help          show help = false
22
23  KEY: N=fileName F=float P=posint S=string
24  ]]
25  local any, asserts, big, cli, csv, fails, firsts, fmt, goalp, ignorep, klassp
26  local lessp, map, main, many, max, merge, min, morep, new, nump, o, oo, per, pop, push
27  local r, rows, slots, sort, sum, thing, things, unpack
28  local CLUSTER, COLS, EGS, NUM, ROWS = {}, {}, {}, {}, {}
29  local SKIP, SOME, SPAN, SYM = {}, {}, {}, {}
30
31  local the={}
32  help:gsub("\n[-](^%s+)[^n]*%s(^%s+)", function(key, x)
33      for n, flag in ipairs(arg) do
34          if flag:sub(1,1)=="-" and key:find("^"..flag:sub(2)..".*") then
35              x = x=="false" and true or arg[n+1] end end
36          if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
37              the[key] = tonumber(x) or x end )
38
39  -----
40  -- this code reads csv files where the words on line1 define column types.
41  function ignorep(x) return x:find"$" end -- columns to ignore
42  function klassp(x) return x:find"$" end -- symbolic goals to achieve
43  function lessp(x) return x:find"$" end -- number goals to minimize
44  function morep(x) return x:find"$" end -- numeric goals to maximize
45  function nump(x) return x:find"^[A-Z]" end -- numeric columns
46  function goalp(x) return morep(x) or lessp(x) or klassp(x) end
47
48  -- strings
49  fmt = string.format
50
51  -- maths
52  big = math.huge
53  max = math.max
54  min = math.min
55  r = math.random
56
57  -- tables
58  pop = table.remove
59  unpack = table.unpack
60  function any(t) return t[r(#t)] end
61  function firsts(a,b) return a[1] < b[1] end
62  function many(t,n, u) u={}; for i=1,n do push(u,any(t)) end; return u end
63  function per(t,p) return t[ (#t*(p or .5))/1 ] end
64  function push(t,x) table.insert(t,x); return x end
65  function sort(t,f) table.sort(t,f); return t end
66
67  -- meta
68  function map(t,f, u) u={}; for k,v in pairs(t) do push(u,f(v)) end; return u end
69  function sum(t,f, n) n=0; for _,v in pairs(t) do n=n+f(v) end; return n end
70  function slots(t, u)
71      u={}
72      for k,v in pairs(t) do k=tostring(k); if k:sub(1,1)=="_" then push(u,k) end end
73      return sort(u) end
74
75  -- print tables, recursively
76  function oo(t) print(o(t)) end
77  function o(t)
78      if type(t)=="table" then return tostring(t) end
79      local key=function(k) return fmt("%s %s", k, o(t[k])) end
80      local u = #t>0 and map(t,o) or map(slots(t),key)
81      return '{ '..table.concat(u, ", ")..' }' end
82
83  -- strings to things
84  function csv(file, x)
85      file = io.input(file)
86      return function()
87          x=io.read(); if x then return things(x) else io.close(file) end end end
88
89  function thing(x)
90      x = x:match"^(%s*)(-)%s*$"
91      if x=="true" then return true elseif x=="false" then return false end
92      return tonumber(x) or x end
93
94  function things(x,sep, t)
95      t={}
96      for y in x:gmatch(sep or "([.]+)") do push(t, thing(y)) end
97      return t end
98
99  CLASSES
100
101  function new(k,t) k.__index=k; k.__tostring=o; return setmetatable(t,k) end
102
103  -- COLS: turns list of column names into NUMs, SYMs, or SKIPs
104  function COLS.new(k,row, i)
105      i = new(k, {all={}, x={}, y={}, names=row})
106      for at,txt in ipairs(row) do push(i.all, i:col(at,txt)) end
107      return i end
108
109  function COLS.add(i,t)
110      for _,col in pairs(i.all) do col:add( t[col.at] ) end
111      return t end
112
113  function COLS.col(i,at,txt, col)
114      if ignorep(txt) then return SKIP:new(at,txt) end
115      col = (nump(txt) and NUM or SYM):new(at,txt)
116      push(goalp(txt) and i.y or i.x, col)
117      if klassp(txt) then i.klass = col end
118      return col end
119
120  -- NUM: summarizes a stream of numbers
121  function NUM.new(k,n,s)
122      return new(k, {n=0, at=n or 0, txt=s or "", has=SOME:new(i), ok=false,
123          w=lessp(s or "") and -1 or 1, lo=big, hi=-big}) end
124
125  function NUM.add(i,x)
126      if x == "?" then
127          i.n = i.n + 1
128          if i.has:add(x) then i.ok=false end
129          i.lo, i.hi = min(x, i.lo), max(x, i.hi); end end
130
131  function NUM.dist(i,x,y)
132      if x=="?" and y=="?" then return 1
133      elseif x=="?" then y=i:norm(y); x=y<0.5 and 1 or 0
134      elseif y=="?" then x=i:norm(x); y=x<0.5 and 1 or 0
135      else x,y = i:norm(x), i:norm(y) end
136      return math.abs(x-y) end
137
138  function NUM.mid(i) return per(i:sorted(), .5) end
139
140  function NUM.norm(i,x)
141      return math.abs(i.hi-i.lo)<1E-9 and 0 or (x-i.lo)/(i.hi - i.lo) end
142
143  function NUM.sorted(i)
144      if i.ok==false then table.sort(i.has.all); i.ok=true end
145      return i.has.all end
146
147  -- ROWS: manages 'rows', summarized in 'cols' (columns).
148  function ROWS.new(k, inits, i)
149      i = new(k, {rows={}, cols=nil})
150      if type(inits)=="string" then for t in csv(inits) do i:add(t) end end
151      if type(inits)=="table" then for t in inits do i:add(t) end end
152      return i end
153
154  function ROWS.add(i,t)
155      if i.cols then push(i.rows, i.cols:add(t)) else i.cols=COLS:new(t) end end
156
157  function ROWS.clone(i, j) j = ROWS:new(i); j:add(i.cols.names); return j end
158
159  function ROWS.dist(i, row1, row2, d, fun)
160      function fun(col) return col:dist(row1[col.at], row2[col.at])^the.p end
161      return (sum(i.cols.x, fun) / #i.cols.x)^(1/the.p) end
162
163  function ROWS.far(i, row1, rows, fun)
164      function fun(row2) return i:dist(row1, row2), row2 end
165      return unpack(per(sort(map(rows, fun), firsts), the.far)) end
166
167  function ROWS.half(i, top)
168      local some, top, c, x, y, tmp, mid, lefts, rights, _
169      some = many(i.rows, the.keep)
170      top = top or 1
171      _x = top:far(any(some), some)
172      c, y = top:far(x, some)
173      tmp = sort(map(i.rows, function(r) return top:project(r, x, y, c) end), firsts)
174      mid = #i.rows/2
175      lefts, rights = i:clone(), i:clone()
176      for at, row in pairs(tmp) do (at <= mid and lefts or rights):add(row[2]) end
177      return lefts, rights, x, y, c, tmp[mid] end
178
179  function ROWS.mid(i, cols)
180      return map(cols or i.cols.all, function(col) return col:mid() end) end
181
182  function ROWS.project(i, r, x, y, c, a,b)
183      a,b = i:dist(r,x), i:dist(r,y); return {(a^2 + c^2 - b^2)/(2*c), r} end
184
185  -- SKIP: summarizes things we want to ignore (so does nothing)
186  function SKIP.new(k,n,s) return new(k, {n=0, at=at or 0, txt=s or ""}) end
187  function SKIP.add(i,x) return x end
188  function SKIP.mid(i) return "?" end
189
190  -- SOME: keeps a random sample on the arriving data
191  function SOME.new(k, keep) return new(k, {n=0, all={}, keep=keep or the.keep}) end
192  function SOME.add(i,x)
193      i.n = i.n + 1
194      if #i.all < i.keep then push(i.all, x) ; return i.all
195      elseif r() < i.keep/i.n then i.all[r(#i.all)]=x; return i.all end end
196
197  -- SYM: summarizes a stream of symbols
198  function SYM.new(k,n,s)
199      return new(k, {n=0, at=n or 0, txt=s or "", has={}, most=0}) end
200
201  function SYM.add(i,x,inc)
202      if x == "?" then
203          inc = inc or 1
204          i.n = i.n + inc
205          i.has[x] = inc + (i.has[x] or 0)
206          if i.has[x] > i.most then i.most, i.mode=i.has[x], x end end
207
208  function SYM.dist(i,x,y) return (x=="?" and y=="?" and 1) or (x==y and 0 or 1) end
209  function SYM.mid(i) return i.mode end
210  function SYM.div(i, p)
211      return sum(i.has, function(k) p=-i.has[k]/i.n; return -p*math.log(p,2) end) end
212
213  function SYM.merge(i,j, k)
214      k = SYM:new(i.at, i.txt)
215      for x,n in pairs(i.has) do k:add(x,n) end
216      for x,n in pairs(j.has) do k:add(x,n) end
217      ei, ej, ek = i:div(), j:div(), k:div()
218      if i.n==0 or j.n==0 or .99*ek <= (i.n*ei + j.n*ej)/k.n then
219          return k end end
220

```

```

221 -- CLUSTER
222 --
223 --
224 --
225 -- CLUSTER: recursively divides data by clustering towards two distant points
226 function CLUSTER.new(k,sample,top)
227   local i,enough,left,right
228   top = top or sample
229   i = new(k, {here=sample})
230   enough = (#top.rows)^the.enough
231   if #sample.rows >= 2*enough then
232     left, right, i.x, i.y, i.c, i.mid = sample:half(top)
233     if #left.rows < #sample.rows then
234       i.left = CLUSTER:new(left, top)
235       i.right = CLUSTER:new(right, top) end end
236   return i end
237
238 function CLUSTER.show(i,pre, here)
239   pre = pre or ""
240   here=""
241   if not i.left and not i.right then here= o(i.here:mid(i.here.cols.y)) end
242   print(fmt("%6s:%-30s%6s",#i.here.rows, pre, here))
243   for _,kid in pairs(i.left, i.right) do
244     if kid then kid:show(pre .. "|. ") end end end
245
246 -- EXPLAIN
247 --
248 --
249 -- SPAN: keeps a random sample on the arriving data
250 function SPAN.new(k, col, lo, hi, has)
251   return new(k, {col=col, lo=lo, hi=hi or lo, has=has or SYM:new()}) end
252
253
254 function SPAN.add(i,x,y,n) i.lo, i.hi=min(x,i.lo),max(x,i.hi); i.has=add(y,n) end
255 function SPAN.merge(i,j)
256   local has = i.has:merge(j.has)
257   if now then return SPAN:new(i.col, i.lo, j.hi, has) end end
258
259 function SPAN.select(i,row, x)
260   x = row[i.col.at]
261   return (x=="?") or (i.lo==i.hi and x==i.lo) or (i.lo <= x and x < i.hi) end
262
263 -- EXPLAIN
264 function EXPLAIN(k, sample,top)
265   i.here = sample
266   top = top or sample
267   enough = (#top.rows)^the.enough
268   if #top.rows >= 2*enough then
269     left, right = sample:half(top)
270     spans = {}
271     for n,col in pairs(i.cols.x) do
272       tmp = col:spans(j.cols.x[n])
273       if #tmp>1 then for _,one in pairs(tmp) do push(spans,one) end end
274       if #spans > 2 then
275         XXXX?
276       end
277     end
278     function SYM.spans(i, j)
279       local xys,all,one,last,xys,x,c,n = {},{}
280       for x,n in pairs(i.has) do push(xys, {x,"this",n}) end
281       for x,n in pairs(j.has) do push(xys, {x,"that",n}) end
282       for _,tmp in pairs(sort(xys,firsts)) do
283         x,c,n = unpack(tmp)
284         if x ~= last then
285           last = x
286           one = push(all, Span(i,x,x)) end
287         one:add(x,y,n) end
288       return all end
289     function NUM.spans(i, j)
290       local xys,all,lo,hi,gap,xys,one,x,c,n = {},{}
291       lo,hi = min(i.lo, j.lo), max(i.hi,j.hi)
292       gap = (hi - lo) / bins
293       for x,n in pairs(i.has) do push(xys, {x,"this",1}) end
294       for x,n in pairs(j.has) do push(xys, {x,"that",1}) end
295       one = Span:new(i,lo,lo)
296       all = {one}
297       for _,tmp in pairs(sort(xys,firsts)) do
298         x,c,n = unpack(tmp)
299         if one.hi - one.lo > gap then one = push(all, Span(i, one.hi, x)) end
300         one:add(x,y) end
301       all = merge(all)
302       all[1].lo = -big
303       all[#all].hi = big
304       return all end
305
306 function merge(b4, j,n,now,a,b,merged)
307   j,n,now = 0,#b4,{}
308   while j < #b4 do
309     j = j+1
310     a, b = b4[j], b4[j+1]
311     if b then
312       merged = a:merge(b)
313       if merged then a,j = merged, j+1 end end
314     push(now,a)
315     j = j+1 end
316   return #now == #b4 and b4 or merge(now) end

```

```

317 -- DEMOS
318 --
319 --
320 --
321 fails=0
322 function asserts(test, msg)
323   print(test and "PASS: "or "FAIL: ",msg or "")
324   if not test then
325     fails=fails+1
326     if the.dump then assert(test,msg) end end end
327
328 function EGS.nothing() return true end
329 function EGS.the() co(the) end
330 function EGS.rand() print(r()) end
331 function EGS.some(s,t)
332   s=SOME:new(100)
333   for i=1,100000 do s:add(i) end
334   for j,x in pairs(sort(s.all)) do
335     --if (j % 10)==0 then print("") end
336     --io.write(fmt("%6s",x)) end end
337     fmt ("%6s",x) end end
338
339 function EGS.clone( r,s)
340   r = ROWS:new(the.data)
341   s = r:clone()
342   for _,row in pairs(r.rows) do s:add(row) end
343   asserts(r.cols.x[1].lo==s.cols.x[1].lo,"clone.lo")
344   asserts(r.cols.x[1].hi==s.cols.x[1].hi,"clone.hi")
345   end
346
347 function EGS.data( r)
348   r = ROWS:new(the.data)
349   asserts(r.cols.x[1].hi == 8, "data.columns") end
350
351 function EGS.dist( r,rows,n)
352   r = ROWS:new(the.data)
353   rows = r.rows
354   n = NUM:new()
355   for _,row in pairs(rows) do n:add(r:dist(row, rows[1])) end
356   --oo(r.cols.x[2]:sorted()) end
357   o(r.cols.x[2]:sorted()) end
358
359 function EGS.many( t)
360   t={} ; for j=1,100 do push(t,j) end
361   --print(oo(many(t, 10))) end
362   o(many(t, 10)) end
363
364 function EGS.far( r,c,row1,row2)
365   r = ROWS:new(the.data)
366   row1 = r.rows[1]
367   c,row2 = r:far(r.rows[1], r.rows) end
368   --print(c,"\n",o(row1),"\n", o(row2)) end
369
370 function EGS.half( r,c,row1,row2)
371   local lefts,rights,x,y,x
372   r = ROWS:new(the.data)
373   r:mid(r.cols.y)
374   lefts,rights,x,y,c = r:half()
375   lefts:mid(lefts.cols.y)
376   rights:mid(rights.cols.y)
377   asserts(true,"half") end
378
379 function EGS.cluster(r)
380   r = ROWS:new(the.data)
381   --CLUSTER:new(r):show() end
382   CLUSTER:new(r) end
383
384 -- start-up
385 if arg[0] == "slua" then
386   oo(the)
387   if the.help then print(help) else
388     local b4={} ; for k,v in pairs(the) do b4[k]=v end
389     for _,todo in pairs(the.todo=="all" and slots(EGS) or (the.todo)) do
390       for k,v in pairs(b4) do the[k]=v end
391       math.randomseed(the.seed)
392       if type(EGS[todo])=="function" then EGS[todo]() end end
393     end
394     for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
395     os.exit(fails)
396   else
397     return {CLUSTER=CLUSTER, COLS=COLS, NUM=NUM, ROWS=ROWS,
398           SKIP=SKIP, SOME=SOME, SYM=SYM,the=the,oo=oo,o=o}
399   end
400   -- git rid of SOME for rows
401   -- nss = NUM | SYM | SKIP
402   -- COLS = all:[nss] +, x:[nss]*, y:[nss]*, klass:col?
403   -- ROWS = cols:COLS, rows:SOME

```