

```

1 #!/usr/bin/env lua
2 -- vim : ft=lua et sts=2 sw=2 ts=2 :
3
4 -----
5
6
7
8
9
10
11
12
13
14 -- keys0: understand "N" items by peeking at at few (maybe zero) items.
15 -- Copyright 2022, Tim Menzies, MIT license
16
17 -- Permission is hereby granted, free of charge, to any person obtaining a copy
18 -- of this software and associated documentation files (the "Software"), to
19 -- deal in the Software without restriction, including without limitation the
20 -- rights to use, copy, modify, merge, publish, distribute, sublicense, and/or
21 -- sell copies of the Software, and to permit persons to whom the Software is
22 -- furnished to do so, subject to the following conditions:
23
24 -- The above copyright notice and this permission notice shall be included in
25 -- all copies or substantial portions of the Software.
26
27 -- THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR
28 -- IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY,
29 -- FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE
30 -- AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER
31 -- LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING
32 -- FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS
33 -- IN THE SOFTWARE.
34 -----
35
36 local your = {} -- user settings (may be changes from command-line)
37 local our = {} -- system settings (controlled internal to code)
38 our.help = {}
39
40 --/keys0 [OPTIONS]
41 Understand "N" items by peeking at at few (maybe zero) items.
42 (c) 2022, Tim Menzies, opensource.org/licenses/MIT
43
44 -ample max items in a 'SAMPLE' : 512
45 -better prune best half of each split : true
46 -Debug one crash, show stackdump : true
47 -dull small effect if 'dull'*sd : .35
48 -far for far, skip after 'far' : .9
49 -file load data from file : ./../data/auto93.csv
50 -h show help : false
51 -goal smile,frown,xplor, doubt : smile
52 -p coefficient on distance calcs : 2
53 -round round numbers to 'round' : 2
54 -seed random number seed : 10019
55 -Some max number items to explore : 512
56 -Tiny bin size = #t~/Tiny' : .5
57 -todo start up action ('all'=every) : -]]
58
59 our.b4={} -- globals known, pre-code. used to find stray globals
60 for k,_ in pairs(_ENV) do our.b4[k]=k end
61
62 local add, any, asserts,coerce, col, copy, csv, defaults, dist
63 local firsts, fmt, klass, map, main, new,o, push, rand, randi, rnd, rnds
64 local same, seconds, slots, sort, xpects
65
66 function klass(s, it)
67 it = {__is=s, __toString=o}
68 it.__index = it
69 return setmetatable(it, {__call=function(_, ...) return it.new(...) end}) end
70
71 local COLS,EG,EGS = klass"COLS", klass"EG", klass"EGS"
72 local NUM,RANGE,SAMPLE,SYM = klass"NUM", klass"RANGE", klass"SAMPLE", klass"S
YM"
73
74 -- TODO:
75 -- - reservoir sampler
76 -- - r.s. added to num
77 -- - mergabel numbers
78 -- - add :add(x,y) to range. updates an N
79 -- - add :div() to range (so now xpect works for those as well)
80
81 -----
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190

```

```

190 -----
191 function EG.new(t) return new(EG, {cooked={}, has=t}) end
192
193 function EG.better(eg1, eg2, eggs)
194   local s1, s2, e, n, a, b = 0, 0, 10, #egs.cols.y
195   for _, col in pairs(egs.cols.y) do
196     a = col:norm(eg1.has[col.at])
197     b = col:norm(eg2.has[col.at])
198     s1 = s1 - e^(col.w * (a-b)/n)
199     s2 = s2 - e^(col.w * (b-a)/n) end
200   return s1/n < s2/n end
201
202 function EG.cols(i, cols) return map(cols, function(x) return i.has[x.at] end) end
203
204 function EG.dist(i, j, eggs, a, b, d, n)
205   d, n = 0, #egs.cols.x + 1E-31
206   for _, col in pairs(egs.cols.x) do
207     a, b = i.has[col.at], j.has[col.at]
208     d = d + col:dist(a, b) ^ your.p end
209   return (d/n) ^ (1/your.p) end
210
211 -----
212 function RANGE.new(col, lo, hi, has)
213   lo = lo or -math.huge
214   return new(RANGE, {n=0, score=nil, col=col, lo=lo, hi=hi or lo, has=has or SYM()
215   }) end
216
217 function RANGE.__tostring(i)
218   if i.lo == i.hi then return fmt("%s==%s", i.col.txt, i.lo) end
219   if i.lo == -math.huge then return fmt("%s<%s", i.col.txt, i.hi) end
220   if i.hi == math.huge then return fmt("%s>=%s", i.col.txt, i.lo) end
221   return fmt("%s<=%s<%s", i.col.txt, i.lo, i.hi) end
222
223 function RANGE.add(i, x, y)
224   i.n = n.n+1
225   i.hi = math.max(x, i.hi)
226   i.lo = math.min(x, i.lo)
227   i.has:add(y) end
228
229 function RANGE.div(i) return i.has:div() end
230
231 function RANGE.select(i, eg, x)
232   x = eg.has[i.col.at]
233   return x=="?" or i.lo <= x and x < i.hi end
234
235 function RANGE.merge(i, j, k)
236   k = RANGE(i.col, i.lo, j.hi, i.has:merged(j.has))
237   k.n = i.n + j.n
238   if k.has:div()*1.01 <= xpects(i, j) then return k end end
239
240 function RANGE.eval(i, goal)
241   local best, rest, goals = 0, 0, {}
242   if not i.score then
243     function goals.smile(b, r) return r>b and 0 or b*b/(b+r +1E-31) end
244     function goals.frown(b, r) return b<r and 0 or r*r/(b+r +1E-31) end
245     function goals.xplor(b, r) return 1/(b+r +1E-31) end
246     function goals.doubt(b, r) return 1/(math.abs(b-r) +1E-31) end
247     for x, n in pairs(i.has) do
248       if x==goal then best = best+n/i.n else rest = rest+n/i.n end end
249     i.score = best + rest < 0.01 and 0 or goals[goal](best, rest) end
250   return i.score end
251
252 -----
253 function COLS.new(eg, i, now, where)
254   i = new(COLS, {all={}, x={}, y={}})
255   for at, s in pairs(eg) do -- First row. Create the right columns
256     now = push(i.all, {s:find"^[A-Z]" and NUM or SYM}(at, s))
257     where = {s:find"-" or s:find"+"} and i.y or i.x
258     if not s:find"." then push(where, now) end end
259   return i end
260
261 function COLS.add(i, eg)
262   return map(i.all, function(col) return col:add(eg[col.at]) end) end
263
264 -----
265 function EGS.new(i) return new(EGS, {rows={}, cols=nil}) end
266
267 function EGS.add(i, eg)
268   eg = eg.has and eg.has or eg -- If eg has data buried inside, expose it.
269   if i.cols then push(i.rows, EG(i.cols:add(eg))) else i.cols=COLS(eg) end end
270
271 function EGS.clone(i, inits, j)
272   j = EGS()
273   j:add(map(i.cols.all, function(col) return col.txt end))
274   for _, x in pairs(inits or {}) do j:add(x) end
275   return j end
276
277 function EGS.cluster(i, rows)
278   local zero, one, two, ones, twos, both, a, b, c
279   zero = any(rows)
280   one = i:far(zero)
281   two, c = i:far(one)
282   ones, twos, both = i:clone(), i:clone(), {}
283   for _, eg in pairs(rows) do
284     a = eg:dist(one, i)
285     b = eg:dist(two, i)
286     push(both, {(a^2 + c^2 - b^2) / (2*c), eg}) end
287   for n, pair in pairs(sort(both, firsts)) do
288     (n <= #both//2 and ones or twos):add(pair[2]) end
289   if your.better and two:better(one, i) then ones, twos=twos, ones end
290   return ones, twos end
291
292 function EGS.bestRanges(i, top)
293   local one, two = top:cluster(i.rows)
294   local best, out, col2, tmp, ranges = math.huge
295   for n, coll in pairs(one.cols.x) do
296     col2 = two.cols.x[n]
297     ranges = coll:bestRanges(col2)
298     if #ranges > 1 then
299       tmp = xpects(ranges)
300       if tmp < best then best, out = tmp, ranges end end end
301   return out, lefts, firsts end
302
303 function EGS.far(i, eg1, fun, tmp)
304   fun = function(eg2) return {eg2, eg1:dist(eg2, i)} end
305   tmp = #i.rows > your.Some and any(i.rows, your.Some) or i.rows
306   tmp = sort(map(tmp, fun), seconds)
307   return table.unpack(tmp[#tmp*your.far//1]) end
308
309 function EGS.from(t, i)
310   i=i or EGS(); for _, eg in pairs(t) do i:add(eg) end; return i end
311
312 function EGS.mid(i, cols)
313   return map(cols or i.all, function(col) return col:mid() end) end
314
315 function EGS.read(file, i)
316   i=i or EGS(); for eg in csv(file) do i:add(eg) end; return i end
317

```

```

317 -----
318 function any(t, n)
319   if not n then return t[randi(1,#t)] end
320   u={};for j=1,n do push(u,any(t)) end; return u end
321
322 our.fails = 0
323 function asserts(test,msg)
324   msg=msg or ""
325   if test then return print(" PASS:".msg) end
326   our.fails = our.fails+1
327   print(" FAIL:".msg)
328   if your.Debug then assert(test,msg) end end
329
330 function coerce(x)
331   if x=="true" then return true elseif x=="false" then return false end
332   return tonumber(x) or x end
333
334 function copy(t,u)
335   u={}; for k,v in pairs(t) do u[k]=v end
336   return setmetatable(u, getmetatable(t)) end
337
338 function csv(file, x,row)
339   function row(x, t)
340     for y in x:gsub("%s+", ""):gmatch("[^,]+") do push(t,coerce(y)) end
341     return t
342   end
343   file = io.input(file)
344   return function()
345     x=io.read(); if x then return row(x,{}) else io.close(file) end end end
346
347 function defaults(help_string, t,fun)
348   function fun(flag,x)
349     for n,txt in ipairs(arg) do
350       if txt:sub(1,1)=="-" and flag:match("^"..txt:sub(2)..".*")
351       then x = x=="false" and"true" or x=="true" and"false" or arg[n+1] end end
352       t[flag] = coerce(x)
353     end
354     t = {}
355     help_string:gsub("\n[-|([%s+)|^%n]*%s([%s+)|", fun)
356     return t end
357
358 function firsts(a,b) return a[1] < b[1] end
359
360 function fmt(...) return string.format(...) end
361
362 function main(our,your, reset,todos)
363   our.defaults = defaults(our.help)
364   reset = function() for k,v in pairs(our.defaults) do your[k] = v end end
365   reset()
366   if your.h
367   then print(our.help)
368   else our.fails = 0
369     todos = your.todo=="all" and slots(our.go) or {your.todo}
370     for _,one in pairs(todos) do
371       if our.go[one] then our.go[one]() end
372       reset() end end
373   for k,v in pairs(_ENV) do
374     if not our.b4[k] then print("?rogues",k,type(v)) end end
375   return our.fails end
376
377 function map(t,f, u)
378   u = {};for k,v in pairs(t) do push(u,(f or same)(v)) end; return u end
379
380 our.oid=0
381 function new(mt,x)
382   our.oid = our.oid+1; x._oid = our.oid -- Everyone gets a unique id.
383   return setmetatable(x,mt) end -- Methods now delegate to 'mt'.
384
385 function o(t)
386   local u,key
387   key= function(k) return fmt(":%s %s", k, o(t[k])) end
388   if type(t) ~= "table" then return tostring(t) end
389   u = #t>0 and map(t,o) or map(slots(t),key)
390   return (t._is or "").."["..table.concat(u, " ")."]" end
391
392 function push(t,x) table.insert(t,x); return x end
393
394 your.seed = your.seed or 10019
395 function rand(lo,hi)
396   your.seed = (16807 * your.seed) % 2147483647
397   return (lo or 0) + ((hi or 1) - (lo or 0)) * your.seed / 2147483647 end
398
399 function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
400
401 function rnd(x,d, n)
402   if type(x)~="number" then return x end
403   n=10^(d or your.round)
404   return math.floor(x*n+0.5)/n end
405
406 function rnds(t,d) return map(t,function(x) return rnd(x,d) end) end
407
408 function same(x,...) return x end
409
410 function seconds(a,b) return a[2] < b[2] end
411
412 function slots(t, u)
413   u={ }
414   for k,_ in pairs(t) do if tostring(k):sub(1,1) ~= "-" then push(u,k) end end
415   return sort(u) end
416
417 function sort(t,f) table.sort(t,f); return t end
418
419 function xpects(t)
420   local sum,n = 0,0
421   for _,z in pairs(t) do n = n + z.n; sum = sum + z.n*z:div() end
422   return sum/n end
423
424 -----
425 our.go={ } -- list of enabled tests
426 our.nogo={ } -- list of disabled test
427 local go, nogo = our.go,our.nogo
428
429 function go.settings()
430   print("our",o(our))
431   print("your",o(your)) end
432
433 function go.range( r)
434   r=RANGE(NUM(10,"fred"),"apple")
435   assert(tostring(r) == "fred == apple", "print ok") end
436
437 function go.num( m,n)
438   m=NUM(); for j=1,10 do m:add(j) end
439   n=copy(m); for j=1,10 do n:add(j) end
440   asserts(2.95 == rnd(n:div()),"sd ok") end
441
442 function go.egs( egs)
443   egs = EGS.read(your.file)
444   asserts(egs.cols.y[1].hi==5140,"most seen") end
445
446 function go.clone( egs1,egs2,s1,s2)
447   egs1 = EGS.read(your.file)
448   s1 = o(egs1.cols.y)
449   egs2 = egs1:clone(egs1.rows)
450   s2 = o(egs2.cols.y)
451   asserts(s1==s2, "cloning works") end
452
453 function go.dist()
454   local egs,egl,dist,tmp,j1,j2,d1,d2,d3,one
455   egs = EGS.read(your.file)
456   egl = egs.rows[1]
457   dist = function(eg2) return {eg2,egl:dist(eg2,egs)} end
458   tmp = sort(map(egs.rows, dist), seconds)
459   one = tmp[1][1]
460   for j=1,10 do
461     j1 = randi(1,#tmp)
462     j2 = randi(1,#tmp)
463     if j1>j2 then j1,j2=j2,j1 end
464     d1 = tmp[j1][1]:dist(one,egs)
465     d2 = tmp[j2][1]:dist(one,egs)
466     asserts(d1 <= d2,"distance") end end
467
468 function go.cluster( top,left,right)
469   top = EGS.read(your.file)
470   left, right = top:cluster()
471   for n,t in pairs{top,left,right} do print(n,o(rnds(t:mid(t.cols.y)))) end
472 end
473
474 os.exit( main(our,your) )

```