

```

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
73 YM"
74
75 -- TODO:
76 -- - resorvoir sampler
77 -- - r.s. added to num
78 -- - mergabel numbers
79 -- - add :add(x,y) to range. updates an N
80 -- - add :div() to range (so now xpect works for those as well)

```

```

80 -----
81 local SAMPLE=klass"SAMPLE"
82 function SAMPLE.new() return new(SAMPLE,{n=0, all={}, max=your.ample}) end
83
84 function SAMPLE.add(i,x, pos)
85   i.n = i.n + 1
86   if #i.all < i.max then pos = #i.all + 1
87   elseif rand() < #i.all/i.n then pos = #i.all * rand() end
88   if pos then i.all[pos//1] = x end end
89
90 -----
91 function NUM.new(at,s, i)
92   i = new(NUM,{n=0,at=at or 0, txt=s or "",_has=SAMPLE(),
93     mu=0,m2=0,lo=math.huge,hi=-math.huge})
94   i.w = i.txt:find "-" and -1 or 1
95   return i end
96
97 function NUM.add(i,x, d)
98   if x=="?" then
99     i.n = i.n + 1
100    d = x - i.mu
101    i.mu = i.mu + d/i.n
102    i.m2 = i.m2 + d*(x-i.mu)
103    i.has:add(x)
104    i.lo = math.min(i.lo,x); i.hi = math.max(i.hi,x) end
105   return x end
106
107 function NUM.dist(i,x,y)
108   if x=="?" and y=="?" then return 1
109   elseif y=="?" then y = i:norm(y); x=y>.5 and 0 or 1
110   elseif x=="?" then x = i:norm(x); y=x>.5 and 0 or 1
111   else x,y = i:norm(x), i:norm(y) end
112   return math.abs(x-y) end
113
114 function NUM.div(i) return i.n<2 and 0 or (i.m2/(i.n-1))^0.5 end
115
116 function NUM.merged(i,j)
117   k = NUM(i.at, i.txt)
118   for _,x in pairs(i._has.all) do k:add(x) end
119   for _,x in pairs(j._has.akk) do k:add(x) end
120   return k end
121
122 function NUM.mid(i) return i.mu end
123
124 function NUM.norm(i,x) return i.hi-i.lo<1E-9 and 0 or (x-i.lo)/(i.hi-i.lo) end
125
126 function NUM.ranges(i,j, yklass)
127   local xys, dull, tiny, range,out
128   yklass = yklass or SYM
129   xys = {}
130   for _,x in pairs(i._has.all) do push(xys, {x=x, y="best"}) end
131   for _,x in pairs(j._has.all) do push(xys, {x=x, y="rest"}) end
132   xys = sort(xys, function(a,b) return a.x < b.x end)
133   dull = xpects(i,j)*your.dull
134   tiny = (#xys)^your.Tiny
135   range = RANGE(i,xys[1].x, xys[1].x, yklass())
136   out = {range}
137   for k,xy in pairs(xys) do
138     if k < #xys - tiny and xy.x == xys[k+1].x and
139       range.has.n > tiny and range.hi - range.lo > dull
140     then range = push(out, RANGE(i, range.hi, xy.x, yklass()))
141     end
142     range:add(xy.x, xy.y) end
143   out[1].lo = -math.huge
144   out[#ranges].hi = math.huge
145   return out end
146
147 function NUM.superRanges(i,b4)
148   local j,tmp,now,after,maybe = 0, {}
149   while j < #b4 do
150     j = j + 1
151     now, after = b4[j], b4[j+1]
152     if after then
153       maybe = now:merge(after)
154       if maybe then now=maybe; j=j+1 end end
155     push(tmp,now) end
156   return #tmp==#b4 and b4 or i:superRanges(tmp) end
157
158 -----
159 function SYM.new(at,s)
160   return new(SYM,{n=0, at=at or 0, txt=s or "", has={}, most=0, mode=nil}) end
161
162 function SYM.add(i,x,count)
163   count = count or 1
164   i.has[x] = count + (i.has[x] or 0)
165   if i.has[x] > i.most then i.most, i.mode = i.has[x], x end
166   return x end
167
168 function SYM.dist(i,x,y) return x=="?" and y=="?" and 1 or x==y and 0 or 1 end
169
170 function SYM.div(i, e)
171   e=0; for _,n in pairs(i.has) do e=e-n/i.n*math.log(n/i.n,2) end; return e end
172
173 function SYM.merged(i,j, k)
174   k = SYM(i.at, i.txt)
175   for x,count in pairs(i.has) do k:add(x,count) end
176   for x,count in pairs(j.has) do k:add(x,count) end
177   return k end
178
179 function SYM.mid(i) return i.mode end
180
181 function SYM.ranges(i,j, ranges,t,n,xpect)
182   t,out = {},{}
183   for x,n in pairs(i.has) do t[x] = t[x] or SYM(); t[x]:add("best",n) end
184   for x,n in pairs(j.has) do t[x] = t[x] or SYM(); t[x]:add("rest",n) end
185   for x,stats in pairs(t) do push(out, RANGE(i,x,x,stats)) end
186   return out end
187
188 function SYM.superRanges(i, ranges) return ranges end
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()}) end
215
216 function RANGE.__tostring(i)
217   if i.lo == i.hi then return fmt("%s==%s", i.col.txt, i.lo) end
218   if i.lo == -math.huge then return fmt("%s<%s", i.col.txt, i.hi) end
219   if i.hi == math.huge then return fmt("%s>=%s", i.col.txt, i.lo) end
220   return fmt("%s<=%s<%s", i.col.txt, i.lo, i.hi) end
221
222 function RANGE.add(i, x, y)
223   i.n = i.n + 1
224   i.hi = math.max(x, i.hi)
225   i.lo = math.min(x, i.lo)
226   i.has:add(y) end
227
228 function RANGE.div(i) return i.has:div() end
229
230 function RANGE.select(i, eg, x)
231   x = eg.has[i.col.at]
232   return x=="?" or i.lo <= x and x < i.hi end
233
234 function RANGE.merge(i, j, k)
235   k = RANGE(i.col, i.lo, j.hi, i.has:merged(j.has))
236   k.n = i.n + j.n
237   if k.has:div()*1.01 <= xpects(i, j) then return k end end
238
239 function RANGE.eval(i, goal)
240   local best, rest, goals = 0, 0, {}
241   if not i.score then
242     function goals.smile(b, r) return r>b and 0 or b*b/(b+r +1E-31) end
243     function goals.frown(b, r) return b<r and 0 or r*r/(b+r +1E-31) end
244     function goals.xplor(b, r) return 1/(b+r +1E-31) end
245     function goals.doubt(b, r) return 1/(math.abs(b-r) +1E-31) end
246     for x, n in pairs(i.has) do
247       if x==goal then best = best+n/i.n else rest = rest+n/i.n end end
248     i.score = best + rest < 0.01 and 0 or goals[goal](best, rest) end
249   return i.score end
250
251
252 function COLS.new(eg, i, now, where)
253   i = new(COLS, {all={}, x={}, y={}})
254   for at, s in pairs(eg) do -- First row. Create the right columns
255     now = push(i.all, (s:find("[A-Z]" and NUM or SYM)(at, s)))
256     where = (s:find("-" or s:find("+") and i.y or i.x)
257       if not s:find("-" then push(where, now) end end
258   return i end
259
260 function COLS.add(i, eg)
261   assert(#eg == #i.all, "expected a different number of cells")
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, EGS(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.far(i, eg1, fun, tmp)
293   fun = function(eg2) return {eg2, eg1:dist(eg2, i)} end
294   tmp = #i.rows > your.Some and any(i.rows, your.Some) or i.rows
295   tmp = sort(map(tmp, fun), seconds)
296   return table.unpack(tmp[#tmp*your.far//1]) end
297
298 function EGS.from(t, i)
299   i=i or EGS(); for _, eg in pairs(t) do i:add(eg) end; return i end
300
301 function EGS.mid(i, cols)
302   return map(cols or i.all, function(col) return col:mid() end) end
303
304 function EGS.read(file, i)
305   i=i or EGS(); for eg in csv(file) do i:add(eg) end; return i end
306
307 function EGS.superRanges(i, top)
308   local one, two = top:cluster(i.rows)
309   local best, out, col2, tmp, ranges = math.huge
310   for n, coll in pairs(one.cols.x) do
311     col2 = two.cols.x[n]
312     ranges = coll:ranges(col2)
313     ranges = coll:superRanges(ranges)
314     if #ranges > 1 then
315       tmp = xpects(ranges)
316       if tmp < best then best, out = tmp, ranges end end end
317   return out, lefts, firsts end
318
319

```

```

319 -----
320 function any(t, n)
321   if not n then return t[randi(1,#t)] end
322   u={};for j=1,n do push(u,any(t)) end; return u end
323
324 our.fails = 0
325 function asserts(test,msg)
326   msg=msg or ""
327   if test then return print(" PASS:".msg) end
328   our.fails = our.fails+1
329   print(" FAIL:".msg)
330   if your.Debug then assert(test,msg) end end
331
332 function coerce(x)
333   if x=="true" then return true elseif x=="false" then return false end
334   return tonumber(x) or x end
335
336 function copy(t,u)
337   u={}; for k,v in pairs(t) do u[k]=v end
338   return setmetatable(u, getmetatable(t)) end
339
340 function csv(file, x,row)
341   function row(x, t)
342     for y in x:gsub("%s+", ""):gmatch("[^,]+") do push(t,coerce(y)) end
343     return t
344   end
345   file = io.input(file)
346   return function()
347     x=io.read(); if x then return row(x, {}) else io.close(file) end end end
348
349 function userSettings(help_string, t,fun)
350   function fun(flag,x)
351     for n,txt in ipairs(arg) do
352       if txt:sub(1,1)=="-" and flag:match("^"..txt:sub(2)..".*")
353       then x = x=="false" and"true" or x=="true" and"false" or arg[n+1] end end
354       t[flag] = coerce(x)
355     end
356     t = {}
357     help_string:gsub("\n [~]([%s+])\n"[%s+])", fun)
358     return t end
359
360 function firsts(a,b) return a[1] < b[1] end
361
362 function fmt(...) return string.format(...) end
363
364 function main(user, system, todos)
365   local function reset()
366     for k,v in pairs(userSettings(system.help)) do user[k]=v end end
367   reset()
368   if user.h
369   then print(system.help)
370   else system.fails = 0
371     todos = user.todos=="all" and slots(system.go) or {user.todos}
372     for _,one in pairs(todos) do
373       if type(system.go[one])=="function" then system.go[one]() end
374     reset() end end
375   for k,v in pairs(_ENV) do
376     if not system.b4[k] then print("?rogues",k,type(v)) end end
377   return system.fails end
378
379 function map(t,f, u)
380   u = {};for k,v in pairs(t) do push(u,(f or same)(v)) end; return u end
381
382 our.oid=0
383 function new(mt,x)
384   our.oid = our.oid+1; x._oid = our.oid -- Everyone gets a unique id.
385   return setmetatable(x,mt) end -- Methods now delegate to 'mt'.
386
387 function o(t)
388   local u,key
389   key= function(k) return fmt(":%s %s", k, o(t[k])) end
390   if type(t) ~= "table" then return tostring(t) end
391   u = #t>0 and map(t,o) or map(slots(t),key)
392   return (t._is or "").."["..table.concat(u, " ")."]" end
393
394 function push(t,x) table.insert(t,x); return x end
395
396 your.seed = your.seed or 10019
397 function rand(lo,hi)
398   your.seed = (16807 * your.seed) % 2147483647
399   return (lo or 0) + ((hi or 1) - (lo or 0)) * your.seed / 2147483647 end
400
401 function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
402
403 function rnd(x,d, n)
404   if type(x)~="number" then return x end
405   n=10^(d or your.round)
406   return math.floor(x*n+0.5)/n end
407
408 function rnds(t,d) return map(t,function(x) return rnd(x,d) end) end
409
410 function same(x,...) return x end
411
412 function seconds(a,b) return a[2] < b[2] end
413
414 function slots(t, u)
415   u={ }
416   for k,_ in pairs(t) do if tostring(k):sub(1,1) ~= "-" then push(u,k) end end
417   return sort(u) end
418
419 function sort(t,f) table.sort(t,f); return t end
420
421 function xpects(t)
422   local sum,n = 0,0
423   for _,z in pairs(t) do n = n + z.n; sum = sum + z.n*z:div() end
424   return sum/n end
425
426 -----
427 our.go={ } -- list of enabled tests
428 our.nogo={ } -- list of disabled test
429 local go, nogo = our.go,our.nogo
430
431 function go.settings()
432   print("our",o(our))
433   print("your",o(your)) end
434
435 function go.range( r)
436   r=RANGE(NUM(10,"fred"),"apple")
437   assert(tostring(r) == "fred == apple", "print ok") end
438
439 function go.num( m,n)
440   m=NUM(); for j=1,10 do m:add(j) end
441   n=copy(m); for j=1,10 do n:add(j) end
442   asserts(2.95 == rnd(n:div()),"sd ok") end
443
444 function go.egs( egs)
445   egs = EGS.read(your.file)
446   asserts(egs.cols.y[1].hi==5140,"most seen") end
447
448 function go.clone( egs1,egs2,s1,s2)
449   egs1 = EGS.read(your.file)
450   s1 = o(egs1.cols.y)
451   egs2 = egs1:clone(egs1.rows)
452   s2 = o(egs2.cols.y)
453   asserts(s1==s2, "cloning works") end
454
455 function go.dist()
456   local egs,egl,dist,tmp,j1,j2,d1,d2,d3,one
457   egs = EGS.read(your.file)
458   egl = egs.rows[1]
459   dist = function(eg2) return {eg2,egl:dist(eg2,egs)} end
460   tmp = sort(map(egs.rows, dist), seconds)
461   one = tmp[1][1]
462   for j=1,10 do
463     j1 = randi(1,#tmp)
464     j2 = randi(1,#tmp)
465     if j1>j2 then j1,j2=j2,j1 end
466     d1 = tmp[j1][1]:dist(one,egs)
467     d2 = tmp[j2][1]:dist(one,egs)
468     asserts(d1 <= d2,"distance") end end
469
470 function go.cluster( top,left,right)
471   top = EGS.read(your.file)
472   left, right = top:cluster()
473   for n,t in pairs{top,left,right} do print(n,o(rnds(t:mid(t.cols.y)))) end
474   end
475
476 -- assuming our.go = demos and our.help==help string and our.fails = 0 then...
477 os.exit( main(your, our))

```