```lua
1  #!/usr/bin/env lua
2  --                                        __
3  --                                       /  \
4  --     _____                       _    /    \/\/\
5  --   /\ '-,'-'\              /'-,  /'-'\  \  \/\/ <
6  --   \ \ \_L\ |    /\/_/     /'-,  \  '-'\   \ \'\
7  --    \ \ \__/    /  \___/   \  \___    \  \__/  \_
8  --     \ \ \ \/   \/\___/    \/\___/     \/_/\/_/
9  --      \  \ \
10 --       \/_/
11
12 local your, our={}, {b4={}, help=[[
13 peek.lua [OPTIONS]
14 (c)2022 Tim Menzies, MIT license (2 clause)
15 Understand N items after log(N) probes, or less.
16
17  -file    ../../data/auto93.csv
18  -ample   512
19  -far     .9
20  -best    .5
21  -help    false
22  -dull    .35
23  -rest    3
24  -seed    10019
25  -rnd     %.2f
26  -task    -
27  -p       2]]}
28
29 for k,_ in pairs(_ENV) do our.b4[k] = k end
30 local any,asserts,cells,copy,firsts,fmt,go,id,main,many,map
31 local merge,new,o,push,rand,randi,ranges,rnd,rogues,rows,same
32 local seconds,settings,slots,sort,super,thing,things,xpect
33 local COLS,EG,EGS,NUM,RANGE,SAMPLE,SYM
34 local class= function(t,  new)
35   function new(_,...) return t.new(...) end
36   t.__index=t
37   return setmetatable(t,{__call=new}) end
38
39 -- Copyright (c) 2022, Tim Menzies
40 --
41 -- Redistribution and use in source and binary forms, with or without
42 -- modification, are permitted provided that the following conditions are met.
43 -- (1) Redistributions of source code must retain the above copyright notice,
44 -- this list of conditions and the following disclaimer.  (2) Redistributions
45 -- in binary form must reproduce the above copyright notice, this list of
46 -- conditions and the following disclaimer in the documentation and/or other
47 -- materials provided with the distribution.
48 --
49 -- THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS
50 -- IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO,
51 -- THE IMPLIED WARRANTIES OF MERCHNTABILITY AND FITNESS FOR A PARTICULAR
52 -- PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR
53 -- CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL,
54 -- EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO,
55 -- PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR
56 -- PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF
57 -- LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING
58 -- NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS
59 -- SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE
60
```

```lua
60  --
61  --       _
62  --      | |
63  --   ___| | __ _ ___ ___  ___  ___
64  --  / __| |/ _` / __/ __|/ _ \/ __|
65  COLS=class{}
66
67  function COLS.new(t,      i,where,now)
68    i = new({all={}, x={}, y={}},COLS)
69    for at,s in pairs(t) do
70      now = push(i.all, (s:find"^[A-Z]" and NUM or SYM)(at,s))
71      if not s:find":" then
72        push((s:find"-" or s:find"+") and i.y or i.x, now) end end
73    return i end
74
75  function COLS.__tostring(i, txt)
76    function txt(c) return c.txt end
77    return fmt("COLS{:all %s\n\t:x %s\n\t:y %s", o(i.all,txt), o(i.x,txt), o(i.y,txt)) end
78
79  function COLS.add(i,t,       add)
80    function add(col,    x) x=t[col.at]; col:add(x);return x end
81    return map(i.all, add) end
82  -- ------------------------------------------------------------------------
83  EG=class{}
84  function EG.new(t) return new({has=t, id=id()},EG) end
85
86  function EG.__tostring(i) return fmt("EG%s%s %s", i.id,o(i.has),#i.has) end
87
88  function EG.better(i,j,cols)
89    local s1,s2,e,n,a,b = 0,0,10,#cols
90    for _,col in pairs(cols) do
91      a  = col:norm(i.has[col.at])
92      b  = col:norm(j.has[col.at])
93      s1 = s1 - e^(col.w * (a-b)/n)
94      s2 = s2 - e^(col.w * (b-a)/n) end
95    return s1/n < s2/n end
96
97  function EG.dist(i,j,egs,      a,b,d,n)
98    d,n = 0, #egs.cols.x + 1E-31
99    for _,col in pairs(egs.cols.x) do
100     a,b = i.has[col.at], j.has[col.at]
101     d   = d + col:dist(a,b) ^ your.p end
102   return (d/n) ^ (1/your.p) end
103 -- ------------------------------------------------------------------------
104 EGS=class{}
105 function EGS.new() return new({rows={}, cols=nil}, EGS) end
106
107 function EGS.__tostring(i) return fmt("EGS{#rows %s:cols %s", #i.rows,i.cols) end
108
109 function EGS.add(i,row)
110   row = row.has and row.has or row
111   if i.cols then push(i.rows,EG(i.cols:add(row))) else i.cols=COLS(row) end end
112
113 function EGS.clone(i,inits,     j)
114   j = EGS()
115   j:add(map(i.cols.all, function(col) return col.txt end))
116   for _,x in pairs(inits or {}) do  j:add(x) end
117   return j end
118
119 function EGS.far(i,eg1,rows,    fun,tmp)
120   fun = function(eg2) return {eg2, eg1:dist(eg2,i)} end
121   tmp = sort(map(rows, fun), seconds)
122   return table.unpack(tmp[#tmp*your.far//1] ) end
123
124 function EGS.file(i,file) for row in rows(file) do i:add(row) end; return i end
125
126 function EGS.mid(i,cols,      mid)
127   function mid(col)  return col:mid() end
128   return map(cols or i.cols.all, mid) end
129
130 function EGS.halve(i,rows)
131   local c,l,r,ls,rs,cosine,some
132   function cosine(row,      a,b)
133     a,b = eg:dist(l,i),eg:dist(r,i); return {(a^2+c^2-b^2)/(2*c),row} end
134   some  = #rows > your.ample and many(rows, your.ample) or rows
135   l     = i:far(any(rows), some)
136   r,c   = i:far(l,          some)
137   ls,rs = i:clone(), i:clone()
138   for n,pair in pairs(sort(map(rows,cosine), firsts)) do
139     (n <= #rows//2 and ls or rs):add(pair[2]) end
140   return ls,rs,l,r,c end
141
142 function EGS.splitter(i,top,    ls,rs,there,ranges)
143   ls,rs = (top or i):halve(i.rows)
144   ranges = {}
145   for n,here in pairs(ls.cols.xs) do
146     there = rs.cols.xs[n]
147     for range in pairs(here:ranges(there)) do
148       push(ranges,range) end end
149   return sort(ranges)[1] end
150
151 function EGS.xcluster(i,top,lvl)
152   local split, left, right,kid1, kid2
153   top, lvl = top or i, lvl or 0
154   if #i.rows >= 2*(#top.rows)^your.small then
155     split, kid1, kid2 = i:splitter(top), i:clone(), i:clone()
156     for _,row in pairs(i.rows) do
157       (split:selects(row) and kid1 or kid2):add(row) end
158     if #kid1.rows ~= #i.rows then left  = kid1:xcluster(top,lvl+1) end
159     if #kid2.rows ~= #i.rows then right = kid2:xcluster(top,lvl+1) end
160   end
161   return {here=i, split=split, left=left, right=right} end
162
```

```lua
162  -- ----------------------------------------------------------------------------
163  NUM=class{}
164  function NUM.new(at,s, big)
165    big = math.huge
166    return new({lo=big, hi=-big, at=at or 0, txt=s or "",
167               n=0, mu=0, m2=0, sd=0,_all=SAMPLE(),
168               w=(s or ""):find"-" and -1 or 1},NUM) end
169
170  function NUM.__tostring(i)
171    return fmt("NUM{:at %s :txt %s :n %s :lo %s :hi %s :mu %s :sd %s}",
172               i.at, i.txt,  i.n, i.lo, i.hi, rnd(i.mu), rnd(i:div())) end
173
174  function NUM.add(i,x,    d,pos)
175    if x~="?" then
176      i.n  = i.n+1
177      d    = x - i.mu
178      i.mu = i.mu + d/i.n
179      i.m2 = i.m2 + d*(x-i.mu)
180      i.lo = math.min(x,i.lo); i.hi = math.max(x,i.hi)
181      i._all:add(x) end
182    return x end
183
184  function NUM.dist(i,a,b)
185    if     a=="?" and b=="?" then a,b =1,0
186    elseif a=="?"            then b   = i:norm(b); a=b>.5 and 0 or 1
187    elseif b=="?"            then a   = i:norm(a); b=a>.5 and 0 or 1
188    else                          a,b = i:norm(a), i:norm(b) end
189    return math.abs(a-b) end
190
191  function NUM.div(i) return i.n <2 and 0 or (i.m2/(i.n-1))^0.5 end
192
193  function NUM.merged(i,j)
194    k= NUM(i.at, i.txt)
195    for _,x in pairs(i._all,it) do k:add(x) end
196    for _,x in pairs(j._all,it) do k:add(x) end
197    return k end
198
199  function NUM.mid(i) return i.mu end
200
201  function NUM.norm(i,x) return i.hi-i.lo < 1E-9 and 0 or (x-i.lo)/(i.hi-i.lo) end
202
203  function NUM.ranges(i,j,ykind,      xys)
204    xys={}
205    for _,x in pairs(i._all,it) do push(xys,{x=x,y="best"}) end
206    for _,x in pairs(j._all,it) do push(xys,{x=x,y="rest"}) end
207    return merge(ranges(xys,i, ykind or SYM,
208                #xys^your.dull, xpect(i,j)*your.small)) end
209  -- ----------------------------------------------------------------------------
210  RANGE=class{}
211  function RANGE.new(col,hi,lo,ys)
212    return new({n=0,cols=-col,lo=lo,hi=hi or lo, ys=ys or SYM()},RANGE) end
213
214  function RANGE.__lt(i,j) return i:div() < j:div() 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.lo, i.col.txt, i.hi) end
221
222  function RANGE.add(i,x,y,inc)
223    inc  = inc or 1
224    i.n  = i.n + inc
225    i.lo = math.min(x,i.lo)
226    i.hi = math.max(x,i.hi)
227    i.ys:add(y, inc) end
228
229  function RANGE.div(i) return i.ys:div() end
230
231  function RANGE.selects(i,row,    x)
232    x=row.has[col.at]; return x=="?" or i.lo<=x and x<i.hi end
233  -- ----------------------------------------------------------------------------
234  SAMPLE=class{}
235  function SAMPLE.new() return new({n=0,it={},ok=false,max=your.ample},SAMPLE) end
236
237  function SAMPLE.add(i,x,    pos)
238    i.n = i.n + 1
239    if     #i.it < i.max      then pos= #i.it + 1
240    elseif rand() < #i.it/i.n then pos= #i.it * rand() end
241    if pos then i.ok = false; i.it[pos//1]= x end end
242
243  function SAMPLE.all(i) if not i.ok then i.ok=true;sort(i.it)end; return i.it end
244  -- ----------------------------------------------------------------------------
245  SYM=class{}
246  function SYM.new(at,s)
247    return new({at=at or 0,txt=s or "",has={},n=0,most=0,mode=nil},SYM) end
248
249  function SYM.__tostring(i)
250    return fmt("SYM{:at %s :txt %s :mode %s :has %s}",
251               i.at, i.txt, i.mode, o(i.has)) end
252
253  function SYM.add(i,x, inc)
254    if x ~= "?" then
255      inc = inc or 1
256      i.n = i.n+inc
257      i.has[x] = inc  + (i.has[x] or 0)
258      if i.has[x] > i.most then i.most, i.mode = i.has[x], x end end
259    return x end
260
261  function SYM.dist(i,a,b) return a=="?" and b=="?" and 1 or a==b and 0 or 1 end
262
263  function SYM.div(i)
264    e=0;for _,v in pairs(i.has) do e=e - v/i.n*math.log(v/i.n,2) end; return e end
265
266  function SYM.merge(i,j,    k)
267    k= SYM(i.at, i.txt)
268    for x,count in pairs(i.has) do k:add(x,count) end
269    for x,count in pairs(j.has) do k:add(x,count) end
270    return k end
271
272  function SYM.mid(i) return i.mode end
273
274  function SYM.ranges(i,j,    t)
275    t = {}
276    for _,pair in pairs{{i.has,"bests"}, {j.has,"rests"}} do
277      for x,inc in pairs(pair[1]) do
278        t[x] = t[x] or RANGE(i,x)
279        t[x]:add(x, pair[2], inc) end end
280    return map(t, same) end
281  -- ----------------------------------------------------------------------------
```

```lua
281  --
282  --        _____ __  _____ __ _
283  --       / __(_)_ _/ /_ __ _/ // /
284  --      / _// / // / _ \(_-</ _  /
285  --     /_/ /_/\_,_/_//_/___/\_,_/
286
287  fmt  = string.format
288  new  = setmetatable
289  same = function(x,...) return x end
290
291  function asserts(test,msg)
292    msg=msg or ""
293    if test then return print("PASS:"..msg) end
294    our.failures = our.failures + 1
295    print("FAIL:"..msg)
296    if your.Debug then assert(test,msg) end end
297
298  function copy(t,    u)
299    if type(t) ~= "table" then return t end
300    u={};for k,v in pairs(t) do u[k]=copy(v) end;return new(u,getmetatable(t)) end
301
302  function firsts(a,b) return a[1] < b[1] end
303
304  function id() our.id = 1+(our.id or 0); return our.id end
305
306  function many(t,n, u) u={};for j=1,n do push(u,any(t)) end; return u end
307
308  function map(t,f,  u)
309    u={};for _,v in pairs(t) do u[1+#u]=(f or same)(v) end; return u end
310
311  function o(t,f,  u,key)
312    key= function(k)
313           if t[k] then return fmt(":%s %s", k, rnd((f or same)(t[k]))) end end
314    u = #t>0 and map(map(t,f),rnd) or map(slots(t),key)
315    return "{"..table.concat(u, "").."}" end
316
317  function rand(lo,hi)
318    your.seed = (16807 * your.seed) % 2147483647
319    return (lo or 0) + ((hi or 1) - (lo or 0)) * your.seed / 2147483647 end
320
321  function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
322
323  function push(t,x)  table.insert(t,x); return x end
324
325  function rnd(x)
326    return fmt(type(x)=="number" and x~=x//1 and your.rnd or"%s",x) end
327
328  function rows(file,    x)
329    file =io.input(file)
330    return function()
331      x=io.read(); if x then return things(x) else io.close(file) end end end
332
333  function main(      defaults,tasks)
334    tasks = your.task=="all" and slots(go) or {your.task}
335    defaults=copy(your)
336    our.failures=0
337    for _,x in pairs(tasks) do
338      if type(our.go[x]) == "function" then our.go[x]() end
339      your = copy(defaults) end
340    rogues()
341    return our.failures end
342
343  function merge(b4,      j,tmp,merged,one,two)
344    j, tmp = 0, {}
345    while j < #b4 do
346      j = j + 1
347      one, two = b4[j], b4[j+1]
348      if two then
349        merged = one.ys:merge(two.ys)
350        if merged:div()*1.01 <= xpect(one.ys, two.ys) then
351          j   = j+1
352          one = RANGE(one.col, one.lo, two.hi, merged) end end
353      push(tmp,one) end
354    return #tmp==#b4 and b4 or merge(tmp) end
355
356  function ranges(xys,col,ykind, dull, small,      one,out)
357    one, xys = {}, sort(xys, function(a,b) return a.x < b.x end)
358    one      = push(out, RANGE(col, xys[1].x, xys[1].x, ykind()))
359    for j,xy in pairs(xys) do
360      if    j < #xys - small   and -- enough items remaining after split
361            xy.x ~= xys[j+1].x and -- next item is different (so can split here)
362            one.n > small      and -- one has enough items
363            one.hi - one.lo > dull -- one is not trivially small
364      then one = push(out, RANGE(col, one.hi, xy.x, ykind())) end
365      one:add(xy,x, xy.y) end
366    out[1].lo    = -math.huge
367    out[#out].hi =  math.huge
368    return out end
369
370  function rogues()
371    for k,v in pairs(_ENV) do
372      if not our.b4[k] then print("??",k,type(v)) end end end
373
374  function seconds(a,b) return a[2] < b[2] end
375
376  function settings(help,   t)
377    t={}
378    help:gsub("\n [-]([^%s]+)[^\n]*%s([^%s]+)", function(slot, x)
379      for n,flag in ipairs(arg) do
380        if   flag:sub(1,1)=="-" and slot:match("^"..flag:sub(2)..".*")
381        then x=x=="false" and "true" or x=="true" and "false" or arg[n+1] end end
382      t[slot] = thing(x) end)
383    if t.help then print(t.help) end
384    return t end
385
386  function slots(t,u) u={};for x,_ in pairs(t) do u[1+#u]=x end;return sort(u) end
387
388  function sort(t,f)  table.sort(t,f); return t end
389
390  function thing(x)
391    x = x:match"^%s*(.-)%s*$"
392    if x=="true" then return true elseif x=="false" then return false end
393    return tonumber(x) or x end
394
395  function things(x,sep,   t)
396    t={};for y in x:gmatch(sep or"([^,]+)") do t[1+#t]=thing(y) end; return t end
397
398  function xpect(...)
399    m,d=0,0; for _,z in pairs{...} do n=n+z.n; d=d+z.n*z:div() end; return d/n end
400
```

```
400  --
401  --    _          _
402  --   | |_ ___ ___| |_ ___
403  --   |  _/ -_|_-<  _/ _/
404  --    \__\___/__/\__/__/

406  our.go, our.no = {},{}; go=our.go
407  function go.settings() print("your",o(your)) end

409  function go.sample() print(EGS():file(your.file)) end

411  function go.clone( a,b)
412    a= EGS():file(your.file)
413    b= a:clone(a.rows)
414    asserts(#a.egs == #b.egs, tostring(a.cols.all[1])==tostring(b.cols.all[1]),"cl
  oning")
415    asserts(tostring(a.cols.all[1])==tostring(b.cols.all[1]),"cloning")
416  end

418  function go.sort(    i,a,b)
419    i   = EGS():file(your.file)
420    a,b = i:bestRest()
421    a,b = i:clone(a), i:clone(b)
422    print(#a.rows)
423    print(a.cols.all[1])
424    print("all",  o(i:mid(i.cols.y)))
425    print("best", o(a:mid(a.cols.y)))
426    print("rest", o(b:mid(b.cols.y)))
427  end

429  your = settings(our.help)
430  os.exit( main() )
```