

```

1 local help = {}
2
3 15 == a little lab of lots of LUA learning algorithms.
4 (c) 2022, Tim Menzies, BSD 2-clause license.
5
6 USAGE:
7 lua 15.lua [OPTIONS]
8
9 OPTIONS:
10 -cohen F Cohen's delta = .35
11 -data N data file = etc/data/auto93.csv
12 -dump stack dump on assert fails = false
13 -furthest F far = 9
14 -Format s format string = %5.2f
15 -keep P max kept items = 512
16 -p P distance coefficient = 2
17 -seed P set seed = 10019
18 -todo S start up action (or 'all') = nothing
19 -help show help = false
20 -want F recurse until rows^want = .5
21
22 KEY: N=fileName F=float P=posint S=string
23
24 NOTES: This code uses Aha's distance measure [Aha91] (that can
25 handle numbers and symbols) to recursively divide data based on 20
26 distant points (these two are found in linear time using the Fastmap
27 heuristic [Fa95]).
28
29 To avoid spurious outliers, this code use the 90% furthest points.
30
31 To avoid long runtimes, uses a subset of the data to learn where
32 to divide data (then all the data gets pushed down first halves).
33
34 To support explanation, optionally, at each level of recursion,
35 this code reports what ranges can best distinguish sibling clusters
36 C1,C2. The discretizer is inspired by the ChiMerge algorithm:
37 numerics are divided into, say, 16 bins. Then, while we can find
38 adjacent bins with the similar distributions in C1,C2, then
39 (a) merge them (b) look for other merges.
40 ]]
41
42 -- ## Namespace
43
44 -- Cache current globals, use at end to find rogue variables
45 local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
46
47 -- Defined local names.
48 local any, asserts, big, cli, csv, fails, firsts, fmt, goalp, ignorep, klassp
49 local lessp, map, main, many, max, merge, min, morep, new, nump, o, oo, per, pop, push
50 local r, rows, rnd, rnds, slots, sort, sum, thing, things, unpack
51
52 -- Classes have UPPER CASE names.
53 local CLUSTER, COLS, EGS, EXPLAIN, NUM, ROWS = {}, {}, {}, {}, {}
54 local SKIP, SOME, SPAN, SYM = {}, {}, {}, {}
55
56 -- ## Settings
57 -- Parse the help text for flags and defaults (e.g. -keep, 512).
58 -- Check for updates on those details from command line
59 -- (and and there,
60 -- some shortcuts are available;
61 -- e.g. _k N _&Arr; 'keep=N';
62 -- and _booleanFlag _&Arr; 'booleanFlag=not default').
63 local the={}
64 help:gsub("\n [-|[%s+]|^%n)*%s([%s+])", function(key, x)
65 for n, flag in ipairs(arg) do
66 if flag:sub(1,1)=="-" and key:find("^"..flag:sub(2)..".") then
67 x = x=="false" and true or x=="true" and "false" or arg[n+1] end and
68 if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
69 the[key] = tonumber(x) or x end end
70
71 -- -----
72 -- this code reads csv files where the words on line1 define column types.
73 function ignorep(x) return x:find"%" end -- columns to ignore
74 function klassp(x) return x:find"$" end -- symbolic goals to achieve
75 function lessp(x) return x:find"- $" end -- number goals to minimize
76 function morep(x) return x:find"+ $" end -- numeric goals to maximize
77 function nump(x) return x:find"^[A-Z]" end -- numeric columns
78 function goalp(x) return morep(x) or lessp(x) or klassp(x) end
79
80 -- strings
81 fmt = string.format
82
83 -- maths
84 big = math.huge
85 max = math.max
86 min = math.min
87 r = math.random
88
89 function rnds(t, f) return map(t, function(x) return rnd(x, f) end) end
90 function rnd(x, f)
91 return fmt(type(x)=="number" and (x~x//1 and f or the.Format) or "%s", x) end
92
93 -- tables
94 pop = table.remove
95 unpack = table.unpack
96 function any(t) return t[r[#t]] end
97 function firsts(a, b) return a[1] < b[1] end
98 function many(t, n, u) u={}; for i=1, n do push(u, any(t)) end; return u end
99 function per(t, p) return t[ (#t*(p or .5))//1 ] end
100 function push(t, x) table.insert(t, x); return x end
101 function sort(t, f) table.sort(t, f); return t end
102
103 -- meta
104 function map(t, f, u) u={}; for k, v in pairs(t) do push(u, f(v)) end; return u end
105 function sum(t, f, n) n=0; for _, v in pairs(t) do n=n+f(v) end; return n end
106 function slots(t, u) u={};
107 for k, v in pairs(t) do k=tostring(k); if k:sub(1,1)=="_" then push(u, k) end end
108 return sort(u) end
109
110 -- print tables, recursively
111 function oo(t) print(o(t)) end
112 function o(t)
113 if type(t)=="table" then return tostring(t) end
114 local key=function(k) return fmt("%s %s", k, o(t[k])) end
115 local u = #t>0 and map(t, o) or map(slots(t), key)
116 return '{ ' .. table.concat(u, ", ") .. " }" end
117
118 -- strings to things
119 function csv(file, x)
120 file = io.input(file)
121 return function()
122 x=io.read(); if x then return things(x) else io.close(file) end end end
123
124 function thing(x)
125 x = x:match"^(%s+)(-)%s*$"
126 if x=="true" then return true elseif x=="false" then return false end
127 return tonumber(x) or x end
128
129 function things(x, sep, t)
130 t={}
131 for y in x:gmatch(sep or "([A+])") do push(t, thing(y)) end
132 return t end
133
134 -- misc
135 function distance2Heaven(t, heaven, num, d)
136 for n, txt in pairs(heaven) do
137 num = Num(at, txt)
138 for _, z in pairs(t) do num:add(z.ys[n]) end
139 for _, z in pairs(t) do z.ys[n] = num:distance2Heaven(z.ys[n]) end end
140 d = function(one) return (sum(one.ys) / #one.ys) ^ .5 end
141 return sort(t, function(a, b) return d(a) < d(b) end) end
142
143 -- CLASSES
144
145 function new(k, t) k.__index=k; k.__tostring=o; return setmetatable(t, k) end
146
147 -- COLS: turns list of column names into NUMs, SYMs, or SKIPs
148 function COLS.new(k, row, i)
149 i = new(k, {all={}, x={}, y={}, names=row})
150 for at, txt in ipairs(row) do push(i.all, i:col(at, txt)) end
151 return i end
152
153 function COLS.add(i, t)
154 for _, col in pairs(i.all) do col:add( t[col.at] ) end
155 return t end
156
157 function COLS.col(i, at, txt, col)
158 if ignorep(txt) then return SKIP:new(at, txt) end
159 col = (nump(txt) and NUM or SYM):new(at, txt)
160 push(goalp(txt) and i.y or i.x, col)
161 if klassp(txt) then i.klass = col end
162 return col end
163
164 -- NUM: summarizes a stream of numbers
165 function NUM.new(k, n, s)
166 return new(k, {n=0, at=n or 0, txt=s or "", has=SOME:new(), ok=false,
167 w=lessp(s or "") and -1 or 1, lo=big, hi=-big}) end
168
169 function NUM.add(i, x)
170 if x == "?" then
171 i.n = i.n + 1
172 if i.has:add(x) then i.ok=false end
173 i.lo, i.hi = min(x, i.lo), max(x, i.hi); end end
174
175 function NUM.dist(i, x, y)
176 if x=="?" and y=="?" then return 1
177 elseif x=="?" then y=i:norm(y); x=y<0.5 and 1 or 0
178 elseif y=="?" then x=i:norm(x); x=y<0.5 and 1 or 0
179 x, y = i:norm(x), i:norm(y) end
180 return math.abs(x-y) end
181
182 function NUM.distance2Heaven(x, w)
183 return ((i.w>0 and 1 or 0) - i:norm(x))^2 end
184
185 function NUM.mid(i) return per(i:sorted(), .5) end
186
187 function NUM.norm(i, x)
188 return math.abs(i.hi-i.lo)<1E-9 and 0 or (x-i.lo)/(i.hi - i.lo) end
189
190 function NUM.sorted(i)
191 if i.ok==false then table.sort(i.has.all); i.ok=true end
192 return i.has.all end
193
194 -- ROWS: manages 'rows', summarized in 'cols' (columns).
195 function ROWS.new(k, inits, i)
196 i = new(k, {rows={}, cols=nil})
197 if type(inits)=="string" then for t in csv(inits) do i:add(t) end end
198 if type(inits)=="table" then for t in inits do i:add(t) end end
199 return i end
200
201 function ROWS.add(i, t)
202 if i.cols then push(i.rows, i.cols:add(t)) else i.cols=COLS:new(t) end end
203
204 function ROWS.clone(i, j) j = ROWS:new(); j:add(i.cols.names); return j end
205
206 function ROWS.dist(i, row1, row2, d, fun)
207 function fun(col) return col:dist(row1[col.at], row2[col.at])^the.p end
208 return (sum(i.cols.x, fun) / #i.cols.x)^(1/the.p) end
209
210 function ROWS.furthest(i, row1, rows, fun)
211 function fun(row2) return i:dist(row1, row2) end
212 return unpack(per(sort(map(rows, fun), firsts), the.furthest)) end
213
214 function ROWS.half(i, top)
215 local some, top, c, x, y, tmp, mid, lefts, rights, _
216 some = many(i.rows, the.keep)
217 top = top or i
218 _, x = top:furthest(any(some), some)
219 c, y = top:furthest(x, some)
220 tmp = sort(map(i.rows, function(r) return top:fastmap(r, x, y, c) end), firsts)
221 mid = #i.rows//2
222 lefts, rights = i:clone(), i:clone()
223 for at, row in pairs(tmp) do at <= mid and lefts or rights:add(row[2]) end
224 return lefts, rights, x, y, c, tmp[mid] end
225
226 function ROWS.mid(i, cols)
227 return map(cols or i.cols.all, function(col) return col:mid() end) end
228
229 function ROWS.fastmap(i, r, x, y, c, a, b)
230 a, b = i:dist(r, x), i:dist(r, y); return ((a^2 + c^2 - b^2) / (2*c), r) end
231
232 -- SKIP: summarizes things we want to ignore (so does nothing)
233 function SKIP.new(k, n, s) return new(k, {n=0, at=at or 0, txt=s or ""}) end
234 function SKIP.add(i, x) return x end
235 function SKIP.mid(i) return "?" end
236
237 -- SOME: keeps a random sample on the arriving data
238 function SOME.new(k, keep) return new(k, {n=0, all={}, keep=keep or the.keep}) end
239 function SOME.add(i, x)
240 i.n = i.n + 1
241 if #i.all < i.keep then push(i.all, x) ; return i.all
242 elseif r() < i.keep/i.n then i.all[r(#i.all)]=x; return i.all end end
243
244 -- SYM: summarizes a stream of symbols
245 function SYM.new(k, n, s)
246 return new(k, {n=0, at=n or 0, txt=s or "", has={}, most=0}) end
247
248 function SYM.add(i, x, inc)
249 if x == "?" then
250 inc = inc or 1
251 i.n = i.n + inc
252 i.has[x] = inc + (i.has[x] or 0)
253 if i.has[x] > i.most then i.most, i.mode=i.has[x], x end end end
254
255 function SYM.dist(i, x, y) return (x=="?" and y=="?" and 1) or (x==y and 0 or 1) end
256 function SYM.mid(i) return i.mode end
257 function SYM.div(i, p)
258 return sum(i.has, function(k) p=-i.has[k]/i.n; return -p*math.log(p, 2) end) end
259
260 function SYM.merge(i, j, k)
261 k = SYM:new(i.at, i.txt)
262 for x, n in pairs(i.has) do k:add(x, n)
263 for x, n in pairs(j.has) do k:add(x, n) end
264 ei, ej, ek = i:div(), j:div(), k:div()
265 if i.n==0 or j.n==0 or .99*ek <= (i.n*ei + j.n*ej)/k.n then
266 return k end end
267

```

```

269 -- CLUSTER
270 --
271 --
272 --
273 -- CLUSTER: recursively divides data by clustering towards two distant points
274 function CLUSTER.new(k, eggs, top)
275   local i, want, left, right
276   i = new(k, {here=egs})
277   top = top or egs
278   want = (#top.rows)*the.want
279   if #egs.rows >= 2*want then
280     left, right, i.x, i.y, i.c, i.mid = egs:half(top)
281     if #left.rows < #egs.rows then
282       i.left = CLUSTER.new(left, top)
283       i.right = CLUSTER.new(right, top) end end
284   return i end
285
286 function CLUSTER.show(i, pre, here)
287   pre = pre or ""
288   here=""
289   if not i.left and not i.right then here= o(i.here:mid(i.here.cols.y)) end
290   print(fmt("%6s: %-30s %s", #i.here.rows, pre, here))
291   for _, kid in pairs(i.left, i.right) do
292     if kid then kid:show(pre .. "|. ") end end end
293
294 -- EXPLAIN
295 --
296 --
297 -- SPAN: keeps a random sample on the arriving data
298 function SPAN.new(k, col, lo, hi, has)
299   return new(k, {col=col, lo=lo, hi=hi or lo, has=has or SYM:new()}) end
300
301 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
302 function SPAN.merge(i, j)
303   local has = i.has:merge(j.has)
304   if now then return SPAN:new(i.col, i.lo, j.hi, has) end end
305
306 function SPAN.select(i, row, x)
307   x = row[i.col.at]
308   return (x=="") or (i.lo==i.hi and x==i.lo) or (i.lo <= x and x < i.hi) end
309
310 function SPAN.score(i) return {i.has.n/i.col.n, i.has:div()} end
311
312 -- EXPLAIN:
313 --
314 function EXPLAIN.new(k, egs, top)
315   local i, top, want, left, right, spans, best, yes, no
316   i = new(k, {here = egs})
317   top = top or egs
318   want = (#top.rows)*the.want
319   if #top.rows >= 2*want then
320     left, right = egs:half(top)
321     spans = {}
322     for n, col in pairs(i.cols.x) do
323       for _, s in pairs(col:spans(j.cols.x[n])) do
324         push(spans, {yws=score(j), it=s}) end end
325     best = distance2heaven(spans, {"", "-"})[1]
326     yes, no = egs:clone(), egs:clone()
327     for _, row in pairs(egs.rows) do
328       (best:selects(row) and yes or no):add(row) end -- divide data in two
329     if #yes.rows < #egs.rows then -- make kids if kid size different to parent siz
330       e
331       if #yes.rows >= want then i.yes=EXPLAIN:new(yes, top) end
332       if #no.rows >= want then i.no=EXPLAIN:new(no, top) end end end
333   return i end
334
335 function EXPLAIN.show(i, pre)
336   pre = pre or ""
337   if not pre then
338     tmp = i.here:mid(i.here.y)
339     print(fmt("%6s: %-30s %s", #i.here.rows, pre, o(i.here:mid(i.here.cols.y))))
340     for _, pair in pairs({true, i.yes}, {false, i.no}) do
341       status, kid = unpack(pair)
342       k:shpw(pre .. "|. ") end end end
343
344 function SYM.spans(i, j)
345   local xys, all, one, last, xys, x, c, n = {}, {}
346   for x, n in pairs(i.has) do push(xys, {x, "this", n}) end
347   for x, n in pairs(j.has) do push(xys, {x, "that", n}) end
348   for _, tmp in pairs(sort(xys, firsts)) do
349     x, c, n = unpack(tmp)
350     if x ~= last then
351       last = x
352       one = push(all, Span(i, x, x)) end
353     one:add(x, y, n) end
354   return all end
355
356 function NUM.spans(i, j)
357   local xys, all, lo, hi, gap, xys, one, x, c, n = {}, {}
358   lo, hi = min(i.lo, j.lo), max(i.hi, j.hi)
359   gap = (hi - lo) / (6/the.cohen)
360   for x, n in pairs(i.has) do push(xys, {x, "this", 1}) end
361   for x, n in pairs(j.has) do push(xys, {x, "that", 1}) end
362   one = Span:new(i, lo, lo)
363   all = {one}
364   for _, tmp in pairs(sort(xys, firsts)) do
365     x, c, n = unpack(tmp)
366     if one.hi - one.lo > gap then one = push(all, Span(i, one.hi, x)) end
367     one:add(x, y) end
368   all = merge(all)
369   all[1].j.lo = -big
370   all[#all].hi = big
371   return all end
372
373 function merge(b4, j, n, now, a, b, merged)
374   j, n, now = 0, #b4, {}
375   while j < #b4 do
376     j = j+1
377     a, b = b4[j], b4[j+1]
378     if b then
379       merged = a:merge(b)
380       if merged then a, j = merged, j+1 end end
381     push(now, a)
382     j = j+1 end
383   return #now == #b4 and b4 or merge(now) end

```

```

384 -- DEMOS
385 --
386 --
387 --
388 fails=0
389 function asserts(test, msg)
390   print(test and "PASS: " or "FAIL: ", msg or "")
391   if not test then
392     fails=fails+1
393     if the.dump then assert(test, msg) end end end
394
395 function EGS.nothing() return true end
396 function EGS.the() co(the) end
397 function EGS.rand() print(r()) end
398 function EGS.some(s, t)
399   s=SOME:new(100)
400   for i=1, 100000 do s:add(i) end
401   for j, x in pairs(sort(s.all)) do
402     --if (j % 10) == 0 then print("") end
403     --io.write(fmt("%6s", x)) end end
404     fmt("%6s", x) end end
405
406 function EGS.clone( r, s)
407   r = ROWS:new(the.data)
408   s = r:clone()
409   for _, row in pairs(r.rows) do s:add(row) end
410   asserts(r.cols.x[1].lo==s.cols.x[1].lo, "clone.lo")
411   asserts(r.cols.x[1].hi==s.cols.x[1].hi, "clone.hi")
412   end
413
414 function EGS.data( r)
415   r = ROWS:new(the.data)
416   asserts(r.cols.x[1].hi == 8, "data.columns") end
417
418 function EGS.dist( r, rows, n)
419   r = ROWS:new(the.data)
420   rows = r.rows
421   n = NUM:new()
422   for _, row in pairs(rows) do n:add(r:dist(row, rows[1])) end
423   --oo(r.cols.x[2]:sorted()) end
424   o(r.cols.x[2]:sorted()) end
425
426 function EGS.many( t)
427   t={}; for j=1, 100 do push(t, j) end
428   --print(oo(many(t, 10))) end
429   o(many(t, 10)) end
430
431 function EGS.far( r, c, row1, row2)
432   r = ROWS:new(the.data)
433   row1 = r.rows[1]
434   c, row2 = r:far(r.rows[1], r.rows) end
435   --print(c, "\n", o(row1), "\n", o(row2)) end
436
437 function EGS.half( r, c, row1, row2)
438   local lefts, rights, x, y, x
439   r = ROWS:new(the.data)
440   r:mid(r.cols.y)
441   lefts, rights, x, y, c = r:half()
442   lefts:mid(lefts.cols.y)
443   rights:mid(rights.cols.y)
444   asserts(true, "half") end
445
446 function EGS.cluster(r)
447   r = ROWS:new(the.data)
448   --CLUSTER:new(r):show() end
449   CLUSTER:new(r) end
450
451 -- start-up
452 if arg[0] == "slua" then
453   if the.help then print(help:gsub("\nNOTES:$", "")) else
454     local b4={}; for k, v in pairs(the) do b4[k]=v end
455     for _, todo in pairs(the.todo=="all" and slots(EGS) or {the.todo}) do
456       for k, v in pairs(b4) do the[k]=v end
457       math.randomseed(the.seed)
458       if type(EGS[todo])=="function" then EGS[todo]() end end
459     end
460     for k, v in pairs(_ENV) do if not b4[k] then print("?", k, type(v)) end end
461     os.exit(fails)
462   else
463     return {CLUSTER=CLUSTER, COLS=COLS, NUM=NUM, ROWS=ROWS,
464            SKIP=SKIP, SOME=SOME, SYM=SYM, the=the, oo=oo, o=o}
465   end
466 -- git rid of SOME for rows
467 -- nss = NUM | SYM | SKIP
468 -- COLS = all: [nss]+, x: [nss]*, y: [nss]*, klass: col?
469 -- ROWS = cols: COLS, rows: SOME

```