```lua
1   -- -  Recursively divide data based on two
2   --   distant points (found in linear time using the Fastmap
3   --   heuristic [Fa95]). Then find and print the attribute range
4   --   that best distinguishes these halves. Recurse on each half.
5   -- - (which is sort of like PDDP [Bo98] but faster; and we
6   --   offers a human-readable description for each division).
7   -- - To find those ranges, this code uses a variant of the ChiMerge
8   --   discretizer (but we select on entropy and size,
9   --   not the Chi statistic)
10  -- - To avoid spurious outliers, this code separates using `-furthest=.9`;
11  --   i.e. the 90% furthest points.
12  -- - To avoid long runtimes, this code only searches at most `-keep=512 `
13  --   randomly selected examples to find those furtherst points.
14  -- - To suport multi-objective optimization, this code reads csv files
15  --   whose headers may contain markers for "minimize this" or "maximize
16  --   that" (see the `lessp, morep` functions).
17  -- - To support explanation, optionally, at each level of recursion,
18  --   this code reports what ranges can best distinguish sibling clusters
19  --   C1,C2. The  discretizer is inspired by the ChiMerge algorithm:
20  --   numerics are divided into, say, 16 bins. Then, while we can find
21  --   adjacent bins with the similar distributions in C1,C2, then
22  --   (a) merge then (b) look for other merges.
23  local help = [[
24
25  l5 == a little lab of lots of LUA learning algorithms.
26  (c) 2022, Tim Menzies, BSD 2-clause license.
27
28  USAGE:
29    lua l5.lua [OPTIONS]
30
31  OPTIONS:
32   -cohen    -c   F   Cohen's delta            = .35
33   -data     -d   N   data file                = etc/data/auto93.csv
34   -Dump     -D       stack dump on assert fails = false
35   -furthest -f   F   far                      = .9
36   -Format   -F   S   format string            = %5.2f
37   -keep     -k   P   max kept items           = 512
38   -p        -p   P   distance coefficient     = 2
39   -seed     -s   P   set seed                 = 10019
40   -todo     -t   S   start up action (or 'all') = nothing
41   -help     -h       show help                = false
42   -want     -w   F   recurse until rows^want   = .5
43
44  KEY: N=fileName F=float P=posint S=string
45
46  ]]
47
48  -- ## Definitions
49
50  -- ### Cache current names (used at end to find rogue variables)
51  local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
52
53  -- ### Define locals.
54  local any,asserts,big,cli,csv,fails,firsts,fmt,goalp,ignorep,klassp
55  local lessp,map,main,many,max,merge,min,morep,new,nump,o,oo,per,pop,push
56  local r,rows,rnd,rnds,slots,sort,sum,thing,things,unpack
57
58  -- ### Define classes
59  local CLUSTER, COLS, EGS,  EXPLAIN, NUM, ROWS = {},{},{},{},{},{}
60  local SKIP,    SOME, SPAN, SYM      = {},{},{},{}
61
62  -- ### Define parameter settings.
63  -- Update parameter defaults from command line. Allow for some shorthand:
64  -- e.g.  _-k N_ &rArr; `keep=N`;
65  -- and  _-booleanFlag_ &rArr; `booleanFlag=not default`).
66  local the={}
67  help:gsub("\n [-]([^%s]+)[%s]+(-[^%s]+)[^\n]*%s([^%s]+)",function(key,flag1,x)
68    for n,flag2 in ipairs(arg) do
69      if flag1==flag2 or "-"..key =="flag2"then
70        x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
71    if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
72      the[key] = tonumber(x) or x end end )
73
74  print(the.help)
75
76  -- ### Define headers for row1 of csv files
77
78  -- Columns to ignore
79  function ignorep(x) return x:find":$" end     -- columns to ignore
80  -- Symbolic classes
81  function klassp(x)  return x:find"!$" end      -- symbolic goals to achieve
82  -- Goals to minimize
83  function lessp(x)   return nump(x) and x:find"-$" end     -- number goals to min
    imize
84  --i Goals to mazumze
85  function morep(x)   return x:find"+$" end      -- numeric goals to maximize
86  function nump(x)    return x:find"^[A-Z]" end -- numeric columns
87  function goalp(x)   return morep(x) or lessp(x) or klassp(x) end


88  -- ## Misc Utils
89
90  -- ### Strings
91  fmt = string.format
92
93  -- ### Maths
94  big = math.huge
95  max = math.max
96  min = math.min
97  r   = math.random
98
99  function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
100 function rnd(x,f)
101   return fmt(type(x)=="number" and (x~=x//1 and f or the.Format) or "%s",x) end
102
103 -- ### Tables
104 pop = table.remove
105 unpack = table.unpack
106 function any(t)       return t[r(#t)] end
107 function firsts(a,b)  return a[1] < b[1] end
108 function many(t,n, u) u={}; for i=1,n do push(u,any(t)) end; return u end
109 function per(t,p)     return t[ (#t*(p or .5))//1 ] end
110 function push(t,x)    table.insert(t,x); return x end
111 function sort(t,f)    table.sort(t,f); return t end
112
113 -- ### Meta
114 function map(t,f, u)  u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
115 function sum(t,f, n)  n=0; for _,v in pairs(t) do n=n+f(v)     end; return n end
116 function slots(t, u)
117   u={}
118   for k,v in pairs(t) do k=tostring(k);if k:sub(1,1)~="_" then push(u,k) end end
119   return sort(u) end
120
121 -- ### Print tables, recursively
122 function oo(t) print(o(t)) end
123 function o(t)
124   if type(t)~="table" then return tostring(t) end
125   local key=function(k) return fmt(":%s %s",k,o(t[k])) end
126   local u = #t>0 and map(t,o) or map(slots(t),key)
127   return '{'..table.concat(u,"")..'}' end
128
129 -- ### Coerce strings to things
130 function csv(file,     x)
131   file = io.input(file)
132   return function()
133     x=io.read(); if x then return things(x) else io.close(file) end end end
134
135 function thing(x)
136   x = x:match"^%s*(.-)%s*$"
137   if x=="true" then return true elseif x=="false" then return false end
138   return tonumber(x) or x end
139
140 function things(x,sep,  t)
141   t={}
142   for y in x:gmatch(sep or"([^,]+)") do push(t,thing(y)) end
143   return t end
144
145 -- ### Misc stuff
146 function distance2Heaven(t,heaven,   num,d)
147   for n,txt in pairs(heaven) do
148     num = Num(at,txt)
149     for _,z in pairs(t) do num:add(z.ys[n]) end
150     for _,z in pairs(t) do z.ys[n] = num:distance2heaven(z.ys[n]) end end
151   d = function(one) return (sum(one.ys)/#one.ys)^.5 end
152   return sort(t, function(a,b) return d(a) < d(b) end) end
153
154 -- objects
155 function new(k,t) k.__index=k; k.__tostring=o; return setmetatable(t,k) end
```

```lua
156  -- ## COLS
157
158  -- Factory. Turns list of column names into NUMs, SYMs, or SKIPs
159  function COLS.new(k,row,    i,create1)
160    create1 = function(i,at,txt,    col)
161      if ignorep(txt) then return SKIP:new(at,txt) end
162      col = (nump(txt) and NUM or SYM):new(at,txt)
163      push(goalp(txt) and i.y or i.x, col)
164      if klassp(txt) then i.klass = col end
165      return col
166    end ------------------
167    i= new(k,{all={},x={},y={},names=row})
168    for at,txt in ipairs(row) do  push(i.all, create1(at,txt)) end
169    return i end
170
171  function COLS.add(i,t)
172    for _,col in pairs(i.all) do col:add( t[col.at] ) end
173    return t end
174
175    -- NUM: summarizes a stream of numbers
176  function NUM.new(k,n,s)
177    return new(k,{n=0,at=n or 0,txt=s or"",has=SOME:new(),ok=false,
178                   w=lessp(s or "") and -1 or 1, lo=big, hi=-big}) end
179
180  function NUM.add(i,x)
181    if x ~= "?" then
182      i.n = i.n + 1
183      if i.has:add(x) then i.ok=false end
184      i.lo,i.hi = min(x,i.lo), max(x,i.hi); end end
185
186  function NUM.dist(i,x,y)
187    if     x=="?" and y=="?" then return 1
188    elseif x=="?" then y=i:norm(y); x=y<0.5 and 1 or 0
189    elseif y=="?" then x=i:norm(x); y=x<0.5 and 1 or 0
190    else   x,y = i:norm(x), i:norm(y) end
191    return math.abs(x-y) end
192
193  function NUM.distance2heaven(x, w)
194    return ((i.w>0 and 1 or 0) - i:norm(x))^2 end
195
196  function NUM.mid(i) return per(i:sorted(), .5) end
197
198  function NUM.norm(i,x)
199    return math.abs(i.hi-i.lo)<1E-9 and 0 or (x-i.lo)/(i.hi - i.lo) end
200
201  function NUM.sorted(i)
202    if i.ok==false then table.sort(i.has.all); i.ok=true end
203    return i.has.all end
204
205  -- ROWS: manages `rows`, summarized in `cols` (columns).
206  function ROWS.new(k,inits,    i)
207    i = new(k,{rows={},cols=nil})
208    if type(inits)=="string" then for t in csv(inits) do i:add(t) end end
209    if type(inits)=="table"  then for t in inits       do i:add(t) end end
210    return i end
211
212  function ROWS.add(i,t)
213    if i.cols then push(i.rows,i.cols:add(t)) else i.cols=COLS:new(t) end end
214
215  function ROWS.clone(i,  j) j= ROWS:new(); j:add(i.cols.names);return j end
216
217  function ROWS.dist(i,row1,row2,    d,fun)
218    function fun(col)  return col:dist(row1[col.at], row2[col.at])^the.p end
219    return (sum(i.cols.x, fun)/ #i.cols.x)^(1/the.p) end
220
221  function ROWS.furthest(i,row1,rows,     fun)
222    function fun(row2)  return {i:dist(row1,row2), row2} end
223    return unpack(per(sort(map(rows,fun),firsts), the.furthest))  end
224
225  function ROWS.half(i, top)
226    local some, top,c,x,y,tmp,mid,lefts,rights,_
227    some= many(i.rows, the.keep)
228    top = top or i
229    _,x = top:furthest(any(some), some)
230    c,y = top:furthest(x,         some)
231    tmp = sort(map(i.rows,function(r) return top:fastmap(r,x,y,c) end),firsts)
232    mid = #i.rows//2
233    lefts, rights = i:clone(), i:clone()
234    for at,row in pairs(tmp) do (at <=mid and lefts or rights):add(row[2]) end
235    return lefts,rights,x,y,c, tmp[mid] end
236
237  function ROWS.mid(i,cols)
238    return map(cols or i.cols.all, function(col) return col:mid() end) end
239
240  function ROWS.fastmap(i, r,x,y,c,       a,b)
241    a,b = i:dist(r,x), i:dist(r,y); return {(a^2 + c^2 - b^2)/(2*c), r} end
242
243  -- SKIP: summarizes things we want to ignore (so does nothing)
244  function SKIP.new(k,n,s) return new(k,{n=0,at=at or 0,txt=s or""}) end
245  function SKIP.add(i,x)       return x end
246  function SKIP.mid(i)       return "?" end
247
248  -- SOME: keeps a random sample on the arriving data
249  function SOME.new(k,keep) return new(k,{n=0,all={}, keep=keep or the.keep}) end
250  function SOME.add(i,x)
251    i.n = i.n+1
252    if     #i.all < i.keep then push(i.all,x)               ; return i.all
253    elseif r()      < i.keep/i.n then i.all[r(#i.all)]=x; return i.all end end
254
255  -- SYM: summarizes a stream of symbols
256  function SYM.new(k,n,s)
257    return new(k,{n=0,at=n or 0,txt=s or"",has={},most=0}) end
258
259  function SYM.add(i,x,inc)
260    if x ~= "?" then
261      inc = inc or 1
262      i.n = i.n + inc
263      i.has[x] = inc + (i.has[x] or 0)
264      if i.has[x] > i.most then i.most,i.mode=i.has[x],x end end end
265
266  function SYM.dist(i,x,y) return(x=="?" and y=="?" and 1) or(x==y and 0 or 1) end
267  function SYM.mid(i)         return i.mode end
268  function SYM.div(i,   p)
269    return sum(i.has,function(k) p=-i.has[k]/i.n;return -p*math.log(p,2) end) end
270
271  function SYM.merge(i,j,    k)
272    k = SYM:new(i.at,i.txt)
273    for x,n in pairs(i.has) do k:add(x,n) end
274    for x,n in pairs(j.has) do k:add(x,n) end
275    ei, ej, ejk= i:div(), j:div(), k:div()
276    if i.n==0 or j.n==0 or .99*ek <= (i.n*ei + j.n*ej)/k.n then
277      return k end end

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

```lua
--    ___    ___   __   __    ___
--   |   \  |_  | |  \ |  | |   _
--   |__ / |__  | |  | |  | |___ |
--
fails=0
function asserts(test, msg)
  print(test and "PASS:"or "FAIL: ",msg or "")
  if not test then
     fails=fails+1
     if the.dump then assert(test,msg) end end end

function EGS.nothing() return true end
function EGS.the()      oo(the) end
function EGS.rand()     print(r()) end
function EGS.some(s,t)
  s=SOME:new(100)
  for i=1,100000 do s:add(i) end
  for j,x in pairs(sort(s.all)) do
    --if (j % 10)==0 then print("") end
    --io.write(fmt("%6s",x))   end end
    fmt("%6s",x)   end end

function EGS.clone( r,s)
  r = ROWS:new(the.data)
  s = r:clone()
  for _,row in pairs(r.rows) do s:add(row) end
  asserts(r.cols.x[1].lo==s.cols.x[1].lo,"clone.lo")
  asserts(r.cols.x[1].hi==s.cols.x[1].hi,"clone.hi")
  end

function EGS.data( r)
  r = ROWS:new(the.data)
  asserts(r.cols.x[1].hi == 8, "data.columns") end

function EGS.dist( r,rows,n)
  r = ROWS:new(the.data)
  rows = r.rows
  n = NUM:new()
  for _,row in pairs(rows) do n:add(r:dist(row, rows[1])) end
  --oo(r.cols.x[2]:sorted()) end
  o(r.cols.x[2]:sorted()) end

function EGS.many(   t)
  t={}; for j=1,100 do push(t,j) end
  --print(oo(many(t, 10))) end
  o(many(t, 10)) end

function EGS.far(   r,c,row1,row2)
  r = ROWS:new(the.data)
  row1   = r.rows[1]
  c,row2 = r:far(r.rows[1], r.rows) end
  --print(c,"\n",o(row1),"\n", o(row2)) end

function EGS.half(   r,c,row1,row2)
  local lefts,rights,x,y,x
  r = ROWS:new(the.data)
  r:mid(r.cols.y)
  lefts,rights,x,y,c = r:half()
  lefts:mid(lefts.cols.y )
  rights:mid(rights.cols.y)
  asserts(true,"half") end

function EGS.cluster(r)
  r = ROWS:new(the.data)
  --CLUSTER:new(r):show() end
  CLUSTER:new(r) end

-- start-up
if arg[0] == "sl.lua" then
  if the.help then print(help:gsub("\nNOTES:*$","")) else
    local b4={}; for k,v in pairs(the) do b4[k]=v end
    for _,todo in pairs(the.todo=="all" and slots(EGS) or {the.todo}) do
      for k,v in pairs(b4) do the[k]=v end
      math.randomseed(the.seed)
      if type(EGS[todo])=="function" then EGS[todo]() end end
  end
  for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
  os.exit(fails)
else
  return {CLUSTER=CLUSTER, COLS=COLS, NUM=NUM, ROWS=ROWS,
          SKIP=SKIP, SOME=SOME, SYM=SYM,the=the,oo=oo,o=o}
end
-- git rid of SOME for rows
-- nss  = NUM | SYM | SKIP
-- COLS = all:[nss]+, x:[nss]*, y:[nss]*, klass;col?
-- ROWS = cols:COLS, rows:SOME
--
-- [Ah91]: Aha, D.W., Kibler, D. & Albert, M.K. Instance-based   learning algori
thms. Mach Learn 6, 37â€”^@M-^S66 (1991).  https://doi.org/10.1007/BF00153759
--
```