```lua
   1  --- --------------------------------------------------------------------------
   2  ---
   3  ---          /\     /\
   4  ---         //\\ ___ //\\
   5  ---         \  '    ' '`
   6  ---          \ \ \L\ \  \/\ \L\ \
   7  ---           \ \____/   \  \____/
   8  ---            \/___/     \/___/
   9  ---
  10  --- --------------------------------------------------------------------------
  11  -- - Recursively divide data based on two
  12  --   distant points (found in linear time using the Fastmap
  13  --   heuristic [Fa95]). Then find and print the attribute range
  14  --   that best distinguishes these halves. Recurse on each half.
  15  -- - (which is sort of like PDDP [Bo98] but faster; and we
  16  --   offers a human-readable description for each division).
  17  -- - To find those ranges, this code uses a variant of the ChiMerge
  18  --   discretizer (but we select on entropy and size,
  19  --   not the Chi statistic)
  20  -- - To avoid spurious outliers, this code separates using `-furthest=.9`;
  21  --   i.e. the 90% furthest points.
  22  -- - To avoid long runtimes, this code only searches at most `-keep=512 `
  23  --   randomly selected examples to find those furtherst points.
  24  -- - To suport multi-objective optimization, this code reads csv files
  25  --   whose headers may contain markers for "minimize this" or "maximize
  26  --   that" (see the `lessp, morep` functions)
  27  -- - To support explanation, optionally, at each level of recursion,
  28  --   this code reports what ranges can best distinguish sibling clusters
  29  --   C1,C2.  The  discretizer is inspired by the ChiMerge algorithm:
  30  --   numerics are divided into, say, 16 bins. Then, while we can find
  31  --   adjacent bins with the similar distributions in C1,C2, then
  32  --   (a) merge then (b) look for other merges.
  33  local help = [[
  34
  35  l5 == a little LUA learning library
  36  (c) 2022, Tim Menzies, BSD 2-clause license.
  37
  38  USAGE:
  39    lua l5.lua [OPTIONS]
  40
  41  OPTIONS:
  42    -cohen    -c   F   Cohen's delta            = .35
  43    -data     -d   N   data file                = ../etc/data/auto93.csv
  44    -Dump     -D       stack dump on assert fails = false
  45    -furthest -f   F   far                      = .9
  46    -Format   -F   S   format string            = %5.2f
  47    -keep     -k   P   max kept items           = 512
  48    -p        -p   P   distance coefficient     = 2
  49    -seed     -s   P   set seed                 = 10019
  50    -todo     -t   S   start up action (or 'all') = nothing
  51    -help     -h       show help                = false
  52    -want     -w   F   recurse until rows^want   = .5
  53
  54  KEY: N=fileName F=float P=posint S=string
  55  ]]
  56
  57  -- ## Definitions
  58
  59  -- Cache current names (used at end to find rogue variables)
  60  local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
  61
  62  -- Define locals.
  63  local any,asserts,big,cli,distance2Heaven
  64  local fails,firsts,fmt,goalp,ignorep,klassp
  65  local lessp,map,main,many,max,merge,min,morep,new,nump,o,oo,per,pop,push
  66  local r,rows,rnd,rnds,slots,sort,sum,thing,things,file2things,unpack
  67
  68  -- Define classes
  69  local CLUSTER, COLS, EGS,  EXPLAIN, NUM, ROWS = {},{},{},{},{},{}
  70  local SKIP,    SOME, SPAN, SYM       = {},{},{},{}
  71
  72  -- Define parameter settings.
  73  -- Update parameter defaults from command line. Allow for some shorthand.
  74  -- e.g.  _-k N_ &rArr; `keep=N`;
  75  -- and  _-booleanFlag_ &rArr; `booleanFlag=not default`).
  76  local the={}
  77  help:gsub("\n ([-]|([^%s]+))[%s]+(-[^%s]+)[^\n]*%s([^%s]+)",
  78    function(long,key,short,x)
  79      for n,flag in ipairs(arg) do
  80        if flag==short or flag==long then
  81          x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
  82      if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
  83        the[key] = tonumber(x) or x end end )
  84
  85  -- ### Define headers for row1 of csv files
  86
  87  -- Columns to ignore
  88  function ignorep(x) return x:find":$" end
  89  -- Symbolic class columns.
  90  function klassp(x)  return not nump(x) and x:find"!$" end
  91  -- Goal columns to minimize
  92  function lessp(x)   return nump(x) and x:find"-$" end
  93  -- Goal columns to maximize
  94  function morep(x)   return nump(x) and x:find"+$" end
  95  -- Numeric columns
  96  function nump(x)    return x:find"^[A-Z]" end
  97  -- Dependent columns
  98  function goalp(x)   return morep(x) or lessp(x) or klassp(x) end

  99  ---
 100  ---
 101  ---          ___
 102  ---          Functions
 103  ---
 104
 105  -- ## Misc Utils
 106
 107  -- Strings
 108  fmt = string.format
 109
 110  -- Maths
 111  big = math.huge
 112  max = math.max
 113  min = math.min
 114  r   = math.random
 115
 116  function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
 117  function rnd(x,f)
 118    return fmt(type(x)=="number" and (x~=x//1 and f or the.Format) or "%s",x) end
 119
 120  -- Tables
 121  pop = table.remove
 122  unpack = table.unpack
 123  function any(t)         return t[r(#t)] end
 124  function firsts(a,b)    return a[1] < b[1] end
 125  function many(t,n, u)   u={}; for i=1,n do push(u,any(t)) end; return u end
 126  function per(t,p)       return t[ (#t*(p or .5))//1 ] end
 127  function push(t,x)      table.insert(t,x); return x end
 128  function sort(t,f)      table.sort(t,f); return t end
 129
 130  -- Meta
 131  function map(t,f, u)  u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
 132  function sum(t,f, n)  n=0; for _,v in pairs(t) do n=n+f(v)    end; return n end
 133  function slots(t, u)
 134    u={}
 135    for k,v in pairs(t) do k=tostring(k);if k:sub(1,1)~="_" then push(u,k) end end
 136    return sort(u) end
 137
 138  -- Print tables, recursively
 139  function oo(t) print(o(t)) end
 140  function o(t)
 141    if type(t)~="table" then return tostring(t) end
 142    local key=function(k) return fmt(":%s %s",k,o(t[k])) end
 143    local u = #t>0 and map(t,o) or map(slots(t),key)
 144    return '{'..table.concat(u," ")..'}' end
 145
 146  -- Coerce strings to things
 147  function thing(x)
 148    x = x:match"^%s*(.-)%s*$"
 149    if x=="true" then return true elseif x=="false" then return false end
 150    return tonumber(x) or x end
 151
 152  function things(x,sep,  t)
 153    t={}; for y in x:gmatch(sep or"([^.]+)") do push(t,thing(y)) end
 154    return t end
 155
 156  function file2things(file,      x)
 157    file = io.input(file)
 158    return function()
 159      x=io.read(); if x then return things(x) else io.close(file) end end end
 160
 161  -- ### Misc stuff
 162
 163  -- Multi-objectives. Normalized, scored  via distance to heaven.
 164  function distance2Heaven(t,heaven,    num,d)
 165    for n,txt in pairs(heaven) do
 166      num = Num(at,txt)
 167      for _,z in pairs(t) do num:add(z.ys[n]) end
 168      for _,z in pairs(t) do z.ys[n] = num:distance2heaven(z.ys[n]) end end
 169    d = function(one) return (sum(one.ys)/#one.ys)^.5 end
 170    return sort(t, function(a,b) return d(a) < d(b) end) end
 171
 172  -- While merges found: merge similar adjacent ranges j and j+1 then jump to j+2.
 173  function merge(b4,       j,n,now,a,b,merged)
 174    j,n,now = 0,#b4,{}
 175    while j < #b4 do
 176      j    = j+1
 177      a, b = b4[j], b4[j+1]
 178      if b then
 179        merged = a:merge(b)
 180        if merged then a,j = merged, j+1 end end
 181      push(now,a)
 182      j = j+1 end
 183    return #now == #b4 and b4 or merge(now) end
 184
 185  -- Objects
 186  function new(k,t) k.__index=k; k.__tostring=o; return setmetatable(t,k) end
```

```
187 ---
188 ---     ___    _
189 ---    / _|   | |    ___     ___   ___   ___   ___
190 ---   | (_    | |   / _` |  (_-<  (_-<  / -_)  (_-<
191 ---    \__|   _|_   \__,_|  /__/_ /__/_ \___| /__/_
192
193 ---
194 ---    ____    ___    _       ___
195 ---   / __|   / _ \  | |     / __|
196
197 -- ## COLS
198 -- Factory. Turns list of column names into NUMs, SYMs, or SKIPs
199 function COLS.new(k,row,    i,create1)
200   create1 = function(at,txt,    col)
201     if ignorep(txt) then return SKIP:new(at,txt) end
202     col = (nump(txt) and NUM or SYM):new(at,txt)
203     push(goalp(txt) and i.y or i.x, col)
204     if klassp(txt) then i.klass = col end
205     return col
206   end ---------------------------------
207   i= new(k,{all={},x={},y={},names=row})
208   for at,txt in ipairs(row) do  push(i.all, create1(at,txt)) end
209   return i end
210
211 function COLS.add(i,t)
212   for _,col in pairs(i.all) do col:add( t[col.at] ) end
213   return t end
214 ---
215 ---    _  _   _   _   __  __
216 ---   | \| | | | | | |  \/  |
217
218 -- NUM: summarizes a stream of numbers
219 function NUM.new(k,n,s)
220   return new(k,{n=0,at=n or 0,txt=s or"",has=SOME:new(),ok=false,
221             w=lessp(s or "") and -1 or 1, lo=big, hi=-big}) end
222
223 function NUM.add(i,x)
224   if x ~= "?" then
225     i.n = i.n + 1
226     if i.has:add(x) then i.ok=false end
227     i.lo,i.hi = min(x,i.lo), max(x,i.hi); end end
228
229 function NUM.dist(i,x,y)
230   if     x=="?" and y=="?" then return 1
231   elseif x=="?" then y=i:norm(y); x=y<0.5 and 1 or 0
232   elseif y=="?" then x=i:norm(x); y=x<0.5 and 1 or 0
233   else   x,y = i:norm(x), i:norm(y) end
234   return math.abs(x-y) end
235
236 function NUM.distance2heaven(x, w)
237   return ((i.w>0 and 1 or 0) - i:norm(x))^2 end
238
239 function NUM.mid(i) return per(i:sorted(), .5) end
240
241 function NUM.norm(i,x)
242   return math.abs(i.hi-i.lo)<1E-9 and 0 or (x-i.lo)/(i.hi - i.lo) end
243
244 function NUM.sorted(i)
245   if i.ok==false then table.sort(i.has.all); i.ok=true end
246   return i.has.all end
247 ---
248 ---    ___    ___   __      __  ___
249 ---   | _ \  / _ \  \ \    / / / __|
250
251 -- ROWS: manages `rows`, summarized in `cols` (columns).
252 function ROWS.new(k,inits,    i)
253   i = new(k,{rows={},cols=nil})
254   if type(inits)=="table" then for t in inits do i:add(t) end end
255   if type(inits)=="string" then for t in file2things(inits) do i:add(t) end end
256   return i end
257
258 function ROWS.add(i,t)
259   if i.cols then push(i.rows,i.cols:add(t)) else i.cols=COLS:new(t) end end
260
261 function ROWS.clone(i,  j) j= ROWS:new(); j:add(i.cols.names);return j end
262
263 function ROWS.dist(i,row1,row2,    d,fun)
264   function fun(col) return col:dist(row1[col.at], row2[col.at])^the.p end
265   return (sum(i.cols.x, fun)/ #i.cols.x)^(1/the.p) end
266
267 function ROWS.furthest(i,row1,rows,    fun)
268   function fun(row2) return {i:dist(row1,row2), row2} end
269   return unpack(per(sort(map(rows,fun),firsts), the.furthest)) end
270
271 function ROWS.half(i, top)
272   local some, top,c,x,y,tmp,mid,lefts,rights,_
273   some= many(i.rows, the.keep)
274   top = top or i
275   _,x = top:furthest(any(some), some)
276   c,y = top:furthest(x,         some)
277   tmp = sort(map(i.rows,function(r) return top:fastmap(r,x,y,c) end),firsts)
278   mid = #i.rows//2
279   lefts, rights = i:clone(), i:clone()
280   for at,row in pairs(tmp) do (at <=mid and lefts or rights):add(row[2]) end
281   return lefts,rights,x,y,c, tmp[mid] end
282
283 function ROWS.mid(i,cols)
284   return map(cols or i.cols.all, function(col) return col:mid() end) end
285
286 function ROWS.fastmap(i, r,x,y,c,    a,b)
287   a,b = i:dist(r,x), i:dist(r,y); return {(a^2 + c^2 - b^2)/(2*c), r} end
288 ---
289 ---    ___   _  __  ___   ___
290 ---   / __| | |/ / |_ _| | _ \
291
292 -- SKIP: summarizes things we want to ignore (so does nothing)
293 function SKIP.new(k,n,s) return new(k,{n=0,at=at or 0,txt=s or""}) end
294 function SKIP.add(i,x)     return x end
295 function SKIP.mid(i)       return "?" end
296 ---
297 ---    ___    ___    __  __   ___
298 ---   / __|  / _ \  |  \/  | | __|
299
300 -- SOME: keeps a random sample on the arriving data
301 function SOME.new(k,keep) return new(k,{n=0,all={}, keep=keep or the.keep}) end
302 function SOME.add(i,x)
303   i.n = i.n+1
304   if     #i.all < i.keep then push(i.all,x)                ; return i.all
305   elseif r()     < i.keep/i.n then i.all[r(#i.all)]=x; return i.all end end
306 ---
307 ---    ___   _   _   __  __
308 ---   / __| | | | | |  \/  |
309
310 -- SYM: summarizes a stream of symbols
311 function SYM.new(k,n,s)
312   return new(k,{n=0,at=n or 0,txt=s or"",has={},most=0}) end
313
314 function SYM.add(i,x,inc)
315   if x ~= "?" then
316     inc = inc or 1
317     i.n = i.n + inc
318     i.has[x] = inc + (i.has[x] or 0)
319     if i.has[x] > i.most then i.most,i.mode=i.has[x],x end end end
320
321 function SYM.dist(i,x,y) return(x=="?" and y=="?" and 1) or(x==y and 0 or 1) end
322 function SYM.mid(i)      return i.mode end
```

```
323 function SYM.div(i,   p)
324   return sum(i.has,function(k) p=-i.has[k]/i.n;return -p*math.log(p,2) end) end
325
326 function SYM.merge(i,j,   k)
327   k = SYM:new(i.at,i.txt)
328   for x,n in pairs(i.has) do k:add(x,n) end
329   for x,n in pairs(j.has) do k:add(x,n) end
330   ei, ej, ejk= i:div(), j:div(), k:div()
331   if i.n==0 or j.n==0 or .99*ek <= (i.n*ei + j.n*ej)/k.n then
332     return k end end
333 ---
334 ---    ___   _       _   _   ___   _____   ___   ___
335 ---   / __| | |     | | | | / __| |_   _| | __| | _ \
336
337 -- CLUSTER: recursively divides data by clustering towards two distant points
338 function CLUSTER.new(k,egs,top)
339   local i,want,left,right
340   i     = new(k, {here=egs})
341   top   = top or egs
342   want  = (#top.rows)^the.want
343   if #egs.rows >= 2*want then
344     left, right, i.x, i.y, i.c, i.mid = egs:half(top)
345     if #left.rows < #egs.rows then
346       i.left = CLUSTER:new(left,   top)
347       i.right= CLUSTER:new(right, top) end end
348   return i end
349
350 function CLUSTER.show(i,pre,  here)
351   pre = pre or ""
352   here=""
353   if not i.left and not i.right then here= o(i.here:mid(i.here.cols.y)) end
354   print(fmt("%6s: %-30s %s",#i.here.rows, pre, here))
355   for _,kid in pairs{i.left, i.right} do
356     if kid then kid:show(pre .. "|.. ") end end end
357 ---
358 ---    ___   ___    _    _  _
359 ---   / __| | _ \  /_\  | \| |
360
361 -- SPAN: keeps a random sample on the arriving data
362 function SPAN.new(k, col, lo, hi, has)
363   return new(k,{col=col,lo=lo,hi=hi or lo,has=has or SYM:new()}) end
364
365 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
366 function SPAN.merge(i,j)
367   local has = i.has:merge(j.has)
368   if now then return SPAN:new(i.col, i.lo, j.hi, has) end end
369
370 function SPAN.select(i,row,    x)
371   x = row[i.col.at]
372   return (x=="?") or (i.lo==i.hi and x==i.lo) or (i.lo <= x and x < i.hi) end
373
374 function SPAN.score(i) return {i.has.n/i.col.n,  i.has:div()} end
375 ---
376 ---    ___  __  __  ___    _     _    ___   _  _
377 ---   | __| \ \/ / | _ \  | |   /_\  |_ _| | \| |
378
379 -- ### EXPLAIN:
380 function EXPLAIN.new(k,egs,top)
381   local i,top,want,left,right,spans,best,yes,no
382   i     = new(k,{here = egs})
383   top   = top or egs
384   want  = (#top.rows)^the.want
385   if #top.rows >= 2*want then
386     left = egs:half(top)
387     spans = {}
388     for n,col in pairs(i.cols.x) do
389       for _,s in pairs(col:spans(j.cols.x[n])) do
390         push(spans,{ys=s:score(),it=s}) end end
391     best  = distance2heaven(spans,{"+","-"})[1]
392     yes,no = egs:clone(), egs:clone()
393     for _,row in pairs(egs.rows) do
394       (best:selects(row) and yes or no):add(row) end -- divide data in two
395     if #yes.rows<#egs.rows then -- make kids if kid size different to parent size
396       if #yes.rows>=want then i.yes=EXPLAIN:new(yes,top) end
397       if #no.rows >=want then i.no =EXPLAIN:new(no, top)  end end end
398   return i end
399
400 function EXPLAIN.show(i,pre)
401   pre = pre or ""
402   if not pre then
403     tmp = i.here:mid(i.here.y)
404     print(fmt("%6s: %-30s %s", #i.here.rows, pre, o(i.here:mid(i.here.cols.y))))
405   for _,pair in pairs{{true,i.yes},{false,i.no}} do
406     status,kid = unpack(pair)
407     k:shpw(pre .. "|.. ") end end end
408 ---
409 ---    ___   ___    _    _  _   ___
410 ---   / __| | _ \  /_\  | \| | / __|
411
412 function SYM.spans(i, j)
413   local xys,all,one,last,xys,x,c n = {},{}
414   for x,n in pairs(i.has) do push(xys, {x,"this",n}) end
415   for x,n in pairs(j.has) do push(xys, {x,"that",n}) end
416   for _,tmp in ipairs(sort(xys,firsts)) do
417     x,c,n = unpack(tmp)
418     if x ~= last then
419       last = x
420       one  = push(all, Span(i,x,x)) end
421     one:add(x,y,n) end
422   return all end
423
424 function NUM.spans(i, j)
425   local xys,all,lo,hi,gap,one,x,c,n = {},{}
426   lo,hi = min(i.lo, j.lo), max(i.hi,j.hi)
427   gap   = (hi - lo) / (6/the.cohen)
428   for x,n in pairs(i.has.all) do push(xys, {x,"this",1}) end
429   for x,n in pairs(j.has.all) do push(xys, {x,"that",1}) end
430   one = SPAN:new(i,lo,lo)
431   all = {one}
432   oo(xys)
433   for _,tmp in ipairs(sort(xys,firsts)) do
434     x,c,n = unpack(tmp)
435     if one.hi - one.lo > gap then one = push(all, SPAN(i, one.hi, x)) end
436     one:add(x,y) end
437   all         = merge(all)
438   all[1   ].lo = -big
439   all[#all].hi =  big
440   return all end
```

```
441  ---
442  ---     _            _
443  ---    | |_ __ _ _ _| |_ _  _ _ __
444  ---    (_-< _/ _` | '_|  _| || | '_ \
445  ---    |_/__,_|_|  \__|\_,_| .__/
446  ---                        |_|
```

```lua
447  fails=0
448  function asserts(test, msg)
449    print(test and "PASS:"or "FAIL: ",msg or "")
450    if not test then
451      fails=fails+1
452      if the.dump then assert(test,msg) end end end
453
454  function EGS.nothing() return true end
455  function EGS.the()      oo(the) end
456  function EGS.rand()     print(r()) end
457  function EGS.some(s,t)
458    s=SOME:new(100)
459    for i=1,100000 do s:add(i) end
460    asserts(100==#s.all,"length")
461    for j,x in pairs(sort(s.all)) do
462      --if (j % 10)==0 then print("") end
463      --io.write(fmt("%6s",x))  end end
464      fmt("%6s",x)   end end
465
466  function EGS.clone( r,s)
467    r = ROWS:new(the.data)
468    s = r:clone()
469    for _,row in pairs(r.rows) do s:add(row) end
470    asserts(r.cols.x[1].lo==s.cols.x[1].lo,"clone.lo")
471    asserts(r.cols.x[1].hi==s.cols.x[1].hi,"clone.hi")
472    end
473
474  function EGS.data( r)
475    r = ROWS:new(the.data)
476    asserts(r.cols.x[1].hi == 8, "data.columns") end
477
478  function EGS.dist( r,rows,n)
479    r = ROWS:new(the.data)
480    rows = r.rows
481    n = NUM:new()
482    for _,row in pairs(rows) do n:add(r:dist(row, rows[1])) end
483    oo(rnds(n:sorted()))
484    --oo(r.cols.x[2]:sorted())
485    o(r.cols.x[2]:sorted()) end
486
487  function EGS.many(   t)
488    t={}; for j=1,1000 do push(t,j) end
489    --print(oo(many(t, 10))) end
490    oo(many(t, 10)) end
491
492  function EGS.far(   r,c,row1,row2)
493    r = ROWS:new(the.data)
494    row1  = r.rows[1]
495    c,row2 = r:far(r.rows[1], r.rows) end
496    --print(c,"\n",o(row1),"\n", o(row2)) end
497
498  function EGS.half(   r,c,row1,row2)
499    local lefts,rights,x,y,x
500    r = ROWS:new(the.data)
501    r:mid(r.cols.y)
502    lefts,rights,x,y,c = r:half()
503    lefts:mid(lefts.cols.y )
504    rights:mid(rights.cols.y)
505    asserts(199==#lefts.rows,"left rows")
506    asserts(199==#rights.rows,"right rows")
507    asserts(true,"half") end
508
509  function EGS.cluster(r)
510    r = ROWS:new(the.data)
511    --CLUSTER:new(r):show() end
512    CLUSTER:new(r):show() end
513
514  function EGS.numspan(   r,c,row1,row2)
515    local lefts,rights,x,y,x
516    r = ROWS:new(the.data)
517    r:mid(r.cols.y)
518    lefts,rights,x,y,c = r:half()
519    lefts.cols.x[1]:spans(rights.cols.x[1]) end
520
521  -- start-up
522  if arg[0] == "l5.lua" then
523    if the.help then print(help) else
524      local b4={}; for k,v in pairs(the) do b4[k]=v end
525      for _,todo in pairs(the.todo=="all" and slots(EGS) or {the.todo}) do
526        for k,v in pairs(b4) do the[k]=v end
527        math.randomseed(the.seed)
528        if type(EGS[todo])=="function" then EGS[todo]() end end
529    end
530    for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
531    os.exit(fails)
532  else
533    return {CLUSTER=CLUSTER, COLS=COLS, NUM=NUM, ROWS=ROWS,
534            SKIP=SKIP, SOME=SOME, SYM=SYM,the=the,oo=oo,o=o}
535  end
536  -- git rid of SOME for rows
537  -- nss  = NUM | SYM | SKIP
538  -- COLS = all:[nss]+, x:[nss]*, y:[nss]*, klass;col?
539  -- ROWS = cols:COLS, rows:SOME
540  --
541  -- [Ah91]: Aha, D.W., Kibler, D. & Albert, M.K. Instance-based   learning algori
      thms. Mach Learn 6, 37â M-^@M-^S66 (1991).  https://doi.org/10.1007/BF00153759
542  --
```