```lua
  1  local help = [[
  2
  3  l5 == a little lab of lots of LUA learning algorithms.
  4  (c) 2022, Tim Menzies, BSD 2-clause license.
  5
  6  USAGE:
  7    lua l5.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  ]]
 25  -- NOTES: This code uses Aha's distance measure [^Ah91] (that can
 26  -- handle numbers and symbols) to recursively divide data based on two
 27  -- distant points (these two are found in linear time using the Fastmap
 28  -- heuristic [Fa95]).
 29  --
 30  -- To avoid spurious outliers, this code use the 90% furthest points.
 31  --
 32  -- To avoid long runtimes, uses a subset of the data to learn where
 33  -- to divide data (then all the data gets pushed down first halves).
 34  --
 35  -- To support explanation, optionally, at each level of recursion,
 36  -- this code reports what ranges can best distinguish sibling clusters
 37  -- C1,C2.  The  discretizer is inspired by the ChiMerge algorithm:
 38  -- numerics are divided into, say, 16 bins. Then, while we can find
 39  -- adjacent bins with the similar distributions in C1,C2, then
 40  -- (a) merge then (b) look for other merges.
 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_ &rArr; 'keep=N';
 62  -- and  _-booleanFlag_ &rArr; '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 end
 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)
107    u={}
108    for k,v in pairs(t) do k=tostring(k);if k:sub(1,1)~="_" then push(u,k) end end
109    return sort(u) end
110
111  -- print tables, recursively
112  function oo(t)  print(o(t)) end
113  function o(t)
114    if type(t)~="table" then return tostring(t) end
115    local key=function(k) return fmt(":%s %s",k,o(t[k])) end
116    local u = #t>0 and map(t,o) or map(slots(t),key)
117    return '{'..table.concat(u,"")..."}" end
118
119  -- strings to things
120  function csv(file,       x)
121    file = io.input(file)
122    return function()
123      x=io.read(); if x then return things(x) else io.close(file) end end end
124
125  function thing(x)
126    x = x:match"^%s*(.-)%s*$"
127    if x=="true" then return true elseif x=="false" then return false end
128    return tonumber(x) or x end
129
130  function things(x,sep,  t)
131    t={}
132    for y in x:gmatch(sep or"([^,]+)") do push(t,thing(y)) end
133    return t end
```

```lua
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  --   |  _ \ |  _ \  | |  |  ___| |  ___|
144  --   | |_) || |_) | | |  | |_    | |_
145  --   |  _ < |  __/  | |  |  _|   |  _|
146  --   |_| \_\|_|     |_|  |_|     |_|
147  function new(k,t) k.__index=k; k.__tostring=o; return setmetatable(t,k) end
148
149  -- COLS: turns list of column names into NUMs, SYMs, or SKIPs
150  function COLS.new(k,row,    i)
151    i= new(k,{all={},x={},y={},names=row})
152    for at,txt in ipairs(row) do push(i.all, i:col(at,txt)) end
153    return i end
154
155  function COLS.add(i,t)
156    for _,col in pairs(i.all) do col:add( t[col.at] ) end
157    return t end
158
159  function COLS.col(i,at,txt,      col)
160    if ignorep(txt) then return SKIP:new(at,txt) end
161    col = (nump(txt) and NUM or SYM):new(at,txt)
162    push(goalp(txt) and i.y or i.x, col)
163    if klassp(txt) then i.klass = col end
164    return col end
165
166  -- NUM: summarizes a stream of numbers
167  function NUM.new(k,n,s)
168    return new(k,{n=0,at=n or 0,txt=s or"",has=SOME:new(),ok=false,
169            w=lessp(s or "") and -1 or 1, lo=big, hi=-big}) end
170
171  function NUM.add(i,x)
172    if x ~= "?" then
173      i.n = i.n + 1
174      if i.has:add(x) then i.ok=false end
175      i.lo,i.hi = min(x,i.lo), max(x,i.hi); end end
176
177  function NUM.dist(i,x,y)
178    if      x=="?" and y=="?" then return 1
179    elseif x=="?" then y=i:norm(y); x=y<0.5 and 1 or 0
180    elseif y=="?" then x=i:norm(x); y=x<0.5 and 1 or 0
181    else    x,y = i:norm(x), i:norm(y) end
182    return math.abs(x-y) end
183
184  function NUM.distance2heaven(x, w)
185    return ((i.w>0 and 1 or 0) - i:norm(x))^2 end
186
187  function NUM.mid(i) return per(i:sorted(), .5) end
188
189  function NUM.norm(i,x)
190    return math.abs(i.hi-i.lo)<1E-9 and 0 or (x-i.lo)/(i.hi - i.lo) end
191
192  function NUM.sorted(i)
193    if i.ok==false then table.sort(i.has.all); i.ok=true end
194    return i.has.all end
195
196  -- ROWS: manages 'rows', summarized in 'cols' (columns).
197  function ROWS.new(k,inits,     i)
198    i = new(k,{rows={},cols=nil})
199    if type(inits)=="string" then for t in csv(inits) do i:add(t) end end
200    if type(inits)=="table"  then for t in inits       do i:add(t) end end
201    return i end
202
203  function ROWS.add(i,t)
204    if i.cols then push(i.rows,i.cols:add(t)) else i.cols=COLS:new(t) end end
205
206  function ROWS.clone(i,  j) j= ROWS:new(); j:add(i.cols.names);return j end
207
208  function ROWS.dist(i,row1,row2,    d,fun)
209    function fun(col) return col:dist(row1[col.at], row2[col.at])^the.p end
210    return (sum(i.cols.x, fun)/ #i.cols.x)^(1/the.p) end
211
212  function ROWS.furthest(i,row1,rows,       fun)
213    function fun(row2) return {i:dist(row1,row2), row2} end
214    return unpack(per(sort(map(rows,fun),firsts), the.furthest)) end
215
216  function ROWS.half(i, top)
217    local some, top,c,x,y,tmp,mid,lefts,rights,_
218    some= many(i.rows, the.keep)
219    top = top or i
220    _,x = top:furthest(any(some), some)
221    c,y = top:furthest(x,          some)
222    tmp = sort(map(i.rows,function(r) return top:fastmap(r,x,y,c) end),firsts)
223    mid = #i.rows//2
224    lefts, rights = i:clone(), i:clone()
225    for at,row in pairs(tmp) do (at <=mid and lefts or rights):add(row[2]) end
226    return lefts,rights,x,y,c, tmp[mid] end
227
228  function ROWS.mid(i,cols)
229    return map(cols or i.cols.all, function(col) return col:mid() end) end
230
231  function ROWS.fastmap(i, r,x,y,c,      a,b)
232    a,b = i:dist(r,x), i:dist(r,y); return {(a^2 + c^2 - b^2)/(2*c), r} end
233
234  -- SKIP: summarizes things we want to ignore (so does nothing)
235  function SKIP.new(k,n,s)  return new(k,{n=0,at=at or 0,txt=s or""}) end
236  function SKIP.add(i,x)    return x end
237  function SKIP.mid(i)      return "?" end
238
239  -- SOME: keeps a random sample on the arriving data
240  function SOME.new(k,keep) return new(k,{n=0,all={}, keep=keep or the.keep}) end
241  function SOME.add(i,x)
242    i.n = i.n+1
243    if      #i.all < i.keep then push(i.all,x)            ; return i.all
244    elseif r()      < i.keep/i.n then i.all[r(#i.all)]=x; return i.all end end
245
246  -- SYM: summarizes a stream of symbols
247  function SYM.new(k,n,s)
248    return new(k,{n=0,at=n or 0,txt=s or"",has={},most=0}) end
249
250  function SYM.add(i,x,inc)
251    if x ~= "?" then
252      inc = inc or 1
253      i.n = i.n + inc
254      i.has[x] = inc + (i.has[x] or 0)
255      if i.has[x] > i.most then i.most,i.mode=i.has[x],x end end end
256
257  function SYM.dist(i,x,y) return(x=="?" and y=="?" and 1) or(x==y and 0 or 1) end
258  function SYM.mid(i)      return i.mode end
259  function SYM.div(i,   p)
260    return sum(i.has,function(k) p=-i.has[k]/i.n;return -p*math.log(p,2) end) end
261
262  function SYM.merge(i,j,    k)
263    k = SYM:new(i.at,i.txt)
264    for x,n in pairs(i.has)  do k:add(x,n) end
265    for x,n in pairs(j.has)  do k:add(x,n) end
266    ei, ej, ejk= i:div(), j:div(), k:div()
267    if i.n==0 or j.n==0 or .99*ek <= (i.n*ei + j.n*ej)/k.n then
268      return k end end
```

page 2

```lua
--
-- CLUSTER
--
--
-- CLUSTER: recursively divides data by clustering towards two distant points
function CLUSTER.new(k,egs,top)
  local i,want,left,right
  i    = new(k, {here=egs})
  top  = top or egs
  want = (#top.rows)^the.want
  if #egs.rows >= 2*want then
    left, right, i.x, i.y, i.c, i.mid = egs:half(top)
    if #left.rows < #egs.rows then
      i.left = CLUSTER:new(left,  top)
      i.right= CLUSTER:new(right, top) end end
  return i end

function CLUSTER.show(i,pre,  here)
  pre = pre or ""
  here=""
  if not i.left and not i.right then here= o(i.here:mid(i.here.cols.y)) end
  print(fmt("%6s :%-30s %s",#i.here.rows, pre, here))
  for _,kid in pairs{i.left, i.right} do
    if kid then kid:show(pre .. "|.. ") end end end

--
-- EXPLAIN
--
-- SPAN: keeps a random sample on the arriving data
function SPAN.new(k, col, lo, hi, has)
  return new(k,{col=col,lo=lo,hi=hi or lo,has=has or SYM:new()}) end

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
function SPAN.merge(i,j)
  local has = i.has:merge(j.has)
  if now then return SPAN:new(i.col, i.lo, j.hi, has) end end

function SPAN.select(i,row,   x)
  x = row[i.col.at]
  return (x=="?") or (i.lo==i.hi and x==i.lo) or (i.lo <= x and x < i.hi) end

function SPAN.score(i) return {i.has.n/i.col.n,  i.has:div()} end


-- EXPLAIN:
function EXPLAIN.new(k,egs,top)
  local i,top,want,left,right,spans,best,yes,no
  i    = new(k,{here = egs})
  top  = top or egs
  want = (#top.rows)^the.want
  if #top.rows >= 2*want then
    left,right = egs:half(top)
    spans  = {}
    for n,col in pairs(i.cols.x) do
      for _,s in pairs(col:spans(j.cols.x[n])) do
        push(spans,{ys=s:score(),it=s}) end end
    best  = distance2heaven(spans,{"+","-"})[1]
    yes,no = egs:clone(), egs:clone()
    for _,row in pairs(egs.rows) do
      (best:selects(row) and yes or no):add(row)  end -- divide data in two
    if #yes.rows<#egs.rows then -- make kids if kid size different to parent size
      if #yes.rows>=want then i.yes=EXPLAIN:new(yes,top) end
      if #no.rows >=want then i.no =EXPLAIN:new(no, top)  end end end
  return i end

function EXPLAIN.show(i,pre)
  pre = pre or ""
  if not pre then
    tmp = i.here:mid(i.here.y)
  print(fmt("%6s :%-30s %s", #i.here.rows, pre, o(i.here:mid(i.here.cols.y))))
  for _,pair in pairs{{true,i.yes},{false,i.no}} do
    status,kid = unpack(pair)
    k:shpw(pre .. "|.. ") end end end

function SYM.spans(i, j)
  local xys,all,one,last,xys,x,c n = {},{}
  for x,n in pairs(i.has) do push(xys, {x,"this",n}) end
  for x,n in pairs(j.has) do push(xys, {x,"that",n}) end
  for _,tmp in ipairs(sort(xys,firsts)) do
    x,c,n = unpack(tmp)
    if x ~= last then
      last = x
      one  = push(all, Span(i,x,x)) end
    one:add(x,y,n) end
  return all end

function NUM.spans(i, j)
  local xys,all,lo,hi,gap,xys,one,x,c,n = {},{}
  lo,hi = min(i.lo, j.lo), max(i.hi,j.hi)
  gap   = (hi - lo) / (6/the.cohen)
  for x,n in pairs(i.has) do push(xys, {x,"this",1}) end
  for x,n in pairs(j.has) do push(xys, {x,"that",1}) end
  one = Span:new(i,lo,lo)
  all = {one}
  for _,tmp in ipairs(sort(xys,first)) do
    x,c,n = unpack(tmp)
    if one.hi - one.lo > gap then one = push(all, Span(i, one.hi, x)) end
    one:add(x,y) end
  all           = merge(all)
  all[1   ].lo = -big
  all[#all].hi =  big
  return all end

function merge(b4,        j,n,now,a,b,merged)
  j,n,now = 0,#b4,{}
  while j < #b4 do
    j    = j+1
    a, b = b4[j], b4[j+1]
    if b then
      merged = a:merge(b)
      if merged then a,j = merged, j+1 end end
    push(now,a)
    j = j+1 end
  return #now == #b4 and b4 or merge(now) end
```

```lua
--
-- DEMOS
--
--
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 algorithms. Mach Learn 6, 37â€"66 (1991).  https://doi.org/10.1007/BF00153759
--
```