```lua
#!/usr/bin/env lua
--      _
--   __| |_   _  ___
--  / _` | | | |/ _ \
-- | (_| | |_| | (_) |
--  \__,_|\__,_|\___/  .lua

local F=require"fun"
local the=F.options[[

lua duo.lua [OPTIONS]
(c)2022 Tim Menzies, MIT license

Data miners using/used by optimizers.
Understand N items after log(N) probes, or less.

OPTIONS:
  -ample   when enough is enough         = 512
  -bins    initial bins size             = 16
  -Debug   on error, dump stack and halt = false
  -enough  use (#t)^enough               = .5
  -far     how far to go                 = .9
  -file    read data from file           = ../etc/data/auto93.csv
  -help    show help                     = false
  -p       distance coefficient          = 2
  -rnd     default round                 = %5.2f
  -seed    random number seed            = 10019
  -task    start up actions              = donothing]]

local EGS, NUM, RANGE, SYM = {}, {}, {}, {}
local    any,  asserts,  brange,  firsts,  fmt,  many,  map,  mapp =
      F.any,F.asserts,F.brange,F.firsts,F.fmt,F.many,F.map,F.mapp
local    new,  o,  oo,  push,  rows,  seconds,  sort =
      F.new,F.o,F.oo,F.push,F.rows,F.seconds,F.sort
-- # RANGE
--
-- |**Does**  | Create  |: models a span from `lo` to `hi`            |
-- |----      |-------  |---------------------------------------------|
-- |          | Sort    | : know the mean and standard deviation      |
-- |          | 3       | : support inference; e.g. distance, likelihood|
-- |**Has**   | n       | : counter of things seen so far             |
-- |          | at      | : column index                              |
-- |          | name    | : column name                               |
-- |          | mu      | : mean seen so far                          |
-- |          | lo      | : smallest number seen so far               |
-- |          | hi      | : largest number  seen so far               |
-- |          | sd      | : standard deviation                        |
-- |          | some    | : stores a sample of the symbols            |
-- |          | w       | : (for minimize) and  1 (for maximize)      |
-- |          | _m2     | : incrementally 2nd moment (internal)       |
-- |**Uses**  |         | : [Some](some.html)                         |
--
-- ## Create
-- **RANGE:new(col:NUM|SYM, lo:num, hi:num, b:num, B:num, r:num, R:num):RANGE**

function RANGE.new(k,col,lo,hi,b,B,r,R)
  return new(k,{col=col,lo=lo,hi=hi or lo,b=b,B=B,r=r,R=R}) end

-- **i:RANGE:merge(j:RANGE):RANGE?**
-- Return a combined range (if it has better value) or return nil.
function RANGE.merge(i,j,k,    lo,hi,z,B,R)
  lo = math.min(i.lo, j.lo)
  hi = math.max(i.hi, j.hi)
  z=1E-31; B,R = i.B+z, i.R+z
  k = RANGE:new(i.col,lo,hi,i.b+j.b,i.B,i.r+j.r, j.R)
  if k.b/B < .01 or k.r/R < .01             then return k end
  if k:val() > i:val() and k:val() > j:val() then return k end end

function RANGE.__lt(i,j) return i:val() < j:val() end

function RANGE.show(i)
  if i.lo == i.hi        then return fmt("%s == %s", i.col.txt, i.lo) end
  if i.lo == -math.huge then return fmt("%s < %s",  i.col.txt, i.hi) end
  if i.hi == math.huge  then return fmt("%s >= %s", i.col.txt, i.lo) end
  return fmt("%s <= %s < %s", i.lo, i.col.txt, i.hi) end

function RANGE.val(i,    z,B,R)
  z=1E-31; B,R = i.B+z, i.R+z; return (i.b/B)^2/( i.b/B + i.r/R) end

function RANGE.selects(i,row,    x)
  x=row.has[col.at]; return x=="?" or i.lo<=x and x<i.hi end

-- Class methods
function RANGE.merged(b4)
  local j,tmp,now,after,maybe = 0, {}
  while j < #b4 do
    j = j + 1
    now, after = b4[j], b4[j+1]
    if after then
      maybe = now:merge(after)
      if maybe then now=maybe; j=j+1 end end
    push(tmp,now) end
  return #tmp==#b4 and b4 or RANGE.merged(tmp) end

function RANGE.uninformative(t)
  return #t == 1 and #t[1].lo == -math.huge and #t[1].hi == math.huge end

-- # NUM
function NUM.new(k,at,s)
  return new(k,{n=0, at=at,txt=s,w=s:find"-" and -1 or 1,_has={},
               ok=false, lo=math.huge, hi=-math.huge}) end

function NUM.add(i,x)
  if x ~= "?" then
    i.ok = false
    i.n = i.n + 1
    push(i._has, x)
    if x < i.lo then i.lo = x end
    if x > i.hi then i.hi = x end end
  return x end

function NUM.dist(i,a,b)
  if     a=="?" and  b=="?" then a,b=1,0
  elseif a=="?" then b    = i:norm(b); a=b>.5 and 0 or 1
  elseif b=="?" then a    = i:norm(a); b=a>.5 and 0 or 1
  else               a, b= i:norm(a), i:norm(b) end
  return math.abs(a-b) end

function NUM.has(i)
  if not i.ok then sort(i._has); i.ok=true end; return i._has end

function NUM.mid(i,    a) a=i:has(); return a[#a//2] end

function NUM.norm(i,x)
  return i.hi - i.lo<1E-9 and 0 or (x - i.lo)/(i.hi - i.lo) end

function NUM.with(i,lo,hi,    t,left,right)
  t=i:has()
  if hi<t[1] or lo>t[#t] then return 0 end
  left= lo < t[1] and 1 or F.bleft(t,lo)
  right= hi > t[#t] and #t or F.bright(t,hi)
  return  right - left end

-- compare to old above
function NUM.ranges(i,j)
  local out,lo,hi,gap = {}
  lo  = math.min(i.lo,j.lo)
  hi  = math.max(i.hi, j.hi)
  gap = (hi - lo) / the.bins
  for x = lo,hi,gap do
    push(out,--col,lo hi        b               B          r               R
      RANGE:new(i, x, x+gap,i:with(x,x+gap),#i:has(), j:with(x,x+gap),#j:has()))
  end
  out = RANGE.merged(out)
  out[1].lo = -math.huge
  out[#out].hi =  math.huge
  return out end

-- # SYM
function SYM.new(k,at,s) return new(k,{n=0,at=at,txt=s,_has={},mode=nil,most=0})
  end
function SYM.add(i,x)
  if x~="?" then
    i.n = i.n + 1
    i._has[x] = 1 + (i._has[x] or 0)
    if i._has[x] > i.most then i.most, i.mode = i._has[x],x end
  end
  return x end

function SYM.dist(i,a,b) return  a=="?" and b=="?" and 1 or a==b and 0 or 1 end
function SYM.has(i)        return i.has end
function SYM.mid(i)        return i.mode end
function SYM.ranges(i,j,     out)
  return mapp(i._has,
    function(x,n) return RANGE:new(i,x,x,n,i.n,(j._has[x] or 0),j.n) end) end

-- # EGS
function EGS.new(k,file,    i)
  i= new(k,{_rows={}, cols=nil, x={},  y={}})
  if file then for row in rows(file) do i:add(row) end end
  return i end

function EGS.add(i,t)
  local add,now,where = function(col) return col:add(t[col.at]) end
  if i.cols then
    push(i._rows, map(i.cols, add))
  else
    i.cols = {}
    for n,x in pairs(t) do
      now = push(i.cols, (x:find"^[A-Z]" and NUM or SYM):new(n,x))
      if not x:find":" then
        push((x:find"+" or x:find"-") and i.y or i.x, now) end end end end

function EGS.clone(i,inits,    j)
  j = EGS:new()
  j:add(map(i.cols, function(col) return col.txt end))
  for _,row in pairs(inits or {}) do j = j:add(row) end
  return j end

function EGS.mid(i,cols)
  return map(cols or i.y, function(col) return col:mid() end) end

function EGS.dist(i,r1,r2)
  local d,n,inc = 0, (#i.x)+1E-31
  for _,col in pairs(i.x) do
    inc = col:dist(r1[col.at], r2[col.at])
    d   = d + inc^the.p end
  return (d/n)^(1/the.p) end

function EGS.far(i,r1,rows,       act,tmp)
  act = function(r2) return {r2, i:dist(r1,r2)} end
  tmp = sort(map(rows,act), seconds)
  return table.unpack(tmp[#tmp*the.far//1] ) end

function EGS.half(i,rows)
  local some,left,right,c,cosine,lefts,rights
  rows   = rows or i._rows
  some   = #rows > the.ample and many(rows, the.ample) or rows
  left   = i:far(any(rows), some)
  right,c = i:far(left,     some)
  function cosine(r,     a,b)
    a, b = i:dist(r,left), i:dist(r,right); return {(a^2+c^2-b^2)/(2*c),r} end
  lefts,rights = i:clone(), i:clone()
  for n,pair in pairs(sort(map(rows,cosine), firsts)) do
    (n <= (#rows)/2 and lefts or rights):add( pair[2] ) end
  return lefts,rights,left,right,c end

local rnd,show
function EGS.cluster(i, top)
  local c,lefts0, rights0, lefts, rights, left, right=0
  top = top or i
  if #i._rows >=  2*(#top._rows)^the.enough then
    lefts0, rights0, left, right, c = top:half(i._rows)
    lefts  = lefts0:cluster( top)
    rights = rights0:cluster(top)
  end
  return {here=i, lefts=lefts, rights=rights, left=left, right=right, c=c} end

function rnd(x)
  return fmt((type(x)=="number" and x~=x//1 and the.rnd or"%s",x) end

function show(t,lvl)
  lvl = lvl or ""
  if t then
    if t.lefts
    then print(fmt("%s%s",lvl,#t.here._rows))
    else print(fmt("%s%s\t%s", lvl,#t.here._rows, o(t.here:mid()))) end
    show(t.lefts, lvl.."|..")
    show(t.rights,lvl.."|..") end end

-- # Tests
local no,go={},{}

function go.any(    t,x,n)
  t={}; for i=1,10 do t[1+#t] = i end
  n=0; for i=1,5000 do x=any(t); n= 1 <= x and x <=10 and n+1 or 0 end
  asserts(n==5000,"any")  end

function go.bleft(    t,x,a,b,bad)
  t,bad = {},0
  for j =1,30 do push(t,100*math.random()//1) end
  table.sort(t)
  for k,v in pairs(t) do print(k,v) end
  for j=1,5 do x=any(t); print(x, F.bleft(t,x)) end
  for j=1,5 do x=100*math.random()//1; print(x, F.bleft(t,x)) end
  x= 200; print(x, F.bleft(t,x))
  x= -1; print(x, F.bleft(t,x))
end

function go.bspan(    t,x,a,b,bad)
  t,bad = {},0
  for j =1,50 do push(t,10*math.random()//1) end
  table.sort(t);
  for k,v in pairs(t) do print(k,v) end
  print""
  for j =1,10 do
    x=any(t); a,b = F.bleft(t,x),F.bright(t,x); print("=",x,a,b) end
  print""
  for j =1,10 do
    x=math.random(100)/10 a,b = F.bleft(t,x),F.bright(t,x); print("ish",x,a,b)
  end
  end
```

```lua
269
270  function no.fail()        asserts(fail,"checking crashes"); print(no.thi.ng) end
271  function go.oo(   u)      asserts("{10 20 30}" == fmt("%s",o{10,20,30}),"table") end
272  function go.rows( t)
273    for row in rows(the.file) do t=row  end
274    asserts(type(t[1])=="number","is number")
275    asserts(t[1]==4, "is four")
276    asserts(#t==8,"is eight") end
277
278  function go.egs(    i,t)
279    i=EGS:new(the.file)
280    asserts(i.y[1].lo==1613,"lo")
281    t=i.y[1]:has(); asserts(1613==t[1],"lo2") asserts(5140== t[#t],"hi");
282    asserts(i.y[1].ok,"ok") end
283
284  function go.dist(  i, t,a,b,d)
285    i=EGS:new(the.file)
286    t= i._rows
287    for j=1,100 do
288      a,b= any(t), any(t)
289      d= i:dist(a,b)
290      assert(0<= d and d <= 1) end end
291
292  function go.half(  a,b,col2,tmp)
293    local top =EGS:new(the.file)
294    local lefts,rights,left,right,c=top:half()
295    asserts(top:dist(left,right) > .75)
296    for n,col1 in pairs(lefts.x) do
297      col2 = rights.x[n]
298      print""
299      tmp= col1:ranges(col2)
300      if #tmp>1 then
301        for n,r in pairs(tmp) do print(0,col1.txt, n,rnd(r.lo),rnd(r.hi),r.b,r.r,
rnd(r:val())) end end
302      print""
303      tmp= col2:ranges(col1)
304      if #tmp>1 then
305        for n,r in pairs(tmp) do print(1,col1.txt, n,rnd(r.lo),rnd(r.hi),r.b,r.r,
rnd(r:val())) end end
306    end
307  end
308
309  function go.cluster(  a)
310    a=EGS:new(the.file):cluster()
311    asserts(49==#a.lefts.lefts.lefts.here._rows)
312    end
313
314  if arg[0] == "duo.lua" then the(go) end -- if called as main function
315  return {the=the, EGS=EGS, NUM=NUM, RANGE=RANGE, SYM=SYM}
```