```lua
#!/usr/bin/env lua
local b4={}; for k,v in pairs(_ENV) do b4[k]=v end;

--
--       _ _      _   _ _
--      | | |    | |     |  
--      | | |    | |   _  |  
--      | | |    | |  | | |  
--                         
--    a little like
--    LUA learning
--    library
--
--

local options={

what  = "Small sample multi-objective optimizer.",
usage= "(c) 2021 Tim Menzies <timm@ieee.org> unlicense.org",
about= [[
Sort N examples on multi-goals using a handful of 'hints'; i.e.

- Evaluate and rank, a few examples (on their y-values);
- Sort other examples by x-distance to the ranked ones;
- Recurse on the better half (so we sample more and more
  from the better half, then quarter, then eighth...).

A regression tree learner then explores the examples (sorted
left to right, worst to best).  By finding branches that
reduce the variance of the index of those examples, this
tree reports what attribute ranges select for the better (or
worse) examples.  ]],

how=  {{"file",    "-f",  "../../data/auto93.csv",  "read data from file"},
       {"cull",    "-c",   .5    ,"cuts per generation"},
       {"help",    "-h",   false ,"show help"                       },
       {"hints",   "-H",   4     ,"hints per generation"          },
       {"p",       "-p",   2     ,"distance calc exponent"       },
       {"small",   "-s",   .5    ,"div list t into t^small"   },
       {"seed",    "-S",   10019 ,"random number seed"          },
       {"train",   "-t",   .5    ,"size of training set"         },
       {"trivial", "-T",   .35   ,"small delta=trivial*sd"      },
       {"todo",    "-T",   "all" ,"run unit test, or 'all'"    },
       {"wild",    "-W",   false ,"run tests, no protection"  }}}

local the={} -- a flat list of key=value options; e.g. {seed=10019,p=2,...}
for _,t in pairs(options.how) do -- update defaults from command line
  the[t[1]] = t[3]
  for n,word in ipairs(arg) do if word==t[2] then
    the[t[1]] = t[3] and (tonumber(arg[n+1]) or arg[n+1]) or true end end end

if the.help then --  print help text
  print(string.format("\n%s [OPTIONS]\n%s\n%s\n\nOPTIONS:\n",
                       arg[0], options.usage, options.what))
  for _,t in pairs(options.how) do
    print(string.format("%4s %-9s%s\t%s %s",
            t[2], t[3] and t[1] or"", t[4], t[3] and"=" or"", t[3] or "")) end
  print("\n"..options.about)
  os.exit() end

--[[
Spans
 Little languages:
    - options
    - data language

Lesson plan
-- w1: ssytems: github. github workplaces. unit tests. doco tools.
-- w2: num,sym
-- W3: sample
-- w4: eval, knn, unfarinessness
-- W5:
```

```lua
--
--      |‾‾|  |   _|  (_       |_|  ‾|‾ |  |  _|
--

-- Random stuff
local Seed,rand,randi
Seed = the.seed or 10019
-- random integers
function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
-- random floats
function rand(lo,hi,     mult,mod)
  lo, hi = lo or 0, hi or 1
  Seed = (16807 * Seed) % 2147483647
  return lo + (hi-lo) * Seed / 2147483647 end

--------------------------------------------------------------------------
-- ## Table Stuff
local cat,map,lap,top,keys,last,copy,pop,push
local sort,firsts,first,second,shuffle,bchop
-- Table to string.
cat     = table.concat
-- Return a sorted table.
sort    = function(t,f) table.sort(t,f); return t end
-- Return first,second, last  item.
first   = function(t) return t[1] end
second  = function(t) return t[2] end
last    = function(t) return t[#t] end
-- Function for sorting pairs of items.
firsts  = function(a,b) return first(a) < first(b) end
-- Add to end, pull from end.
pop     = table.remove
push    = function(t,x) table.insert(t,x); return x end

-- Random order of items in a list (sort in place).
function shuffle(t,   j)
  for i=#t,2,-1 do j=randi(1,i); t[i],t[j]=t[j],t[i] end; return t end

-- Collect values, passed through 'f'.
function lap(t,f)  return map(t,f,1) end
-- Collect key,values, passed through 'f'.
-- If 'f' returns two values, store as key,value.
-- If 'f' returns one values, store at index value.
-- If 'f' return nil then add nothing (so 'map' is also 'select').
function map(t,f,one,     u)
  u={}; for x,y in pairs(t) do
    if one then x,y=f(y) else x,y=f(x,y) end
    if x ~= nil then
      if y then u[x]=y else u[1+#u]=x end end end
  return u end

-- Shallow copy
function copy(t,   u) u={}; for k,v in pairs(t) do u[k]=v end; return u end

function top(t,n,      u)
  u={};for k,v in pairs(t) do if k>n then break end; push(u,v) end; return u;end

--- Return a table's keys (sorted).
function keys(t,u)
  u={}
  for k,_ in pairs(t) do if tostring(k):sub(1,1)~="_" then push(u,k) end end
  return sort(u) end

-- Binary chop (assumes sorted lists)
function bchop(t,val,lt,lo,hi,     mid)
  lt = lt or function(x,y) return x < y end
  lo,hi = lo or 1, hi or #t
  while lo <= hi do
    mid =(lo+hi) // 2
    if lt(t[mid],val) then lo=mid+1 else hi= mid-1 end end
  return math.min(lo,#t)   end

--------------------------------------------------------------------------
-- ## Maths Stuff
local abs,sum,rnd,rnds
abs = math.abs
-- Round 'x' to 'd' decimal places.
function rnd(x,d,  n) n=10^(d or 0); return math.floor(x*n+0.5) / n end
-- Round list of items to  'd' decimal places.
function rnds(t,d) return lap(t, function(x) return rnd(x,d or 2) end) end

-- Sum items, filtered through 'f'.
function sum(t,f)
  f= f or function(x) return x end
  out=0; for _,x in pairs(f) do out = out + f(x) end; return out end

--------------------------------------------------------------------------
-- ## Printing Stuff
local out,shout,red,green,yellow,blue,color,fmt
fmt = string.format
-- Print as red, green, yellow, blue.
function color(s,n) return fmt("\27[1m\27[%sm%s\27[0m",n,s) end
function red(s)     return color(s,31) end
function green(s)   return color(s,32) end
function yellow(s)  return color(s,34) end
function blue(s)    return color(s,36) end

-- Printed string from a nested structure.
shout = function(x) print(out(x)) end
-- Generate string from a nested structures
-- (and don't print any contents more than once).
function out(t,seen,     u,key,value,public)
  function key(k)   return fmt(":%s %s",blue(k),out(t[k],seen)) end
  function value(v) return out(v,seen) end
  if type(t) == "function" then return "(...)" end
  if type(t) ~= "table"    then return tostring(t) end
  seen = seen or {}
  if seen[t] then return "..." else seen[t] = t end
  u = #t>0 and lap(t, value) or lap(keys(t), key)
  return red((t._is or"")..."{")..cat(u,"")..red("}") end

--------------------------------------------------------------------------
-- ## File i/o Stuff
-- Return one table per line, split on commans.
local csv
function csv(file,   line)
  file = io.input(file)
  line = io.read()
  return function(   t,tmp)
    if line then
      t={}
      for cell in line:gsub("[\t\r ]*","");gsub("#.*","");gmatch("([^,]+)") do
        push(t, tonumber(cell) or cell) end
      line = io.read()
      if #t>0 then return t end
    else io.close(file) end end end

--------------------------------------------------------------------------
-- ## OO Stuff
local has,obj
-- Create an instance
function has(mt,x) return setmetatable(x,mt) end

-- Create a clss
function obj(s, o,new)
  o = {_is=s, __tostring=out}
  o.__index = o
  return setmetatable(o,{__call = function(_,...) return o.new(...) end}) end
```

```lua
216  --
217  --        |¯|  |_|  |¯|        _\  \/  |¯|¯|
218  --                          /
219
220  -- ## Stuff for tracking `Num`bers.
221  -- `Num`s track a list of number, and can report  it sorted.
222  local Num=obj"Num"
223  function Num.new(inits,     self)
224    self= has(Num,{has={}, n=0, lo=1E32, hi =1E-32, ready=true})
225    for _,one in pairs(inits or {}) do self:add(one) end
226    return self end
227
228  function Num:add(x)
229    if    x>self.hi then self.hi = x
230    elseif x<self.lo then self.lo = x end
231    push(self.has,x); self.n=self.n+1; self.ready=false end
232
233  -- Ensure that the returned list of numbers is sorted.
234  function Num:all(x)
235    if not self.ready then table.sort(self.has) end
236    self.ready = true
237    return self.has end
238
239  function Num:dist(a,b)
240    if    a=="?" then b=self:norm(b); a = b>.5 and 0 or 1
241    elseif b=="?" then a=self:norm(a); b = a>.5 and 0 or 1
242    else   a,b = self:norm(a), self:norm(b) end
243    return abs(a-b) end
244
245  -- Combine two `num`s.
246  function Num:merge(other,    new)
247    new = Num.new(self.has)
248    for _,x in pairs(other.has) do new:add(x) end
249    return new end
250
251  -- Return a merged item if that combination
252  -- is simpler than its parts.
253  function Num:mergeable(other,   new,b4)
254    new = self:merge(other)
255    b4  = (self.n*self:sd() + other.n*other:sd()) / new.n
256    if b4 >= new:sd() then return new end end
257
258  -- The `mid` is the 50th percentile.
259  function Num:mid() return self:per(.5) end
260
261  -- Return `x` normalized 0..1, lo..hi.
262  function Num:norm(x,     lo,hi)
263    if x=="?" then return x end
264    lo,hi = self.lo, self.hi
265    return abs(hi - lo) < 1E-32 and 0 or (x - lo)/(hi - lo) end
266
267  -- Return the `p`-th percentile number.
268  function Num:per(p,     t)
269    t = self:all()
270    p = p*#t//1
271    return #t<2 and t[1] or t[p < 1 and 1 or p>#t and #t or p] end
272
273  -- The 10th to 90th percentile is 2.56 times the standard deviation.
274  function Num:sd() return (self:per(.9) - self:per(.1))/ 2.56 end
275
276  -- Create one span holding  row indexes associated with each number
277  local div -- defined below
278  function Num:spans(col,egs)
279    local xys,xs = {},  Num()
280    for pos,eg in pairs(egs) do
281      x = eg[col]
282      if x ~= "?" then
283        xs:add(x)
284        push(xys, {x=x,y=pos}) end end
285    return div(xys,                   -- split xys into spans...
286             #xs^the.small,           -- ..where spans are of size sqrt(#xs)..
287             xs:sd()*the.trivial) end -- ..and spans have (last-first)>trivial
288  -----------------------------------------------------------------------------
289  -- ## Stuff for tracking `Sym`bol Counts.
290  -- `Sym`s track symbol counts and the `mode` (most frequent symbol).
291  local Sym=obj"Sym"
292  function Sym.new(inits,     self)
293    self= has(Sym,{has={}, n=0, mode=nil, most=0})
294    for _,one in pairs(inits or {}) do self:add(one) end
295    return self end
296
297  function Sym:add(x)
298    self.n = self.n + 1
299    self.has[x] = 1 + (self.has[x] or 0)
300    if self.has[x] > self.most then self.most, self.mode = self.has[x], x end end
301
302  function Sym:dist(a,b) return a==b and 0 or 1 end
303  function Sym:mid() return self.mode end
304
305  -- Create one span holding  row indexes associated with each symbol
306  function Sym:spans(col,egs,...)
307    local xys,x = {}
308    for pos,eg in pairs(egs) do
309      x = eg[col]
310      if x ~= "?" then
311        xys[x] = xys[x] or {}
312        push(xys[x], pos)  end end
313    return map(xys, function(x,t) return {lo=x, hi=x, has=Num(t)} end) end
314
315
316
```

```lua
317  --
318  --       _\  (_|  |¯|¯|  |_)    |   (/_
319  --                             |
320
321  -- Samples store examples. Samples know about
322  -- (a) lo,hi ranges on the numerics
323  -- and (b) what  are independent `x` or dependent `y` columns.
324  local Sample = obj"Sample"
325  function Sample.new(     src,self)
326    self = has(Sample,{names=nil, all={}, ys={}, xs={}, egs={}})
327    if src then
328      if type(src)=="string" then for x   in csv(src) do self:add(x)   end end
329      if type(src)=="table" then for _,x in pairs(src) do self:add(x) end end end
330    return self end
331
332  function Sample:add(eg,     name,datum)
333    function name(col,new,    weight, where, what)
334      if new:find":" then return end
335      weight= new:find"-" and -1 or 1
336      what  = {col=col, w=weight, txt=new,
337               seen=(new:match("^[A-Z]",x) and Num() or Sym())}
338      where = (new:find("+") or new:find("-")) and self.ys or self.xs
339      push(self.all,  what)
340      push(where,     what) end
341    function datum(one,new)
342      if new ~= "?" then one.seen:add(new) end
343    end ---------------
344    if   not self.names
345    then self.names = eg
346         map(eg, function(col,x) name(col,x)  end)
347    else push(self.egs, eg)
348         map(self.all, function(_,col) datum(col,eg[col.col]) end) end
349    return self end
350
351  function Sample:better(eg1,eg2,     e,n,a,b,s1,s2)
352    n,s1,s2,e = #self.ys, 0, 0, 2.71828
353    for _,num in pairs(self.ys) do
354      a = num.seen:norm(eg1[num.col])
355      b = num.seen:norm(eg2[num.col])
356      s1 = s1 - e^(num.w * (a-b)/n)
357      s2 = s2 - e^(num.w * (b-a)/n) end
358    return s1/n < s2/n end
359
360  function Sample:betters(egs)
361    return sort(egs or self.egs,function(a,b) return self:better(a,b) end) end
362
363  function Sample:clone(      inits,out)
364    out = Sample.new():add(self.names)
365    for _,eg in pairs(inits or {}) do out:add(eg) end
366    return out end
367
368  function Sample:dist(eg1,eg2,     a,b,d,n,inc)
369    d,n = 0,0
370    for _,x in pairs(self.xs) do
371      a,b = eg1[x.col], eg2[x.col]
372      inc = a=="?" and b=="?" and 1 or x.seen:dist(a,b)
373      d   = d + inc^the.p
374      n   = n + 1 end
375    return (d/n)^(1/the.p) end
376
377  -- Report mid of the columns
378  function Sample:mid(cols)
379    return lap(cols or self.ys,function(col) return col.seen:mid() end) end
380
381  local div -- defined below
382  function Sample:tree(min,      node,min,sub,splitter, splitter1)
383    function splitter1(_,col,     out,xpect)
384      out   = col:spans(col,sample.eg, div)
385      xpect = sum(out, function(x) return x.has.n*x:sd() end)/#sample.egs
386      out   = map(out, function(_,x) x.has=x.has:all(); x.col= col end)
387      return {xpect,out} end
388    function splitter()
389      return first(sort(lap(sample.xs, splitter1), firsts))[2]
390    end ------------------------
391    node = {node=self, kids={}}
392    min  = min  or (#self.egs)^the.small
393    if #self.egs >= 2*min then
394      for _,span in pairs(splitter()) do
395        sub = self:clone()
396        for _,at in pairs(span.has) do sub:add(self.egs[at]) end
397        push(node.kids, span)
398        span.has = sub:tree(min) end end
399    return node end
400
401  -- at node
402  function Sample:where(tree,eg,    max,x,default)
403    if #kid.has==0 then return tree end
404    max = 0
405    for _,kid in pairs(tree.node) do
406      if #kid.has > max then default,max = kid,#kid.has end
407      x = eg[kid.col]
408      if x ~= "?" then
409        if x <= kid.hi and x >= kid.lo then
410          return self:where(kid.has.eg) end end end
411    return self:where(default, eg) end
412
413  -----------------------------------------------------------------------------
414  -- discretization tricks
415  -- Input a list of {{x,y}..} values. Return spans that divide the `x` values
416  -- to minimize variance on the `y` values.
417  function div(xys, tiny, dull,          now,out,x,y)
418    function merge(b4) -- merge adjacent spans if whole is simpler than the parts
419      local j, tmp = 0, {}
420      while j < #b4 do
421        j = j + 1
422        local now, after, simpler = b4[j], b4[j+1]
423        if after then
424          simpler = now.has:mergeable(after.has)
425          if simpler then
426            now = {lo=now.lo, hi= after.hi, has=simpler}
427            j = j + 1 end end
428        push(tmp,now)  end
429      return #tmp==#b4 and b4 or merge(tmp) -- recurse until nothing merged
430    end -------------------
431    local spans,span,out,x,y
432    xys   = sort(xys, function(a,b) return a.x < b.x end)
433    span  = {lo=xys[1].x, hi=xys[1].x, has=Num()}
434    spans = {span}
435    for j,xy in pairs(xys) do
436      x, y = xy.x, xy.y
437      if   j<#xys - tiny   and -- if enough items remaining after split
438           x~=xys[j+1].x   and -- if the next item is different (so we split here)
439           span.has.n>tiny and -- if span has enough items
440           span.hi - span.lo>dull -- if span is not trivially small
441      then now = push(spans, {lo=x, hi=x, has=Num()})  -- then new span
442      end
443      span.hi = x
444      span.has:add(y) end
445    return merge(spans) end
446
```

```
447  --     ┌─┐ ┬ ┌─┐ ┬─ ┬ ┌─┐ ┌─┐
448  --     │├┤ │ │ │ ┌┴─ │ │ │ │
449  --                              ─┘

450
451  -- Sorting on a few y values
452  local hints={}
453  function hints.default(eg) return eg end
454
455  function hints.sort(sample,scorefun,    test,train,egs,scored,small)
456    sample = Sample.new(the.file)
457    train,test = {}, {}
458    for i,eg in pairs(shuffle(sample.egs)) do
459      push(i<= the.train*#sample.egs and train or test, eg) end
460    egs = copy(train)
461    small = (#egs)^the.small
462    local i=0
463    scored = {}
464    while #egs >= small do
465      local tmp ={}
466      i = i + 1
467      io.stderr:write(fmt("%s",string.char(96+i)))
468      for j=1,the.hints do
469        egs[j] = (scorefun or hints.default)(egs[j])
470        push(tmp, push(scored, egs[j]))
471      end
472      egs = hints.ranked(scored,egs,sample)
473      for i=1,the.cull*#egs//1 do pop(egs) end
474    end
475    io.stderr:write("\n")
476    train=hints.ranked(scored, train, sample)
477    return #scored, sample:clone(train), sample:clone(test) end
478
479  function hints.ranked(scored,egs,sample,worker,   some)
480    function worker(eg) return {hints.rankOfClosest(scored,eg,sample),eg} end
481    scored = sample:betters(scored)
482    return  lap(sort(lap(egs, worker),firsts),second) end
483
484  function hints.rankOfClosest(scored,eg1,sample,        worker,closest)
485    function worker(rank,eg2) return {sample:dist(eg1,eg2),rank} end
486    closest = first(sort(map(scored, worker),firsts))
487    return  closest[2] end --+ closest[1]/10^8 end
488
```

```
489  --
490  -- ┌─┐ (7_ ┌┌┐ (_) _\
491
492  local eg={}
493  function eg.shuffle(   t)
494    t={}
495    for i=1,100 do push(t,i) end
496    assert(#t == #shuffle(t) and t[1] ~= shuffle(t)[1]) end
497
498  function eg.lap()
499    assert(3==lap({1,2},function(x) return x+1 end)[2]) end
500
501  function eg.map()
502    assert(3==map({1,2},function(_,x) return x+1 end)[2]) end
503
504  function eg.tables()
505    assert(20==sort(shuffle({{10,20},{30,40},{40,50}}),firsts)[1][2]) end
506
507  function eg.csv(   n,z)
508    n=0
509    for eg in csv(the.file) do n=n+1; z=eg end
510    assert(n==399 and z[#z]==50) end
511
512  function eg.rnds(    t)
513    assert(10.2 == first(rnds({10.22,81.22,22.33},1))) end
514
515  function eg.sym(    s)
516    s=Sym{"a","a","a","a","b","b","c"}
517    assert("a"==s.mode) end
518
519  function eg.num1(    n)
520    n=Num{10,20,30,40,50,10,20,30,40,50,10,20,30,40,50}
521    assert(.375 == n:norm(25))
522    assert(15.625 == n:sd()) end
523
524  function eg.num2(    n1,n2,n3,n4)
525    n1=Num{10,20,30,40,50,10,20,30,40,50,10,20,30,40,50}
526    n2=Num{10,20,30,40,50,10,20,30,40,50,10,20,30,40,50}
527    assert(n1:mergeable(n2)~=nil)
528    n3=Num{10,20,30,40,50,10,20,30,40,50,10,20,30,40,50}
529    n4=Num{100,200,300,400,500,100,200,300,400,500,100,200,300,400,500}
530    assert(n3:mergeable(n4)==nil) end
531
532  function eg.sample(    s,tmp,d1,d2,n)
533    s=Sample(the.file)
534    assert(2110 == last(s.egs)[s.all[3].col])
535    local sort1= s:betters(s.egs)
536    local lo, hi = s:clone(), s:clone()
537    for i=1,20              do lo:add(sort1[i]) end
538    for i=#sort1,#sort1-30,-1 do hi:add(sort1[i]) end
539    shout(s:mid())
540    shout(lo:mid())
541    shout(hi:mid())
542    for m,eg in pairs(sort1) do
543      n = bchop(sort1, eg,function(a,b) return s:better(a,b) end)
544      assert(m-n <=2) end end
545
546  function eg.dists(    s,tmp,d1,d2,n)
547    s=Sample(the.file)
548    tmp = sort(lap(shuffle(s.egs),
549                   function(eg2) return {s:dist(eg2,s.egs[1]), eg2} end),
550                firsts)
551    d1=s:dist(tmp[1][2], tmp[10][2])
552    d2=s:dist(tmp[1][2], tmp[#tmp][2])
553    assert(d1*10<d2) end
554
555  function eg.binsym(   s)
556    s=Sample(the.file)
557    print(s.all[6].seen._is=="Sym")
558    end
559
560  function eg.hints(    s,_,__,evals,sort1,train,test,n)
561    s=Sample(the.file)
562    --for _,eg in pairs(sort1) do lap(s.ys, function(col) return eg[col.col] end ) end
563    -- assert(s.ys[4].lo==1613)
564    evals, train,test = hints.sort(s)
565    test.egs = test:betters()
566    for m,eg in pairs(test.egs) do
567      n = bchop(train.egs, eg,function(a,b) return s:better(a,b) end)
568      print(n) end end
569
570  -------------------------------------------------------------------------------
571  -- startup
572  local fails, defaults = 0, copy(the)
573  local function example(k,      f,ok,msg)
574    f= eg[k]; assert(f,"unknown action "..k)
575    the=copy(defaults)
576    Seed=the.seed
577    if the.wild then return f() end
578    ok,msg = pcall(f)
579    if ok then print(green("PASS"),k)
580    else      print(red("FAIL"),  k,msg); fail=fail+1 end end
581
582  -- run one or more examples
583  if the.todo=="all" then lap(keys(eg),example) else example(the.todo) end
584  -- print any rogue global variables
585  for k,v in pairs(_ENV) do if not b4[k] then print("?rogue: ",k,type(v)) end end
586  -- exit, return  our test failure count.
587  os.exit(fail)
588
```

```lua
--[[

    ‾|‾ ‾    ‾| ‾
     |‾(_)   (_| (_)



--  seems to be  a revers that i  need to do .... but dont
-- check if shuffle is working

teaching:
- sample is v.useful
--]]
```