```lua
1   #!/usr/bin/env lua
2   --
3   --
4   --       a little lite
5   --
6   --
7   --
8   --       lua learning
9   --
10  --
11  --
12  --
13  --       library
14  --
15  --
16  --
17  --
18  local it=require"options"{
19  what = "Small sample multi-objective optimizer.",
20  who  = "(c) 2021 Tim Menzies <timm@ieee.org> unlicense.org",
21  why  = [[
22  Sort N examples on multi-goals using a handful of 'hints'; i.e.
23
24  - Evaluate and rank, a few examples (on their y-values);
25  - Sort other examples by x-distance to the ranked ones;
26  - Recurse on the better half (so we sample more and more
27    from the better half, then quarter, then eighth...).
28
29  A regression tree learner then explores the examples (sorted
30  left to right, worst to best).  By finding branches that
31  reduce the variance of the index of those examples, this
32  tree reports what attribute ranges select for the better (or
33  worse) examples.  ]],
34
35  how={{"FILE",        "-f",   "../../data/auto93.csv",  "read data from file"},
36        {"CULL",        "-c",    .5    , "cuts per generation"      },
37        {"HELP",        "-h",  false , "show help"              },
38        {"HINTS",       "-H",   4    ,"hints per generation"      },
39        {"P",           "-p",   2    ,"distance calc exponent"   },
40        {"TINY",        "-s",   .5   ,"div list t into t^small"    },
41        {"SEED",        "-S",  10019  ,"random number seed"       },
42        {"TRAIN",       "-t",   .5   , "size of training set"      },
43        {"TODO",        "-T",  "all" ,"run unit test, or 'all'"  },
44        {"TRIVIAL",     "-v",  .35   ,"small delta=trivial*sd"     },
45        {"WILD",        "-W",  false  ,"run tests, no protection"  }}}
46
47  local _=require"lib"
48  local abs,bchop,cat,copy        = _.abs,    _.bchop, _.cat,    _.copy
49  local csv,first,firsts,fmt,has  = _.csv,    _.first, _.firsts, _.fmt,   _.has
50  local keys,last,lap,map,obj     = _.keys,   _.last,  _.lap,    _.map,   _.obj
51  local out,pop,push,rand,shout   = _.out,    _.pop,   _.push,   _.rand,  _.shout
52  local rnd,rnds,rogues,second    = _.rnd,    _.rnds,  _.rogues, _.second
53  local shuffle,sort,sum,top      = _.shuffle, _.sort,  _.sum,    _.top
54
55  --[[
56  Spans
57   Little languages:
58      - options
59      - data language
60
61  Lesson plan
62  -- w1: ssytems: github. github workplaces. unit tests. doco tools.
63  -- w2: num,sym
64  -- W3: sample
65  -- w4: eval, knn, unfarinessness
66  -- W5:
67  --]]
68
```

```lua
69   -- NUM ----------------------------------------------------------------
70   --
71   -- ## Stuff for tracking `Num`bers.
72   -- `Num`s track a list of number, and can report  it sorted.
73   local Num=obj"Num"
74   function Num.new(inits,at, txt,      self)
75     self= has(Num,{at=at or 0, txt=txt or"", w=(txt or""):find"-" and -1 or 1,
76               has={}, n=0, lo=1E32, hi =1E-32, ready=true})
77     for _,one in pairs(inits or {}) do self:add(one) end
78     return self end
79
80   function Num:add(x)
81     if    x>self.hi then self.hi = x
82     elseif x<self.lo then self.lo = x end
83     push(self.has,x); self.n=self.n+1; self.ready=false end
84
85   -- Ensure that the returned list of numbers is sorted.
86   function Num:all(x)
87     if not self.ready then table.sort(self.has) end
88     self.ready = true
89     return self.has end
90
91   function Num:dist(a,b)
92     if    a=="?" then b=self:norm(b); a = b>.5 and 0 or 1
93     elseif b=="?" then a=self:norm(a); b = a>.5 and 0 or 1
94     else  a,b = self:norm(a), self:norm(b) end
95     return abs(a-b) end
96
97   -- Combine two `num`s.
98   function Num:merge(other,     new)
99     new = Num()
100    new.at, new.txt = self.at, self.txt
101    for _,x in pairs(self.has)  do new:add(x) end
102    for _,x in pairs(other.has) do new:add(x) end
103    return new end
104
105  -- The 'mid' is the 50th percentile.
106  function Num:mid() return self:per(.5) end
107
108  -- Return 'x' normalized 0..1, lo..hi.
109  function Num:norm(x,      lo,hi)
110    if x=="?" then return x end
111    lo,hi = self.lo, self.hi
112    return abs(hi - lo) < 1E-32 and 0 or (x - lo)/(hi - lo) end
113
114  -- Return the 'p'-th percentile number.
115  function Num:per(p,     t)
116    t = self:all()
117    p = p*#t//1
118    return #t<2 and t[1] or t[p < 1 and 1 or p>#t and #t or p] end
119
120  -- The 10th to 90th percentile range is 2.56 times the standard deviation.
121  function Num:sd() return math.abs(self:per(.9) - self:per(.1))/ 2.56 end
122  function Num:spread() return self:sd() end
123
124  -- Create one span (each has the  row indexes of the rows)
125  -- where each span has at least `tiny` items and  span is more than
126  -- `tirvial`ly small.
127  local div -- defined below
128  function Num:spans(sample,tiny,trivial)
129    local xys = {}
130    for _,eg in pairs(sample.egs) do
131      local x = eg[self.at]
132      if x ~= "?" then push(xys, {col=col, x=x, y=eg[sample.klass.at]}) end end
133    return div(xys, tiny, trivial,  self, getmetatable(sample.klass)) end
134
135  -- SYM ----------------------------------------------------------------
136  --
137  -- Stuff for tracking `Sym`bol Counts.
138  -- `Sym`s track symbol counts and the `mode` (most frequent symbol).
139  local Sym=obj"Sym"
140  function Sym.new(inits,at,txt,       self)
141    self= has(Sym,{at=at or 0, txt=txt or "", has={}, n=0, mode=nil, most=0})
142    for _,one in pairs(inits or {}) do self:add(one) end
143    return self end
144
145  function Sym:add(x,n)
146    n = n or 1
147    self.n = self.n + n
148    self.has[x] = n + (self.has[x] or 0)
149    if self.has[x] > self.most then self.most, self.mode = self.has[x], x end end
150
151  function Sym:dist(a,b) return a==b and 0 or 1 end
152
153  function Sym:merge(other)
154    new=Sym()
155    new.at, new.txt = self.at, self.txt
156    for k,n in pairs(self.has)  do new:add(k,n) end
157    for k,n in pairs(other.has) do new:add(k,n) end
158    return new end
159
160  function Sym:mid() return self.mode end
161
162  -- Create one span holding  row indexes associated with each symbol
163  function Sym:spans(sample,...)
164    local xys,yklass = {}, getmetatable(sample.klass)
165    for pos,eg in pairs(sample.egs) do
166      local x = eg[self.at]
167      if x ~= "?" then
168        xys[x] = xys[x] or yklass()
169        xys[x]:add( eg[sample.klass.at] )   end end
170    return map(xys, function(x,ys) return {col=self, lo=x, hi=x, has=ys} end) end
171
172  function Sym:spread()
173    return sum(self.has,
174               function(n1)  return  -n1/self.n * math.log(n1/self.n,2) end) end
175
176  -- SKIP ----------------------------------------------------------------
177  --
178  -- ## Stuff for skipping all things sent to a column
179  local Skip=obj"Skip"
180  function Skip.new(_,at,txt) return has(Skip,{at=at or 0, txt=txt or"", n=0}) end
181  function Skip:add(x) self.n = self.n + 1; return  x end
182  function Skip:mid() return "?" end
183
```

```lua
184  -- SAMPLE -------------------------------------------------------------
185  --
186  -- Samples store examples. Samples know about
187  -- (a) lo,hi ranges on the numerics
188  -- and (b) what  are independent 'x' or dependent 'y' columns.
189  local Sample = obj"Sample"
190  function Sample.new(    src,self)
191    self = has(Sample,{names=nil, klass=nil, all={}, ys={}, xs={}, egs={}})
192    if src then
193      if type(src)=="string" then for x  in csv(src) do self:add(x)  end end
194      if type(src)=="table" then for _,x in pairs(src) do self:add(x) end end end
195    return self end
196
197  function Sample:add(eg,      ako,what,xy)
198    if not self.names
199    then -- create the column headers
200      self.names = eg
201      for at,x in pairs(eg) do
202        ako  = (x:find":"        and Skip) or
203               (x:match"^[A-Z]" and Num ) or
204               Sym
205        what = push(self.all, ako({}, at, x))
206        if not x:find":" then
207          if x:find"!" then self.klass = what end
208          xy = (x:find"+") or x:find"-") or x:find"!") and self.ys or self.xs
209          push(xy, what) end end
210    else -- store another example; update column headers
211      push(self.egs, eg)
212      for at,x in pairs(eg) do if x ~= "?" then self.all[at]:add(x) end end end
213    return self end
214
215  function Sample:better(eg1,eg2,      e,n,a,b,s1,s2)
216    n,s1,s2,e = #self.ys, 0, 0, 2.71828
217    for _,num in pairs(self.ys) do
218      a  = num:norm(eg1[num.at])
219      b  = num:norm(eg2[num.at])
220      s1 = s1 - e^(num.w * (a-b)/n)
221      s2 = s2 - e^(num.w * (b-a)/n) end
222    return s1/n < s2/n end
223
224  function Sample:betters(egs)
225    return sort(egs or self.egs,function(a,b) return self:better(a,b) end) end
226
227  function Sample:clone(      inits,out)
228    out = Sample.new():add(self.names)
229    for _,eg in pairs(inits or {}) do out:add(eg) end
230    return out end
231
232  function Sample:dist(eg1,eg2,      a,b,d,n,inc)
233    d,n = 0,0
234    for _,col in pairs(self.xs) do
235      a,b = eg1[col.at], eg2[col.at]
236      inc = a=="?" and b=="?" and 1 or col:dist(a,b)
237      d   = d + inc^it.P
238      n   = n + 1 end
239    return (d/n)^(1/it.P) end
240
241  -- Report mid of the columns
242  function Sample:mid(cols)
243    return lap(cols or self.ys,function(col) return col:mid() end) end
244
245  -- Return spans of the column that most reduces variance
246  function Sample:bestSplits(tiny, trivials)
247    local function column1(col,    total,xpect,spans,total,xpect)
248      local function xpect1(span) return span.has.n/total * span.has:spread() end
249      spans = col:spans(self, tiny,trivials[col.at])
250      total = sum(spans,function(span) return span.has.n end)
251      xpect = sum(spans, xpect1)
252      return {xpect, spans}
253    end ------------------------------
254    return first(sort(lap(self.xs, column1), firsts))[2]   end
255
256  -- Split on column with best span, recurse on each split.
257  function Sample:tree(tiny,trivials,pre,        node,new,x)
258    pre=pre or ""
259    print(pre .."::"..#self.egs)
260    tiny    = tiny    or (#self.egs)^it.TINY
261    trivials = trivials or map(self.xs,
262                            function(_,x)
263                              return x.at,it.TRIVIAL*x:spread() end)
264    node    = {node=self, kids={}}
265    if #self.egs <= 2*tiny then print(333333);return node end
266    for _,span in pairs(self:bestSplits(tiny,trivials)) do
267      new = self:clone()
268      for _,eg in pairs(self.egs) do
269        x = eg[span.col.at]
270        if x=="?" or (span.lo <= x and x <= span.hi) then new:add(eg) end end
271      if #new.egs < #self.egs then
272        push(node.kids, {txt = span.col.txt, txt= span.col.at,
273                          lo = span.lo,       hi = span.hi,
274                          sub = new:tree(tiny,trivials,pre.."|.. ")}) end end
275      --os.exit()
276      --end end
277    return node end
278
279  -- Find which leaf best matches an example 'eg'.:w
280
281  function Sample:where(tree,eg,      max,x,default)
282    if #kid.has==0 then return tree end
283    max = 0
284    for _,kid in pairs(tree.node) do
285      if #kid.has > max then default,max = kid,#kid.has end
286      x = eg[kid.at]
287      if x ~= "?" then
288        if x <= kid.hi and x >= kid.lo then
289          return self:where(kid.has.eg) end end end
290    return self:where(default, eg) end
291
```

```lua
292  -- DISCRIMINATION ----------------------------------------------------
293  --
294  -- Input a list of {{x,y}..} values. Return spans that divide the 'x' values
295  -- to minimize variance on the 'y' values.
296  -- local div -- do not uncomment. 'div' was declared local above for 'Num:spans'.
297  local mergeable,merge
298
299  -- Return a list of 'spans' {lo=,hi=,col=col}.
300  -- Sort the list of pairs 'xys' then split it into 'spans' of cardinally at
301  -- least 'tiny'. Ensure that the max-min of each span is more that 'trivial'.
302  function div(xys, tiny, trivial,col,yklass)
303    xys    = sort(xys, function(a,b) return a.x < b.x end)
304    local tenth=#xys/10
305    trvial = trivial or it.TRIVIAL*math.abs(xys[9*tenth][1] - xys[tenth][1])/2.56
306    tiny   = tiny    or it.TINY*#xys
307    yklass = yklass  or Num
308    local spans,span
309    span  = {col=col,lo=xys[1].x, hi=xys[1].x, has=yklass()}
310    spans = {span}
311    for j,xy in pairs(xys) do
312      local x, y = xy.x, xy.y
313      if  j < #xys - tiny    and    -- enough items remaining after split
314          x ~= xys[j+1].x    and    -- next item is different (so can split here)
315          span.has.n > tiny  and    -- span has enough items
316          span.hi - span.lo > trivial -- span is not trivially small
317      then span = push(spans, {col=col, lo=span.hi, hi=x, has=yklass()})  -- then new span
318      end
319      span.hi = x
320      span.has:add(y) end
321    first(spans).lo = -math.huge
322    last(spans).hi  =  math.huge
323    return merge(spans) end
324
325  function mergeable(a,b,    new,b4)
326    new = a:merge(b)
327    b4  = (a.n*a:spread() + b.n*b:sd()) / new.n
328    if new:spread() <= b4 then return new end
329  end
330
331  --  Merge adjacent spans if the combo is simpler than the parts.
332  function merge(b4)
333    local j, tmp = 0, {}
334    while j < #b4 do
335      j = j + 1
336      local now, after = b4[j], b4[j+1]
337      if after then
338        local simpler = mergeable(now.has, after.has)
339        if simpler then
340          now = {col=col, lo=now.lo, hi= after.hi, has=simpler}
341          j = j + 1 end end
342      push(tmp,now) end
343    return #tmp==#b4 and b4 or merge(tmp) -- recurse until nothing merged
344  end
345
346
```

```lua
346  -- HINTING ------------------------------------------------------------------
347  --
348  -- Sorting on a few y values
349  local hints={}
350  function hints.default(eg) return eg end
351
352  function hints.sort(sample,scorefun,    test,train,egs,scored,small)
353    sample = Sample.new(it.FILE)
354    train,test = {}, {}
355    for i,eg in pairs(shuffle(sample.egs)) do
356      push(i<= it.TRAIN*#sample.egs and train or test, eg) end
357    egs = copy(train)
358    small = (#egs)^it.TINY
359    local i=0
360    scored = {}
361    while #egs >= small do
362      local tmp ={}
363      i = i + 1
364      io.stderr:write(fmt("%s",string.char(96+i)))
365      for j=1,it.HINTS do
366        egs[j] = (scorefun or hints.default)(egs[j])
367        push(tmp, push(scored, egs[j])) end
368      end
369      egs = hints.ranked(scored,egs,sample)
370      for i=1,it.CULL*#egs//1 do pop(egs) end
371    end
372    io.stderr:write("\n")
373    train=hints.ranked(scored, train, sample)
374    return #scored, sample:clone(train), sample:clone(test) end
375
376  function hints.ranked(scored,egs,sample,worker,   some)
377    function worker(eg) return {hints.rankOfClosest(scored,eg,sample),eg} end
378    scored = sample:betters(scored)
379    return  lap(sort(lap(egs, worker),firsts),second) end
380
381  function hints.rankOfClosest(scored,eg1,sample,          worker,closest)
382    function worker(rank,eg2) return {sample:dist(eg1,eg2),rank} end
383    closest = first(sort(map(scored, worker),firsts))
384    return  closest[2] end --+ closest[1]/10^8 end
385
```

```lua
386  -- demos --------------------------------------------------------------------
387  --
388  it._eg={}
389  it._no={}
390  function it._eg.shuffle(   t,u,v)
391    t={}
392    for i=1,32 do push(t,i) end
393    u = shuffle(copy(t))
394    v = shuffle(copy(t))
395    assert(#t == #u and u[1] ~= v[1]) end
396
397  function it._eg.lap()
398    assert(3==lap({1,2},function(x) return x+1 end)[2]) end
399
400  function it._eg.map()
401    assert(3==map({1,2},function(_,x) return x+1 end)[2]) end
402
403  function it._eg.tables()
404    assert(20==sort(shuffle({{10,20},{30,40},{40,50}}),firsts)[1][2]) end
405
406  function it._eg.csv(   n,z)
407    n=0
408    for eg in csv(it.FILE) do n=n+1; z=eg end
409    assert(n==399 and z[#z]==50) end
410
411  function it._eg.rnds(    t)
412    assert(10.2 == first(rnds({10.22,81.22,22.33},1))) end
413
414  function it._eg.sym(    s)
415    s=Sym{"a","a","a","a","b","b","c"}
416    assert("a"==s.mode) end
417
418  function it._eg.num1(    n)
419    n=Num{10,20,30,40,50,10,20,30,40,50,10,20,30,40,50}
420    assert(.375 == n:norm(25))
421    assert(15.625 == n:sd()) end
422
423  function it._eg.sample(    s,tmp,d1,d2,n)
424    s=Sample(it.FILE)
425    assert(2110 == last(s.egs)[s.all[4].at])
426    local sort1= s:betters(s.egs)
427    local lo, hi = s:clone(), s:clone()
428    for i=1,20             do lo:add(sort1[i]) end
429    for i=#sort1,#sort1-20,-1 do hi:add(sort1[i]) end
430    shout(s:mid())
431    shout(lo:mid())
432    shout(hi:mid())
433    for m,eg in pairs(sort1) do
434      n = bchop(sort1, eg,function(a,b) return s:better(a,b) end)
435      assert(m-n <=2) end end
436
437  function it._eg.dists(    s,tmp,d1,d2,n)
438    s=Sample(it.FILE)
439    tmp = sort(lap(shuffle(s.egs),
440                   function(eg2) return {s:dist(eg2,s.egs[1]), eg2} end),
441               firsts)
442    d1=s:dist(tmp[1][2], tmp[10][2])
443    d2=s:dist(tmp[1][2], tmp[#tmp][2])
444    assert(d1*10 < d2) end
445
446  function it._eg.hints(    s,_,__,evals,sort1,train,test,n)
447    s = Sample(it.FILE)
448    evals, train,test = hints.sort(s)
449    test.egs = test:betters()
450    for m,eg in pairs(test.egs) do
451      n = bchop(train.egs, eg,function(a,b) return s:better(a,b) end); end end
452
453  function it._eg.dump()
454    shout(it) end
455  function it._eg.tree(    s,t,u,eg1,evals,ordered,rest)
456    s = Sample(it.FILE)
457    t = copy(s.names)
458    push(t,"Rank!")
459    u = Sample.new():add(t)
460    evals, ordered, rest = hints.sort(s)
461    for m,eg in pairs(ordered.egs) do
462      eg1 = copy(eg)
463      push(eg1,m)
464      u:add(eg1) end
465    print(1)
466    u:tree() end
467
468  -- START-UP -----------------------------------------------------------------
469  --
470  it{demos=it._eg, nervous=true}
471
```

```
472  --[[
473
474       -| - _    _| _
475        | (_)   (_| (_)
476
477  Spans
478   Little languages:
479      - options
480      - data language
481
482  Lesson plan
483  - w1: ssytems: github. github workplaces. unit tests. doco tools.
484
485  - w2: num,sym
486  - W3: sample
487  - w4: eval, knn, unfarinessness
488  - W5:
489
490  -  seems to be  a revers that i  need to do .... but dont
491  - check if shuffle is working
492
493  teaching:
494  - sample is v.useful
495  --]]
```