

```

1  #!/usr/bin/env lua
2  local b4={}; for k,v in pairs(_ENV) do b4[k]=v end; --[[
3
4  a little lile
5  zwa learning
6  library
7
8
9
10
11
12
13
14
15
16
17
18 --]] local options={
19
20 what = "Small sample multi-objective optimizer.",
21 usage= "(c) 2021 Tim Menzies <timm@ieee.org> unlicense.org",
22 about= [[
23 Sort N examples on multi-goals using a handful of 'hints'; i.e.
24
25 - Evaluate and rank, a few examples (on their y-values);
26 - Sort other examples by x-distance to the ranked ones;
27 - Recurse on the better half (so we sample more and more
28   from the better half, then quarter, then eighth...).
29
30 A regression tree learner then explores the examples (sorted
31 left to right, worst to best). By finding branches that
32 reduce the variance of the index of those examples, this
33 tree reports what attribute ranges select for the better (or
34 worse) examples. ]],
35
36 how= {{{"file",      "-f",      "../data/auto93.csv",   "read data from file"},
37        {"cull",       "-c",       .5,                   "cuts per repeat"},
38        {"help",       "-h",       false,                 "show help"},
39        {"hints",      "-H",       4,                     "hints per generation"},
40        {"p",          "-p",       2,                     "distance calc exponent"},
41        {"small",      "-s",       .5,                   "div list t into t^small"},
42        {"seed",       "-S",       10019,                 "random number seed"},
43        {"train",      "-t",       .5,                   "size of training set"},
44        {"trivial",    "-T",       .35,                  "small delta-trivial^sd"},
45        {"todo",       "-t",       "all",                 "run unit test, or 'all'"},
46        {"wild",       "-W",       false,                 "run tests, no protection"} }}
47
48 local Seed,cli,the
49 Seed=10019
50 -- If '-x X' appears on command line and '-x default' is in 'how'
51 -- then update default from the command line (and if 'default'
52 -- is false, then set it to true. Also, maybe
53 -- set random number seed and maybe show help string.
54 function cli(opt, u)
55   u={}
56   for _,t in pairs(opt.how) do
57     u[t[1]] = t[3]
58     for n,word in ipairs(arg) do if word==t[2] then
59       u[t[1]] = t[3] and (tonumber(arg[n+1]) or arg[n+1]) or true end end end
60   if u.help
61   then print(string.format("\n%s [OPTIONS]n%s\n%s\n\nOPTIONS:n",
62                             arg[0],opt.usage,opt.what))
63   else
64     for _,t in pairs(opt.how) do print(string.format("%4s %-9s%t%s %s",
65                                                       t[2], t[3] and t[1] or "", t[4], t[3] and "=" or "", t[3] or "")) end
66     print("\n"..opt.about)
67     os.exit()
68   end
69   if u.seed then Seed = u.seed end
70   return u end
71
72 -- Make a global for our options e.g. the = {seed=10019, help=false, p=2...}
73 the = cli(options)

```

```

74 --
75 --
76 --
77 --
78 --
79 -- ## Table Stuff
80 local randi -- defined later, needed now in "shuffle"
81 local cat,map,lap,keys, last,copy,pop,push,sort,firsts,first,second,shuffle,bchop
82 -- Table to string.
83 cat = table.concat
84 -- Return a sorted table.
85 sort = function(t,f) table.sort(t,f); return t end
86 -- Add to end, pull from end.
87 push = table.insert
88 pop = table.remove
89 -- Return first,second, last item.
90 first = function(t) return t[1] end
91 second = function(t) return t[2] end
92 last = function(t) return t[#t] end
93 -- Function for sorting pairs of items.
94 firsts = function(a,b) return first(a) < first(b) end
95
96 -- Random order of items in a list (sort in place).
97 function shuffle(t, j)
98   for i=#t,2,-1 do j=randi(1,i); t[i],t[j]=t[j],t[i] end; return t end
99
100 -- Collect values, passed through 'f'.
101 function lap(t,f) return map(t,f,1) end
102 -- Collect key,values, passed through 'f'.
103 -- If 'f' returns two values, store as key,value.
104 -- If 'f' returns one values, store at index value.
105 function map(t,f,one, u)
106   u={}
107   for x,y in pairs(t) do
108     if one then x,y=f(y) else u,x=f(x,y) end
109     if x ~= nil then
110       if y then u[x]=y else u[1+#u]=x end end end
111   return u end
112
113 -- Return a table's keys (sorted).
114 function keys(t,u)
115   u={}
116   for k,_ in pairs(t) do if tostring(k):sub(1,1)~="_" then push(u,k) end end
117   return sort(u) end
118
119 -- Binary chop (assumes sorted lists)
120 function bchop(t,val,lt,lo,hi, mid)
121   lt = lt or function(x,y) return x < y end
122   lo,hi = lo or 1, hi or #t
123   while lo <= hi do
124     mid = (lo+hi) // 2
125     if lt(t[mid],val) then lo=mid+1 else hi= mid-1 end end
126   return math.min(lo,#t) end
127
128 -----
129 -- ## Maths Stuff
130 local abs,norm,sum,rnd,rnds,Seed,rand
131 abs = math.abs
132 -- Round 'x' to 'd' decimal places.
133 function rnd(x,d, n) n=10^(d or 0); return math.floor(x*n+0.5) / n end
134 -- Round list of items to 'd' decimal places.
135 function rnds(t,d) return lap(t, function(x) return rnd(x,d or 2) end) end
136
137 -- Sum items, filtered through 'f'.
138 function sum(t,f)
139   f= f or function(x) return x end
140   out=0; for _,x in pairs(f) do out = out + f(x) end; return out end
141
142 Seed=937162211
143 function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
144 function rand(lo,hi, mult,mod)
145   lo,hi = lo or 0, hi or 1
146   Seed = (16807 * Seed) % 2147483647
147   return lo + (hi-lo) * Seed / 2147483647 end
148
149 -----
150 -- ## Printing Stuff
151 local out,shout,red,green,yellow,blue,color,fmt
152 fmt = string.format
153 -- Print as red, green, yellow, blue.
154 function color(s,n) return fmt("%27[1m27[%sm%s27[0m",n,s) end
155 function red(s) return color(s,31) end
156 function green(s) return color(s,32) end
157 function yellow(s) return color(s,34) end
158 function blue(s) return color(s,36) end
159
160 -- Printed string from a nested structure.
161 shout = function(x) print(out(x)) end
162 -- Generate string from a nested structures
163 -- (and don't print any contents more than once).
164 function out(t,seen, u,key,value,public)
165   function key(k) return fmt("%s%s",blue(k),out(t[k],seen)) end
166   function value(v) return out(v,seen) end
167   if type(t) == "function" then return "(...)" end
168   if type(t) ~= "table" then return tostring(t) end
169   seen = seen or {}
170   if seen[t] then return "..." else seen[t] = t end
171   u = #t>0 and lap(t, value) or lap(keys(t), key)
172   return red((t._is or "")..{"("}.cat(u,"").red(")") end
173
174 -----
175 -- ## File i/o Stuff
176 -- Return one table per line, split on commas.
177 local csv
178 function csv(file, line)
179   file = io.input(file)
180   line = io.read()
181   return function( t,tmp)
182     if line then
183       t={}
184       for cell in line:gsub("[\r\n]*",""):gsub("#.*",""):gmatch("(^[^\,]+)") do
185         push(t, tonumber(cell) or cell) end
186       line = io.read()
187       if #t>0 then return t end
188       else io.close(file) end end end
189
190 -----
191 -- ## OO Stuff
192 local has,obj
193 -- Create an instance
194 function has(mt,x) return setmetatable(x,mt) end
195
196 -- Create a class
197 function obj(s, o,new)
198   o = {_is=s, __tostring=out}
199   o._index = o
200   return setmetatable(o,{__call = function(_,...) return o.new(...) end}) end
201

```

```

202 -----
203 -- ## Stuff for tracking 'Sym'bol Counts.
204 -- 'Sym's track symbol counts and the 'mode' (most frequent symbol).
205 local Sym=obj"Sym"
206 function Sym.new(inits, self)
207   self=has(Sym,{has={}, n=0, mode=nil, most=0})
208   for _,one in pairs(inits or {}) do self:add(one) end
209   return self end
210
211 function Sym:add(x)
212   self.n = self.n + 1
213   self.has[x] = 1 + (self.has[x] or 0)
214   if self.has[x] > self.most then self.most, self.mode = self.has[x], x end end
215
216 function Sym:dist(a,b) return a==b and 0 or 1 end
217 function Sym:mid() return self.mode end
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 -----
277 -- discretization tricks
278 local splits={}
279 function splits.best(sample, best,tmp,xpect,out)
280   best = maths.huge
281   for _,x in pairs(sample.xls) do
282     tmp, xpect = splits.whatif(x.at,self)
283     if xpect < best
284     then out,best = tmp,xpect end end
285   return out end
286
287 function splits.whatif(col,sample, out)
288   out = splits.spans(col,sample)
289   xpect = sum(out, function(x) return x.has.n*x:sd() end)/#sample.egs
290   out = map(out, function(_,x) x.has=x.has:all(); x.col= col end)
291   return out, xpect end
292
293 function splits.spans(col,sample, xs, symbolic,x)
294   xys,xs, symbolic = {}, Num(), sample.nums[col]
295   for rank,eg in pairs(sample.egs) do
296     x = eg[col]
297     if x ~="?" then
298       xs:add(x)
299       if symbolic
300       then -- in symbolic columns, xys are the indexes seen with each symbol
301         xys[x] = xys[x] or {}
302         push(xys[x], rank)
303       else -- in numeric columns, xys are each number paired with its row id
304         push(xys, {x=x,y=rank}) end end
305   end
306   if symbolic
307   then return map(xys, function(x,t) return {lo=x, hi=x, has=Num(t)} end)
308   else return splits.merge(
309     splits.div(xys, #xs*the.small, sd(sort(xs))*the.trivial)) end end
310
311 -- Generate a new range when
312 -- 1. there is enough left for at least one more range; and
313 -- 2. the lo,hi delta in current range is not boringly small; and
314 -- 3. there are enough x values in this range; and
315 -- 4. there is natural split here
316 -- Fuse adjacent ranges when:
317 -- 5. the combined class distribution of two adjacent ranges
318 -- is just as simple as the parts.
319 function splits.div(xys, tiny, dull, now,out,x,y)
320   xys = sort(xys, function(a,b) return a.x < b.x end)
321   now = {lo=xys[1].x, hi=xys[1].x, has=Num()}
322   out = {now}
323   for j,xy in pairs(xys) do
324     x, y = xy.x, xy.y
325     if j<#xys-tiny and x==xys[j+1].x and now.has.n>tiny and now.hi-now.lo>dull
326     then now = {lo=x, hi=x, has=Num()}
327     push(out, now) end
328     now.hi = x
329     now.has:add(y) end
330   return out end
331
332 function splits.merge(b4, j,tmp,a,n,simpler)
333   j, n, tmp = 0, #b4, {}
334   while j<n do
335     j = j + 1
336     a = b4[j]
337     if j < n-1 then
338       simpler = a.has:mergeable(b4[j+1].has)
339       if simpler then
340         j = j + 1
341         a = {lo=a.lo, hi= b4[j+1].hi, has=simpler} end end
342     push(tmp,a) end
343   return #tmp==#b4 and b4 or merge(tmp) end
344
345

```

```

345 -----
346 -- Samples store examples. Samples know about
347 -- (a) lo,hi ranges on the numerics
348 -- and (b) what's there, independent 'x' or dependent 'y' columns.
349 local Sample=obj"Sample"
350 function Sample.new(src,self)
351   self = has(Sample,{names=nil, all={}, ys={}, xs={}, egs={}})
352   if src then
353     if type(src)=="string" then for x in csv(src) do self:add(x) end end
354     if type(src)=="table" then for _,x in pairs(src) do self:add(x) end end end
355   return self end
356
357 function Sample:clone(inits,out)
358   out = Sample.new():add(self.names)
359   for _,eg in pairs(inits or {}) do out:add(eg) end
360   return out end
361
362 function Sample:add(eg, name,datum)
363   function name(col,new, weight, where, what)
364     if new:find"." then return end
365     weight= new:find"." and -1 or 1
366     what = (col=col, w=weight, seen=(new:match("[A-Z]",x) and Num() or Sym()))
367     where = (new:find("+") or new:find("-")) and self.ys or self.xs
368     push(self.all, what)
369     push(where, what)
370   end
371   function datum(one,new)
372     if new ~="?" then one.seen:add(new) end
373   end
374   if not self.names
375   then self.names = eg
376   map(eg, function(col,x) name(col,x) end)
377   else push(self.egs, eg)
378   map(self.all, function(_,col) datum(col,eg[col.col]) end)
379   end
380   return self end
381
382 function Sample:better(eg1,eg2, e,n,a,b,s1,s2)
383   n,s1,s2,e = #self.ys, 0, 0, 2.71828
384   for _,num in pairs(self.ys) do
385     a = num.seen:norm(eg1[num.col])
386     b = num.seen:norm(eg2[num.col])
387     s1 = s1 - e^(num.w * (a-b)/n)
388     s2 = s2 - e^(num.w * (b-a)/n) end
389   return s1/n < s2/n end
390
391 function Sample:betters(egs)
392   return sort(egs or self.egs,function(a,b) return self:better(a,b) end) end
393
394 function Sample:dist(eg1,eg2, a,b,d,n,inc)
395   d,n = 0,0
396   for _,x in pairs(self.xls) do
397     a,b = eg1[x.col], eg2[x.col]
398     inc = a=="?" and b=="?" and 1 or x.seen:dist(a,b)
399     d = d + inc*the.p
400     n = n + 1 end
401   return (d/n)^(1/the.p) end
402
403 function Sample:stats(cols)
404   return lap(cols or self.ys,function(col) return col.seen:mid() end) end
405
406 -- bins his
407 -- bins sorts
408 function Sample:tree(min, node,min,sub)
409   node = {node=self, kids={}}
410   min = min or (#self.egs)*the.small
411   if #self.egs >= 2*min then
412     -- here
413     for _,span in pairs(splits.best(sample)) do
414       sub = self:clone()
415       for _,at in pairs(span.has) do sub:add(self.egs[at]) end
416       push(node.kids, span)
417       span.has = sub:tree(min) end end
418   return node end
419
420 -- at node
421 function Sample:where(tree,eg, max,x,default)
422   if #kid.has==0 then return tree end
423   max = 0
424   for _,kid in pairs(tree.node) do
425     if #kid.has > max then default,max = kid,#kid.has end
426     x = eg[kid.col]
427     if x ~="?" then
428       if x <= kid.hi and x >= kid.lo then
429         return self:where(kid.has.eg) end end end
430   return self:where(default, eg) end
431
432

```

```

433 -----
434 -- sample sample sorting
435 local hints={}
436 function hints.default(eg) return eg end
437
438 function hints.sort(sample,score, test,train,evals)
439 sample = Sample.new(the.file)
440 train,test = {}, {}
441 for i,eg in pairs(shuffle(sample.egs)) do
442   push(i<= the.train*#sample.egs and train or test, eg) end
443   evals,train = hints.recurse(sample, train,0,
444     score or hints.default, {}, (#train^the.small)
445   return evals,sample:clone(train), sample:clone(test) end
446
447 function hints.recurse(sample, eggs, evals, scorefun, out, small, worker)
448   if #egs < small then
449     for i=1, #egs do push(out, pop(egs)) end
450     return evals,out
451   end
452   local scores = {}
453   function worker(eg) return hints.locate(scores,eg,sample) end
454   for j=1,the.hints do evals=evals+1;
455     push(scores, scorefun(pop(egs))) end
456   scores = sample:betters(scores)
457   eggs = lap(sort(lap(egs, worker),firsts),second)
458   print(the.cull*#egs)
459   for i=1,the.cull*#egs//1 do push(out, pop(egs)) end
460   return hints.recurse(sample, eggs,evals, scorefun, out, small) end
461
462 function hints.locate(scores,eg,sample, closest,rank,tmp)
463   closest, rank, tmp = IE32, IE32, nil
464   for rank0, scored in pairs(scores) do
465     tmp = sample:dist(eg, scored)
466     if tmp < closest then closest,rank = tmp,rank0 end end
467   return {rank, eg} end
468   --return {rank+closest/10^6, eg} end
469
470
471

```

```

472 --
473
474
475
476 -----
477 local eg,fail,example={},0
478 function example(k, f,ok,msg)
479   f= eg[k]; assert(f,"unknown action"..k)
480   the=c1i(options)
481   if the.wild then return f() end
482   ok,msg = pcall(f)
483   if ok then print(green("PASS"),k)
484   else print(red("FAIL"), k,msg); fail=fail+1 end end
485
486 function eg.shuffle( t)
487   t={}
488   for i=1,32 do push(t,i) end
489   assert(#t == #shuffle(t) and t[1] ~= shuffle(t)[1]) end
490
491 function eg.lap()
492   assert(3==lap({1,2},function(x) return x+1 end)[2]) end
493
494 function eg.map()
495   assert(3==map({1,2},function(_,x) return x+1 end)[2]) end
496
497 function eg.tables()
498   assert(20==sort(shuffle({{10,20},{30,40},{40,50}}),firsts)[1][2]) end
499
500 function eg.csv( n,z)
501   n=0
502   for eg in csv(the.file) do n=n+1; z=eg end
503   assert(n==399 and z[#z]==50) end
504
505 function eg.rnds( t)
506   assert(10.2 == first(rnds({10.22,81.22,22.33},1))) end
507
508 function eg.sym( s)
509   s=Sym("a","a","a","a","b","b","c")
510   assert("a"==s.mode) end
511
512 function eg.num1( n)
513   n=Num(10,20,30,40,50,10,20,30,40,50,10,20,30,40,50)
514   assert(.375 == n:norm(25))
515   assert(15.625 == n:sd()) end
516
517 function eg.num2( n1,n2,n3,n4)
518   n1=Num(10,20,30,40,50,10,20,30,40,50,10,20,30,40,50)
519   n2=Num(10,20,30,40,50,10,20,30,40,50,10,20,30,40,50)
520   assert(n1:mergeable(n2)==nil)
521   n3=Num(10,20,30,40,50,10,20,30,40,50,10,20,30,40,50)
522   n4=Num(100,200,300,400,500,100,200,300,400,500,100,200,300,400,500)
523   assert(n3:mergeable(n4)==nil) end
524
525 function eg.sample( s,tmp,d1,d2,n)
526   s=Sample(the.file)
527   assert(2110 == last(s.egs)[s.all[3].col])
528   local sortl= s:betters(s.egs)
529   local lo, hi = s:clone(), s:clone()
530   for i=1,20 do lo:add(sortl[i]) end
531   for i=#sortl,#sortl-30,-1 do hi:add(sortl[i]) end
532   shout(s:stats())
533   shout(lo:stats())
534   shout(hi:stats())
535   for m,eg in pairs(sortl) do
536     n = bchop(sortl, eg,function(a,b) return s:better(a,b) end)
537     assert(m-n <=2) end end
538
539 function eg.dists( s,tmp,d1,d2,n)
540   s=Sample(the.file)
541   tmp = sort(lap(shuffle(s.egs),
542     function(eg2) return {s:dist(eg2,s.egs[1]), eg2} end),
543     firsts)
544   d1=s:dist(tmp[1][2], tmp[10][2])
545   d2=s:dist(tmp[1][2], tmp[#tmp][2])
546   assert(d1*10<d2)
547 end
548
549 function eg.hints( s,_,_,evals,sortl,train)
550   s=Sample(the.file)
551   --for _,eg in pairs(sortl) do lap(s.ys, function(col) return eg[col.col] end ) end
552   -- assert(s.ys[4].lo==1613)
553   evals, train,test = hints.sort(s)
554   print("=",evals)
555   test.egs = test:betters()
556   for m,eg in pairs(test.egs) do
557     n = bchop(train.egs, eg,function(a,b) return s:better(a,b) end)
558     print(m,n) end
559   end
560
561 if the.todo=="all" then lap(keys(eg),example) else example(the.todo) end
562
563 -----
564 -- trick for checking for rogues.
565 for k,v in pairs(_ENV) do if not b4[k] then print("?rogue: ",k,type(v)) end end
566 os.exit(fail)
567
568
569

```

```
570 --[[
571 -- seems to be a revers that i need to do .... but dont
572
573 teaching:
574 - sample is v.useful
575
576
577 --]]
```