

```

1  #!/usr/bin/env lua
2  --
3  --
4  -- a little lile
5  -- ZWA learning
6  -- library
7  --
8  --
9  --
10 --
11 --
12 --
13 --
14 --
15 --
16 --
17 --
18 local it=require"z"
19 what = "Small sample multi-objective optimizer.",
20 who = "(c) 2021 Tim Menzies <tmn@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", "-f", ".data/au93.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 {"SMALL", "-s", ".5", "div list into '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 abs,bchop,cat,copy = it.abs, it.bchop, it.cat, it.copy
48 local csv,first,firsts,fmt,has = it.csv, it.first, it.firsts, it.fmt, it.has
49 local keys,last,lap,map,obj = it.keys, it.last, it.lap, it.map, it.obj
50 local out,pop,push,rand,rnd = it.out, it.pop, it.push, it.rand, it.rnd
51 local rnds,rogues,second = it.rnds, it.rogues, it.second, it.shuffle
52 local sort,sum,top = it.sort, it.sum, it.top
53
54 --[[
55 Spans
56 Little languages:
57 - options
58 - data language
59
60 Lesson plan
61 -- w1: ssystems: github. github workplaces. unit tests. doco tools.
62 -- w2: num,sym
63 -- w3: sample
64 -- w4: eval, knn, unfairnessness
65 -- w5:
66 --]]
67

```

```

68 --
69 --
70 --
71 --
72 -- ## Stuff for tracking 'Num'bers.
73 -- 'Num's track a list of number, and can report it sorted.
74 local Num=obj"Num"
75 function Num.new(inits,at,txt, self)
76   self=has(Num,(at=at or 0, txt=txt or "", w=(txt or ""):find("-" and -1 or 1,
77     has={}, n=0, lo=1E32, hi=1E-32, ready=true))
78   for _,one in pairs(inits or {}) do self:add(one) end
79   return self end
80
81 function Num:add(x)
82   if x>self.hi then self.hi = x
83   elseif x<self.lo then self.lo = x end
84   push(self.has,x); self.n=self.n+1; self.ready=false end
85
86 -- Ensure that the returned list of numbers is sorted.
87 function Num:all(x)
88   if not self.ready then table.sort(self.has) end
89   self.ready = true
90   return self.has end
91
92 function Num:dist(a,b)
93   if a=="?" then b=self:norm(b); a = b*.5 and 0 or 1
94   elseif b=="?" then a=self:norm(a); b = a*.5 and 0 or 1
95   else a,b = self:norm(a), self:norm(b) end
96   return abs(a-b) end
97
98 -- Combine two 'num's.
99 function Num:merge(other, new)
100   new = Num.new(self.has)
101   for _,x in pairs(other.has) do new:add(x) end
102   return new end
103
104 -- Return a merged item if that combination
105 -- is simpler than its parts.
106 function Num:mergeable(other, new,b4)
107   new = self:merge(other)
108   b4 = (self.n*self:sd() + other.n*other:sd()) / new.n
109   if b4 >= new:sd() then return new end end
110
111 -- The 'mid' is the 50th percentile.
112 function Num:mid() return self:per(.5) end
113
114 -- Return 'x' normalized 0..1, lo..hi.
115 function Num:norm(x, lo,hi)
116   if x=="?" then return x end
117   lo,hi = self.lo, self.hi
118   return abs(hi - lo) < 1E-32 and 0 or (x - lo)/(hi - lo) end
119
120 -- Return the 'p'-th percentile number.
121 function Num:per(p, t)
122   t = self:all()
123   p = p*#t//1
124   return #t<2 and t[1] or t[p < 1 and 1 or p>#t and #t or p] end
125
126 -- The 10th to 90th percentile range is 2.56 times the standard deviation.
127 function Num:sd() return (self:per(.9) - self:per(.1))/ 2.56 end
128
129 -- Create one span holding row indexes associated with each number
130 local div -- defined below
131 function Num:spans(egs, lo,spans,fin)
132   local xys,xs = {}, Num()
133   for pos,eg in pairs(egs) do
134     local x = eg[self.at]
135     if x ~= "?" then
136       xs:add(x)
137       push(xys, {x=x,y=pos}) end end
138   return div(xys -- split xys into spans...
139     ,math.min(10,xs.n*it.SMALL) -- ..where spans are of size sqrt(#xs)..
140     ,xs:sd()*it.TRIVIAL) end -- ..and spans have (last-first)>trivial
141
142 -----
143 -- ## Stuff for tracking 'Sym'bol Counts.
144 -- 'Sym's track symbol counts and the 'mode' (most frequent symbol).
145 local Sym=obj"Sym"
146 function Sym.new(inits,at,txt, self)
147   self=has(Sym,(at=at or 0, txt=txt or "", has={}, n=0, mode=nil, most=0))
148   for _,one in pairs(inits or {}) do self:add(one) end
149   return self end
150
151 function Sym:add(x)
152   self.n = self.n + 1
153   self.has[x] = 1 + (self.has[x] or 0)
154   if self.has[x] > self.most then self.most, self.mode = self.has[x], x end end
155
156 function Sym:dist(a,b) return a==b and 0 or 1 end
157 function Sym:mid() return self.mode end
158
159 -- Create one span holding row indexes associated with each symbol
160 function Sym:spans(egs, xys,x)
161   xys = {}
162   for pos,eg in pairs(egs) do
163     x = eg[self.at]
164     if x ~= "?" then
165       xys[x] = xys[x] or {}
166       push(xys[x], pos) end end
167   return map(xys, function(x,t) return {lo=x, hi=x, has=Num(t)} end) end end
168
169 -----
170 -- ## Stuff for skipping all things sent to a column
171 local Skip=obj"Skip"
172 function Skip.new(_,at,txt) return has(Skip,(at=at or 0, txt=txt or "", n=0)) end
173 function Skip:add(x) self.n = self.n + 1; return x end
174 function Skip:mid() return "?" end
175

```

Dec 11, 21 21:01

I5.lua

Page 3/6

```

176 --
177 --
178 --
179
180 -- Samples store examples. Samples know about
181 -- (a) lo,hi ranges on the numerics
182 -- and (b) what are independent 'x' or dependent 'y' columns.
183 local Sample = obj{"Sample"}
184 function Sample.new(      src,self)
185   self = has(Sample,{names=nil, all={}, ys={}, xs={}, eggs={}})
186   if src then
187     if type(src)=="string" then for x in csv(src) do self:add(x) end end
188     if type(src)=="table" then for _,x in pairs(src) do self:add(x) end end end
189   return self end
190
191 function Sample:add(eg,      ako,what,where)
192   if not self.names
193   then -- create the column headers
194     self.names = eg
195     for at,x in pairs(eg) do
196       ako = x:find"*" and Skip or x:match("[A-Z]" and Num or Sym
197       what = push(self.all, ako({}, at, x))
198       if not x:find"." then
199         where = (x:find(".*") or x:find("-")) and self.ys or self.xs
200         push(where, what) end end
201     else -- store another example; update column headers
202       push(self.egs, eg)
203       for at,x in pairs(eg) do if x ~= "?" then self.all[at]:add(x) end end end
204     return self end
205
206 function Sample:better(eg1,eg2,      e,n,a,b,s1,s2)
207   n,s1,s2,e = #self.ys, 0, 0, 2.71828
208   for _,num in pairs(self.ys) do
209     a = num:norm(eg1[num.at])
210     b = num:norm(eg2[num.at])
211     s1 = s1 - e^(num.w * (a-b)/n)
212     s2 = s2 - e^(num.w * (b-a)/n) end
213   return s1/n < s2/n end
214
215 function Sample:betters(egs)
216   return sort(egs or self.egs,function(a,b) return self:better(a,b) end) end
217
218 function Sample:clone(      inits,out)
219   out = Sample.new():add(self.names)
220   for _,eg in pairs(inits or {}) do out:add(eg) end
221   return out end
222
223 function Sample:dist(eg1,eg2,      a,b,d,n,inc)
224   d,n = 0,0
225   for _,col in pairs(self.xs) do
226     a,b = eg1[col.at], eg2[col.at]
227     inc = a=="?" and b=="?" and 1 or col:dist(a,b)
228     d = d + inc*it.P
229     n = n + 1 end
230   return (d/n)^(1/it.P) end
231
232 -- Report mid of the columns
233 function Sample:mid(cols)
234   return lap(cols or self.ys,function(col) return col:mid() end) end
235
236 -- Return spans of the column that most reduces variance
237 function Sample:splitter(cols)
238   function worker(col) return self:splitter1(col) end
239   return first(sort(lap(cols or self.xs, worker), firsts))[2] end
240
241 -- Return a column's spans, and the expected sd value of those spans.
242 function Sample:splitter1(col,      spans,xpect)
243   spans = col:spans(self.egs)
244   --spans = lap(spans, shout)
245   xpect = sum(spans,
246     function(span) span.col=col; return span.has.n*span.has.sd() end)
247   return (xpect/#self.egs, spans) end
248
249 -- Split on column with best span, recurse on each split.
250 function Sample:tree(min,      node,min,sub,splitter, splitter1)
251   node = (node=self, kids={})
252   min = min or (#self.egs)*it.SMALL
253   if #self.egs >= 2*min then
254     for _,span in pairs(self:splitter()) do
255       sub = self:clone()
256       for _,at in pairs(span.has) do sub:add(self.egs[at]) end
257       push(node.kids, span)
258       span.has = sub:tree(min) end end
259   return node end
260
261 -- Find which leaf best matches an example 'eg'.
262 function Sample:where(tree,eg,      max,x,default)
263   if #kid.has==0 then return tree end
264   max = 0
265   for _,kid in pairs(tree.node) do
266     if #kid.has > max then default,max = kid,#kid.has end
267     x = eg[kid.at]
268     if x == "?" then
269       if x <= kid.hi and x >= kid.lo then
270         return self:where(kid.has.eg) end end end
271   return self:where(default, eg) end
272
273 -----
274 -- Discretization tricks
275 -- Input a list of {(x,y)..} values. Return spans that divide the 'x' values
276 -- to minimize variance on the 'y' values.
277 function div(xys, tiny, dull,      merge,coverGaps)
278   function merge(b4) -- merge adjacent spans if combo simpler to he parts
279     local j, tmp = 0, {}
280     while j < #b4 do
281       j = j + 1
282       local now, after = b4[j], b4[j+1]
283       if after then
284         local simpler = now.has:mergeable(after.has)
285         if simpler then
286           now = {lo=now.lo, hi=after.hi, has=simpler}
287           j = j + 1 end end
288       push(tmp,now) end
289   return #tmp==#b4 and b4 or merge(tmp) -- recurse until nothing merged
290 end
291
292 function coverGaps(spans,      b4) -- cover gaps in number line
293   spans[1].lo = -math.huge
294   spans[#spans].hi = math.huge
295   b4 = spans[1].hi
296   for _,span in pairs(spans) do span.lo = b4; b4 = span.hi end
297   return spans
298
299 local spans,span
300 xys = sort(xys, function(a,b) return a.x < b.x end)
301 span = {lo=xys[1].x, hi=xys[1].x, has=Num()}
302 spans = {span}
303 for j,xy in pairs(xys) do
304   local x, y = xy.x, xy.y
305   if j < #xys - tiny and -- enough items remaining after split
306   x ~= xys[j+1].x and -- next item is different (so can split here)
307   span.has.n > tiny and -- span has enough items
308   span.hi - span.lo > dull -- span is not trivially small
309   then span = push(spans, {lo=x, hi=x, has=Num()}) -- then new span
310   end
311   span.hi = x
312   span.has:add(y) end
313 return coverGaps(merge(spans)) end

```

Dec 11, 21 21:01

I5.lua

Page 4/6

```

314 --
315 --
316 --
317
318 -- Sorting on a few y values
319 local hints={}
320 function hints.default(eg) return eg end
321
322 function hints.sort(sample,scorefun,      test,train,egs,scored,small)
323   sample = Sample.new(it.FILE)
324   train,test = {}, {}
325   for i,eg in pairs(shuffle(sample.egs)) do
326     push(i<= it.TRAIN*#sample.egs and train or test, eg) end
327   egs = copy(train)
328   small = (#egs)*it.SMALL
329   local i=0
330   scored = {}
331   while #egs >= small do
332     local tmp={}
333     i = i + 1
334     io.stderr:write(fmt("%s",string.char(96+i)))
335     for j=1,it.HINTS do
336       egs[j] = (scorefun or hints.default)(egs[j])
337       push(tmp, push(scored, egs[j]))
338     end
339     egs = hints.ranked(scored,egs,sample)
340     for i=1,it.CULL*#egs/1 do pop(egs) end
341   end
342   io.stderr:write("\n")
343   train=hints.ranked(scored, train, sample)
344   return #scored, sample:clone(train), sample:clone(test) end
345
346 function hints.ranked(scored,egs,sample,worker,      some)
347   function worker(eg) return (hints.rankOfClosest(scored,eg,sample),eg) end
348   scored = sample:betters(scored)
349   return lap(sort(lap(egs, worker),firsts),second) end
350
351 function hints.rankOfClosest(scored,eg1,sample,      worker,closest)
352   function worker(rank,eg2) return {sample:dist(eg1,eg2),rank} end
353   closest = first(sort(map(scored, worker),firsts))
354   return closest[2] end --> closest[1]/10^8 end
355

```

Dec 11, 21 21:01

l5.lua

Page 5/6

```

355 -- [[ a m n o s
356
357
358 it.eg={}
359 function it.eg.shuffle( t,u,v)
360 t={}
361 for i=1,32 do push(t,i) end
362 u = shuffle(copy(t))
363 v = shuffle(copy(t))
364 assert(#t == #u and u[1] ~= v[1]) end
365
366 function it.eg.lap()
367 assert(3==lap({1,2},function(x) return x+1 end)[2]) end
368
369 function it.eg.map()
370 assert(3==map({1,2},function(_,x) return x+1 end)[2]) end
371
372 function it.eg.tables()
373 assert(20==sort(shuffle({{10,20},{30,40},{40,50}}),firsts)[1][2]) end
374
375 function it.eg.csv( n,z)
376 n=0
377 for eg in it.csv(it.FILE) do n=n+1; z=eg end
378 assert(n==399 and z[#z]==50) end
379
380 function it.eg.rnds( t)
381 assert(10.2 == first(rnds({10.22,81.22,22.33},1))) end
382
383 function it.eg.sym( s)
384 s=sym("a","a","a","a","b","b","b","c")
385 assert("#a"==s.mode) end
386
387 function it.eg.num1( n)
388 n=Num(10,20,30,40,50,10,20,30,40,50,10,20,30,40,50)
389 assert(.375 == n:norm(25))
390 assert(15.625 == n:sd()) end
391
392 function it.eg.num2( n1,n2,n3,n4)
393 n1=Num(10,20,30,40,50,10,20,30,40,50,10,20,30,40,50)
394 n2=Num(10,20,30,40,50,10,20,30,40,50,10,20,30,40,50)
395 assert(n1:mergeable(n2)==nil)
396 n3=Num(10,20,30,40,50,10,20,30,40,50,10,20,30,40,50)
397 n4=Num(100,200,300,400,500,100,200,300,400,500,100,200,300,400,500)
398 assert(n3:mergeable(n4)==nil) end
399
400
401 function it.eg.sample( s,tmp,d1,d2,n)
402 s=Sample(it.FILE)
403 assert(2110 == last(s.egs)[s.all[4].at])
404 local sort1= s:betters(s.egs)
405 local lo, hi = s:clone(), s:clone()
406 for i=1,20 do lo:add(sort1[i]) end
407 for i=#sort1,#sort1-20,-1 do hi:add(sort1[i]) end
408 shout(s:mid())
409 shout(lo:mid())
410 shout(hi:mid())
411 for m,eg in pairs(sort1) do
412 n = bchop(sort1, eg,function(a,b) return s:better(a,b) end)
413 assert(m-n <=2) end end
414
415 function it.eg.dists( s,tmp,d1,d2,n)
416 s=Sample(it.FILE)
417 tmp = sort(lap(shuffle(s.egs),
418 firsts),function(eg2) return {s:dist(eg2,s.egs[1]), eg2} end),
419 firsts)
420 d1=s:dist(tmp[1][2], tmp[10][2])
421 d2=s:dist(tmp[1][2], tmp[#tmp][2])
422 assert(d1*10 < d2) end
423
424 function it.eg.binsym( s,col,tmp)
425 s=Sample(it.FILE)
426 col = s.all[4]
427 local function show(v) return out(rnds({v.n, v:mid(), v:sd()},0)) end
428 print(show(col))
429 tmp = s:splitter()
430 print(100,tmp[1])
431 for k,v in pairs(tmp[2]) do print(k,v.lo,v.hi,v.has.n, show(v.has)) end
432 end
433
434 function it.eg.hints( s,_,__,evals,sort1,train,test,n)
435 s = Sample(it.FILE)
436 evals, train,test = hints.sort(s)
437 test.egs = test:betters()
438 for m,eg in pairs(test.egs) do
439 n = bchop(train.egs, eg,function(a,b) return s:better(a,b) end) end end
440
441 ---| start-up | -----
442 it{demos=it.eg, nervous=true}
443

```

Dec 11, 21 21:01

l5.lua

Page 6/6

```

444 --[[
445
446
447
448
449 Spans
450 Little languages:
451 - options
452 - data language
453
454 Lesson plan
455 - w1: ssytems: github. github workplaces. unit tests. doco tools.
456
457 - w2: num,sym
458 - W3: sample
459 - w4: eval, knn, unfaressness
460 - W5:
461
462 - seems to be a revers that i need to do .... but dont
463 - check if shuffle is working
464
465 teaching:
466 - sample is v.useful
467 --]]

```