

```

1  #!/usr/bin/env lua
2  local b4={}; for k,v in pairs(_ENV) do b4[k]=v end;
3
4  --
5  --
6  -- a little lite
7  --
8  --
9  --
10 --
11 --
12 --
13 --
14 --
15 --
16 --
17 --
18 --
19 --
20 --
21
22 local options={
23   what = "Small sample multi-objective optimizer.",
24   usage= "(c) 2021 Tim Menzies <tim@ieee.org> unlicense.org",
25   about= [[
26     Sort N examples on multi-goals using a handful of 'hints'; i.e.
27
28     - Evaluate and rank, a few examples (on their y-values);
29     - Sort other examples by x-distance to the ranked ones;
30     - Recurse on the better half (so we sample more and more
31       from the better half, then quarter, then eighth...).
32
33     A regression tree learner then explores the examples (sorted
34     left to right, worst to best). By finding branches that
35     reduce the variance of the index of those examples, this
36     tree reports what attribute ranges select for the better (or
37     worse) examples.  ]],
38
39   how= {{"file",      "-f",      "-f./data/auto93.csv",  "read data from file"},
40         {"cull",      "-c",      ".5",                "cuts per generation"},
41         {"help",      "-h",      false,                "show help"},
42         {"hints",     "-H",      4,                    "hints per generation"},
43         {"p",         "-p",      2,                    "distance calc exponent"},
44         {"small",     "-s",      .5,                  "div list into 'small'"},
45         {"seed",      "-S",      10019,                "random number seed"},
46         {"train",     "-t",      .5,                  "size of training set"},
47         {"trivial",   "-T",      .35,                 "small delta=trivial*sd"},
48         {"todo",      "-T",      "all",                "run unit test, or 'all'"},
49         {"wild",      "-W",      false,                "run tests, no protection" }}
50
51 local the={} -- a flat list of key=value options; e.g. {seed=10019,p=2,...}
52 for _,t in pairs(options.how) do -- update defaults from command line
53   the[t[1]] = t[3]
54   for n,word in ipairs(arg) do if word==t[2] then
55     the[t[1]] = t[3] and (tonumber(arg[n+1]) or arg[n+1]) or true end end end
56
57 if the.help then -- print help text
58   print(string.format("%n%[OPTIONS]n%[options]n%[OPTIONS]n",
59     arg[0], options.usage, options.what))
60   for _,t in pairs(options.how) do
61     print(string.format("%4s%-9s%t%[options]n",
62       t[2], t[3] and t[1] or "", t[4], t[3] and"" or "", t[3] or "")) end
63   print("%n".options.about)
64   os.exit() end
65
66 --[[
67 Spans
68 Little languages:
69   - options
70   - data language
71
72 Lesson plan
73 -- w1: sysyems: github. github workplaces. unit tests. doco tools.
74 -- w2: num,sym
75 -- W3: sample
76 -- w4: eval, knn, unfairnessness
77 -- W5:
78 --]]
79

```

```

80 --
81 --
82 --
83 --
84 -- Random stuff
85 local Seed,rand,randi
86 Seed = the.seed or 10019
87 -- random integers
88 function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
89 -- random floats
90 function rand(lo,hi, mult,mod)
91   lo, hi = lo or 0, hi or 1
92   Seed = (16807 * Seed) % 2147483647
93   return lo + (hi-lo) * Seed / 2147483647 end
94
95 -----
96 -- ## Table Stuff
97 local cat,map,lap,top,keys,last,copy,pop,push
98 local sort,firsts,first,second,shuffle,bchop
99 -- Table to string.
100 cat = table.concat
101 -- Return a sorted table.
102 sort = function(t,f) table.sort(t,f); return t end
103 -- Return first,second, last item.
104 first = function(t) return t[1] end
105 second = function(t) return t[2] end
106 last = function(t) return t[#t] end
107 -- Function for sorting pairs of items.
108 firsts = function(a,b) return first(a) < first(b) end
109 -- Add to end, pull from end.
110 pop = table.remove
111 push = function(t,x) table.insert(t,x); return x end
112
113 -- Random order of items in a list (sort in place).
114 function shuffle(t, j)
115   for i=#t,2,-1 do j=randi(1,i); t[i],t[j]=t[j],t[i] end; return t end
116
117 -- Collect values, passed through 'f'.
118 function lap(t,f) return map(t,f,1) end
119 -- Collect key,values, passed through 'f'.
120 -- If 'f' returns two values, store as key,value.
121 -- If 'f' returns one values, store at index value.
122 -- If 'f' return nil then add nothing (so 'map' is also 'select').
123 function map(t,f,one, u)
124   u={}; for x,y in pairs(t) do
125     if one then x,y=f(y) else x,y=f(x,y) end
126     if x ~= nil then
127       if y then u[x]=y else u[1+#u]=x end end end
128   return u end
129
130 -- Shallow copy
131 function copy(t, u) u={}; for k,v in pairs(t) do u[k]=v end; return u end
132
133 function top(t,n, u)
134   u={};for k,v in pairs(t) do if k>n then break end; push(u,v) end; return u,end
135
136 -- Return a table's keys (sorted).
137 function keys(t,u)
138   u={};
139   for k, _ in pairs(t) do if tostring(k):sub(1,1)~="_" then push(u,k) end end
140   return sort(u) end
141
142 -- Binary chop (assumes sorted lists)
143 function bchop(t,val,lt,lo,hi, mid)
144   lt = lt or function(x,y) return x < y end
145   lo,hi = lo or 1, hi or #t
146   while lo <= hi do
147     mid=(lo+hi) // 2
148     if lt(t[mid],val) then lo=mid+1 else hi= mid-1 end end
149   return math.min(lo,#t) end
150
151 -----
152 -- ## Maths Stuff
153 local abs,sum,rand,rnds
154 abs = math.abs
155 -- Round 'x' to 'd' decimal places.
156 function rnd(x,d, n) n=10^(d or 0); return math.floor(x*n+0.5) / n end
157 -- Round list of items to 'd' decimal places.
158 function rnds(t,d) return lap(t, function(x) return rnd(x,d or 2) end) end
159
160 -- Sum items, filtered through 'f'.
161 function sum(t,f)
162   f= f or function(x) return x end
163   out=0; for _,x in pairs(f) do out = out + f(x) end; return out end
164
165 -----
166 -- ## Printing Stuff
167 local out,shout,red,green,yellow,blue,color,fmt
168 fmt = string.format
169 -- Print as red, green, yellow, blue.
170 function color(s,n) return fmt("%27[1m27[%sm%[s27[0m",n,s) end
171 function red(s) return color(s,31) end
172 function green(s) return color(s,32) end
173 function yellow(s) return color(s,34) end
174 function blue(s) return color(s,36) end
175
176 -- Printed string from a nested structure.
177 shout = function(x) print(out(x)) end
178 -- Generate string from a nested structures
179 -- (and don't print any contents more than once).
180 function out(t,seen, u,key,value,public)
181   function key(k) return fmt("%s%[blue(k),out(t[k],seen)) end
182   function value(v) return out(v,seen) end
183   if type(t) == "function" then return "(...)" end
184   if type(t) ~= "table" then return tostring(t) end
185   seen = seen or {}
186   if seen[t] then return "..." else seen[t] = t end
187   u = #t>0 and lap(t, value) or lap(keys(t), key)
188   return red(t._is or "").."[ " ..cat(u," " ..red(")") end
189
190 -----
191 -- ## File i/o Stuff
192 -- Return one table per line, split on commas.
193 local csv
194 function csv(file, line)
195   file = io.input(file)
196   line = io.read()
197   return function( t,tmp)
198     if line then
199       t={}
200       for cell in line:gsub("[\r\n]", ""):gsub("#.", ""):gmatch("[^\r\n]+") do
201         push(t, tonumber(cell) or cell) end
202       line = io.read()
203       if #t>0 then return t end
204       else io.close(file) end end end
205
206 -----
207 -- ## OO Stuff
208 local has,obj
209 -- Create an instance
210 function has(mt,x) return setmetatable(x,mt) end
211
212 -- Create a class
213 function obj(s, o,new)
214   o = {__is=s, __tostring=out}
215   o.__index = o
216   return setmetatable(o,{__call = function(_,...) return o.new(...) end}) end
217

```

```

217 --
218 --
219 --
220 --
221 -- ## Stuff for tracking 'Num'bers.
222 -- 'Num's track a list of number, and can report it sorted.
223 local Num=obj"Num"
224 function Num:new(inits, self)
225 self=has(Num,{has={}, n=0, lo=1E32, hi=-1E-32, ready=true})
226 for _,one in pairs(inits or {}) do self:add(one) end
227 return self end
228
229 function Num:add(x)
230 if x>self.hi then self.hi = x
231 elseif x<self.lo then self.lo = x end
232 push(self.has,x); self.n=self.n+1; self.ready=false end
233
234 -- Ensure that the returned list of numbers is sorted.
235 function Num:all(x)
236 if not self.ready then table.sort(self.has) end
237 self.ready = true
238 return self.has end
239
240 function Num:dist(a,b)
241 if a=="" then b=self:norm(b); a = b>.5 and 0 or 1
242 elseif b=="" then a=self:norm(a); b = a>.5 and 0 or 1
243 else a,b = self:norm(a), self:norm(b) end
244 return abs(a-b) end
245
246 -- Combine two 'num's.
247 function Num:merge(other, new)
248 new = Num:new(self.has)
249 for _,x in pairs(other.has) do new:add(x) end
250 return new end
251
252 -- Return a merged item if that combination
253 -- is simpler than its parts.
254 function Num:mergeable(other, new,b4)
255 new = self:merge(other)
256 b4 = (self.n*self:sd() + other.n*other:sd()) / new.n
257 if b4 >= new:sd() then return new end end
258
259 -- The 'mid' is the 50th percentile.
260 function Num:mid() return self:per(.5) end
261
262 -- Return 'x' normalized 0..1, lo..hi.
263 function Num:norm(x, lo,hi)
264 if x=="" then return x end
265 lo,hi = self.lo, self.hi
266 return abs(hi - lo) < 1E-32 and 0 or (x - lo)/(hi - lo) end
267
268 -- Return the 'p'-th percentile number.
269 function Num:per(p, t)
270 t = self:all()
271 p = p*#t/1
272 return #t<2 and t[1] or t[p < 1 and 1 or p>#t and #t or p] end
273
274 -- The 10th to 90th percentile range is 2.56 times the standard deviation.
275 function Num:sd() return (self:per(.9) - self:per(.1))/ 2.56 end
276
277 -- Create one span holding row indexes associated with each number
278 local div -- defined below
279 function Num:spans(col,egs)
280 local xys,xs = {}, Num()
281 for pos,eg in pairs(egs) do
282 x = eg[col]
283 if x ~= "" then
284 xs:add(x)
285 push(xys, {x=x,y=pos}) end end
286 return div(xys, -- split xys into spans...
287 #xs*the.small, -- ..where spans are of size sqrt(#xs)..
288 xs:sd()*the.trivial) end -- ..and spans have (last-first)*trivial
289
290 -----
291 -- ## Stuff for tracking 'Sym'bol Counts.
292 -- 'Sym's track symbol counts and the 'mode' (most frequent symbol).
293 local Sym=obj"Sym"
294 function Sym:new(inits, self)
295 self=has(Sym,{has={}, n=0, mode=nil, most=0})
296 for _,one in pairs(inits or {}) do self:add(one) end
297 return self end
298
299 function Sym:add(x)
300 self.n = self.n + 1
301 self.has[x] = 1 + (self.has[x] or 0)
302 if self.has[x] > self.most then self.most, self.mode = self.has[x], x end end
303
304 function Sym:dist(a,b) return a==b and 0 or 1 end
305 function Sym:mid() return self.mode end
306
307 -- Create one span holding row indexes associated with each symbol
308 function Sym:spans(col,egs,...)
309 local xys,x = {}
310 for pos,eg in pairs(egs) do
311 x = eg[col]
312 if x ~= "" then
313 xys[x] = xys[x] or {}
314 push(xys[x], pos) end end
315 return map(xys, function(x,t) return {lo=x, hi=x, has=Num(t)} end) end
316
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(ones,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 b,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 -- Return spans of the column that most reduces variance
382 function Sample:splitter()
383 function worker(col) return self:splitter(col) end
384 return first(sort(lap(sample.xs, worker), firsts))[2] end
385
386 -- Return a column's spans, and the expected sd value of those spans.
387 function Sample:splitter1(col, out,xpect)
388 out = col:spans(col,sample.eg)
389 xpect = sum(out, function(x) return x.has.n*x:sd() end)/#sample.egs
390 out = map(out, function(_,x) x.has=x.has:all(); x.col= col end)
391 return {xpect,out} end
392
393 -- Split on column with best span, recurse on each split.
394 function Sample:tree(min, node,min,sub,splitter, splitter1)
395 node = {node=self, kids={}}
396 min = min or (#self.egs)*the.small
397 if #self.egs >= 2*min then
398 for _,span in pairs(self:splitter()) do
399 sub = self:clone()
400 for _,at in pairs(span.has) do sub:add(self.egs[at]) end
401 push(node.kids, span)
402 span.has = sub:tree(min) end end
403 return node end
404
405 -- Find which leaf best matches an example 'eg'.::w
406
407 function Sample:where(tree,eg, max,x,default)
408 if #kid.has==0 then return tree end
409 max = 0
410 for _,kid in pairs(tree.node) do
411 if #kid.has > max then default,max = kid,#kid.has end
412 x = eg[kid.col]
413 if x ~= "" then
414 if x <= kid.hi and x >= kid.lo then
415 return self:where(kid.has,eg) end end end
416 return self:where(default, eg) end
417
418 -----
419 -- discretization tricks
420 -- Input a list of {(x,y,...) values. Return spans that divide the 'x' values
421 -- to minimize variance on the 'y' values.
422 function div(xys, tiny, dull, now,out,x,y)
423 function merge(b4) -- merge adjacent spans if whole is simpler than the parts
424 local j, tmp = 0, {}
425 while j < #b4 do
426 j = j + 1
427 local now, after, simpler = b4[j], b4[j+1]
428 if after then
429 simpler = now.has:mergeable(after.has)
430 if simpler then
431 now = {lo=now.lo, hi= after.hi, has=simpler}
432 j = j + 1 end end
433 push(tmp,now) end
434 return #tmp==#b4 and b4 or merge(tmp) -- recurse until nothing merged
435 end -----
436 local spans,span,out,x,y
437 xys = sort(xys, function(a,b) return a.x < b.x end)
438 span = {lo=xys[1].x, hi=xys[1].x, has=Num({})}
439 spans = {span}
440 for j,xy in pairs(xys) do
441 x,y = xy.x, xy.y
442 if j<#xys- tiny and -- if enough items remaining after split
443 x=xys[j+1].x and -- if the next item is different (so we split here)
444 span.has.n>tiny and -- if span has enough items
445 span.hi - span.lo>dull -- if span is not trivially small
446 then
447 now = push(spans, {lo=x, hi=x, has=Num({)}) -- then new span
448 span.hi = x
449 span.has:add(y) end
450 return merge(spans) end
451

```

```
452 -- hinting
453 --
454 --
```

```
455 -- Sorting on a few y values
456 local hints={}
457 function hints.default(eg) return eg end
458
459 function hints.sort(sample,scorefun, test,train,egs,scored,small)
460 sample = Sample.new(the.file)
461 train,test = {}, {}
462 for i,eg in pairs(shuffle(sample.egs)) do
463   push(i<= the.train*#sample.egs and train or test, eg) end
464 egs = copy(train)
465 small = (#egs)^the.small
466 local i=0
467 scored = {}
468 while #egs >= small do
469   local tmp = {}
470   i = i + 1
471   io.stderr:write(fmt("%s",string.char(96+i)))
472   for j=1,the.hints do
473     egs[j] = (scorefun or hints.default)(egs[j])
474     push(tmp, push(scored, egs[j]))
475   end
476   egs = hints.ranked(scored,egs,sample)
477   for i=1,the.cull*#egs//1 do pop(egs) end
478 end
479 io.stderr:write("\n")
480 train=hints.ranked(scored, train, sample)
481 return #scored, sample:clone(train), sample:clone(test) end
482
483 function hints.ranked(scored,egs,sample,worker, some)
484 function worker(eg) return (hints.rankOfClosest(scored,eg,sample),eg) end
485 scored = sample:betters(scored)
486 return lap(sort(lap(egs, worker),firsts),second) end
487
488 function hints.rankOfClosest(scored,egl,sample, worker,closest)
489 function worker(rank,eg2) return (sample:dist(egl,eg2),rank) end
490 closest = first(sort(map(scored, worker),firsts))
491 return closest[2] end --+ closest[i]/10^8 end
492
493
```

```
494 -- d(7,17) 0.5
495 --
```

```
496 local eg={}
497 function eg.shuffle( t)
498 t={}
499 for i=1,100 do push(t,i) end
500 assert(#t == #shuffle(t) and t[1] ~= shuffle(t)[1]) end
501
502 function eg.lap()
503   assert(3==lap({1,2},function(x) return x+1 end)[2]) end
504
505 function eg.map()
506   assert(3==map({1,2},function(_,x) return x+1 end)[2]) end
507
508 function eg.tables()
509   assert(20==sort(shuffle({{10,20},{30,40},{40,50}}),firsts)[1][2]) end
510
511 function eg.csv( n,z)
512   n=0
513   for eg in csv(the.file) do n=n+1; z=eg end
514   assert(n==399 and z[#z]==50) end
515
516 function eg.rnds( t)
517   assert(10.2 == first(rnds({10.22,81.22,22.33},1))) end
518
519 function eg.sym( s)
520   s=Sym{"a","a","a","a","b","b","c"}
521   assert("a"==s.mode) end
522
523 function eg.num1( n)
524   n=Num{10,20,30,40,50,10,20,30,40,50,10,20,30,40,50}
525   assert(.375 == n:norm(25))
526   assert(15.625 == n:sd()) end
527
528 function eg.num2( n1,n2,n3,n4)
529   n1=Num{10,20,30,40,50,10,20,30,40,50,10,20,30,40,50}
530   n2=Num{10,20,30,40,50,10,20,30,40,50,10,20,30,40,50}
531   assert(n1:mergeable(n2)~=nil)
532   n3=Num{10,20,30,40,50,10,20,30,40,50,10,20,30,40,50}
533   n4=Num{100,200,300,400,500,100,200,300,400,500,100,200,300,400,500}
534   assert(n3:mergeable(n4)==nil) end
535
536 function eg.sample( s,tmp,d1,d2,n)
537 s=Sample(the.file)
538 assert(2110 == last(s.egs)[s.all[3].col])
539 local sortl= s:betters(s.egs)
540 local lo, hi = s:clone(), s:clone()
541 for i=1,20 do lo:add(sortl[i]) end
542 for i=#sortl,#sortl-30,-1 do hi:add(sortl[i]) end
543 shout(s:mid())
544 shout(lo:mid())
545 shout(hi:mid())
546 for m,eg in pairs(sortl) do
547   n = bchop(sortl, eg,function(a,b) return s:better(a,b) end)
548   assert(m-n <=2) end end
549
550 function eg.dists( s,tmp,d1,d2,n)
551 s=Sample(the.file)
552 tmp = sort(lap(shuffle(s.egs),
553   firsts),function(eg2) return {s:dist(eg2,s.egs[1]), eg2} end),
554   d1=s:dist(tmp[1][2], tmp[10][2])
555   d2=s:dist(tmp[1][2], tmp[#tmp][2])
556   assert(d1*10<d2) end
557
558 function eg.binsym( s)
559 s=Sample(the.file)
560 print(s.all[6].seen._is=="Sym")
561 end
562
563 function eg.hints( s,_,_,evals,sortl,train,test,n)
564 s=Sample(the.file)
565 --for _,eg in pairs(sortl) do lap(s.ys, function(col) return eg[col.col] end ) end
566 -- assert(s.ys[4].lo==1613)
567 evals, train,test = hints.sort(s)
568 test.egs = test:betters()
569 for m,eg in pairs(test.egs) do
570   n = bchop(train.egs, eg,function(a,b) return s:better(a,b) end)
571   print(n) end end
572
573 -----
574
575 local fails, defaults = 0, copy(the)
576 local function example(k, f,ok,msg)
577   f= eg[k]; assert(f,"unknown action"..k)
578   the=copy(defaults)
579   Seed=the.seed
580   if the.wild then return f() end
581   ok,msg = pcall(f)
582   if ok then print(green("PASS"),k)
583   else print(red("FAIL"), k,msg); fail=fail+1 end end
584
585 -- run one or more examples
586 if the.todo=="all" then lap(keys(eg),example) else example(the.todo) end
587 -- print any rogue global variables
588 for k,v in pairs(_ENV) do if not b4[k] then print("?rogue: ",k,type(v)) end end
589 -- exit, return our test failure count.
590 os.exit(fail)
591
```

```
584 --[[
585      +---+ +---+
586      |   | |   |
587      +---+ +---+
588
589 -- seems to be a revers that i need to do .... but dont
600 -- check if shuffle is working
601
602 teaching:
603 - sample is v.useful
604 --]]
```