

```

1  -- **[Repo](https://github.com/timm/luam) &M-^@& [Issues](https://github.com/timm
2  /lua/issues) &M-^@& [License](LICENSE.md)**
3  --
4  -- The next generation of AI-literature software engineers need a deep
5  -- understanding of AI tools. To that end, I've been refactoring the
6  -- work of my AI graduate students (3 dozen over 20 years) into a
7  -- tool kit small enough to build in a semester, and which can be
8  -- refactored many ways. So my standard "intro to AI" exercise is six
9  -- weeks of homeworks where students rebuild the following code, from
10 -- scratch, in any language they like (except LUA).
11 --
12 -- <hr>
13 --
14 -- Standard supervised learners assume that all examples have labels.
15 -- When this is not true, then we need tools to incrementally
16 -- (a) summarize what has been seen so far; (b) find and focus
17 -- on the most interesting part of that summary, (c) collect
18 -- more data in that region, then (d) repeat.
19 --
20 -- <a href="div.png"></a>
21 -- To make that search manageable, it is useful to exploit a
22 -- manifold assumption; i.e.
23 -- higher-dimensional data can be approximated in a lower dimensional
24 -- manifold without loss of signal [Ch05,Le05].
25 -- Manifolds lead to _continuity_
26 -- effects; i.e. if there are fewer dimensions, then there are more
27 -- similarities between examples.
28 -- Continuity simplifies _clustering_
29 -- (and any subsequent reasoning). More similarities means easier
30 -- clustering. And after clustering, reasoning just means reason about
31 -- a handful of examples (maybe even just one) from each cluster.
32 --
33 -- **ASSIGNMENTS**
34 -- **Instance selection**: filter the data down to just a few samples per
35 -- cluster, the reason using just those.
36 -- **Anomaly detection**
37 -- **Explanation**
38 -- Discretize the numeric ranges (\*) at each level of the recursion,
39 -- then divide the data according what range best selects for one half, or the o
40 -- ther
41 -- at the data at this level of recursion.
42 -- **Multi-objective optimization**: This code
43 -- can apply Zitzler's multi-objective ranking predicate [Zit04] to prune the
44 -- worst
45 -- half of the data, then recurs on the rest [Ch18]. Assuming a large over-gener
46 -- ation
47 -- of the initial population (to say, 10,000, examples), this can be just as eff
48 -- ective
49 -- as genetic optimization [Ch18], but runs much faster.
50 -- **Semi-supervised learning**: these applications require only the _2.log(N)
51 -- labels at
52 -- of the pair of furthest points seen at each level of recursion.
53 --
54 -- local help = [[
55 --
56 -- l4 == a little LUA learner laboratory.
57 -- (c) 2022, Tim Menzies, BSD 2-clause license.
58 --
59 --
60 -- USAGE:
61 -- lua l4.lua [OPTIONS]
62 --
63 -- OPTIONS:
64 -- -Dump          stack dump on assert fails = false
65 -- -data N        data file = etc/data/auto93.csv
66 -- -enough F      recurse until rows^enough = 5
67 -- -furthest F    far = .9
68 -- -keep P        max kept items = 512
69 -- -p P           distance coefficient = 2
70 -- -seed P        set seed = 10019
71 -- -todo S        start up action (or 'all') = nothing
72 -- -help          show help = false
73 --
74 -- KEY: N=fileName F=float P=posint S=string
75 --
76 -- NOTES: This code uses Aha's distance measure [Aha91] (that can
77 -- handle numbers and symbols) to recursively divide data based on two
78 -- distant points (these two are found in linear time using the Fastmap
79 -- heuristic [Fa95]).
80 --
81 -- To avoid spurious outliers, this code use the 90% furthest points.
82 --
83 -- To avoid long runtimes, uses a subset of the data to learn where
84 -- to divide data (then all the data gets pushed down first halves).
85 --
86 -- To support explanation, optionally, at each level of recursion,
87 -- this code reports what ranges can best distinguish sibling clusters
88 -- C1,C2. The discretizer is inspired by the ChiMerge algorithm:
89 -- numerics are divided into, say, 16 bins. Then, while we can find
90 -- adjacent bins with the similar distributions in C1,C2, then
91 -- (a) merge then (b) look for other merges.
92 --
93 -- ]]
94 --
95 -- ## Namespace
96 --
97 -- Cache current globals, use at end to find rogue variables
98 -- local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
99 --
100 -- Defined local names.
101 -- local any,asserts,big,cli, csv, fails, firsts, fmt, goalp, ignorep, klassp
102 -- local lessp, map, main, many, max, merge, min, morep, new, nump, o, oo, per, pop, push
103 -- local r, rows, slots, sort, sum, thing, things, unpack
104 --
105 -- Classes have UPPER CASE names.
106 -- local CLUSTER, COLS, EGS, NUM, ROWS = {}, {}, {}, {}, {}
107 -- local SKIP, SOME, SPAN, SYM = {}, {}, {}, {}
108 --
109 -- ## Settings
110 -- Parse the help text for flags and defaults (e.g. -keep, 512).
111 -- Check for updates on those details from command line
112 -- (and and there,
113 -- some shortcuts are available;
114 -- e.g. -k N &rArr; 'keep=N';
115 -- and -booleanFlag &rArr; 'booleanFlag=not default').
116 -- local the={}
117 -- help:gsub("\n [^\s%+)]^\n"*(^%s+)", function(key, x)
118 -- for n, flag in ipairs(arg) do
119 -- if flag:sub(1,1)=="-" and key:find("^"..flag:sub(2)..".") then
120 -- x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
121 -- if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
122 -- the[key] = tonumber(x) or x end end )
123 --
124 -- =====
125 -- this code reads csv files where the words on line1 define column types.
126 -- function ignorep(x) return x:find"%$" end -- columns to ignore
127 -- function klassp(x) return x:find"%$" end -- symbolic goals to achieve
128 -- function lessp(x) return x:find"%-$" end -- number goals to minimize
129 -- function morep(x) return x:find"%+$" end -- numeric goals to maximize
130 -- function nump(x) return x:find"%[A-Z]" end -- numeric columns
131 -- function goalp(x) return morep(x) or lessp(x) or klassp(x) end
132 --
133 -- strings
134 -- fmt = string.format
135 --
136 -- maths
137 -- big = math.huge
138 -- max = math.max
139 -- min = math.min
140 --
141 -- r = math.random
142 --
143 -- tables
144 -- pop = table.remove
145 -- unpack = table.unpack
146 -- function any(t) return t[r(#t)] end
147 -- function firsts(a,b) return a[1] < b[1] end
148 -- function many(t,n, u) u={}; for i=1,n do push(u,any(t)) end; return u end
149 -- function per(t,p) return t[ (#t*(p or .5))/1 ] end
150 -- function push(t,x) table.insert(t,x); return x end
151 -- function sort(t,f) table.sort(t,f); return t end
152 --
153 -- meta
154 -- function map(t,f, u) u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
155 -- function sum(t,f, n) n=0; for _,v in pairs(t) do n=n+f(v) end; return n end
156 -- function slots(t, u)
157 -- u={}
158 -- for k,v in pairs(t) do k=tostring(k);if k:sub(1,1)~="_" then push(u,k) end end
159 -- return sort(u) end
160 --
161 -- print tables, recursively
162 -- function oo(t) print(o(t)) end
163 -- function o(t)
164 -- if type(t)~="table" then return tostring(t) end
165 -- local key=function(k) return fmt("%s %s",k,o(t[k])) end
166 -- local u = #t>0 and map(t,o) or map(slots(t),key)
167 -- return '{ '..table.concat(u, " " )..' }' end
168 --
169 -- strings to things
170 -- function csv(file, x)
171 -- file = io.input(file)
172 -- return function()
173 -- x=io.read(); if x then return things(x) else io.close(file) end end end
174 --
175 -- function thing(x)
176 -- x = x:match("^%s*(.)%s*$")
177 -- if x=="true" then return true elseif x=="false" then return false end
178 -- return tonumber(x) or x end
179 --
180 -- function things(x,sep, t)
181 -- t={}
182 -- for y in x:gmatch(sep or "[^,]+") do push(t,thing(y)) end
183 -- return t end

```

```

174 -- CLASSES
175
176
177
178 function new(k,t) k.__index=k; k.__tostring=o; return setmetatable(t,k) end
179
180 -- COLS: turns list of column names into NUMs, SYMs, or SKIPs
181 function COLS.new(k,row, i)
182   i = new(k, {all={}, x={}, y={}, names=row})
183   for at,txt in ipairs(row) do push(i.all, i.col(at,txt)) end
184   return i end
185
186 function COLS.add(i,t)
187   for _,col in pairs(i.all) do col:add( t[col.at] ) end
188   return t end
189
190 function COLS.col(i,at,txt, col)
191   if ignorep(txt) then return SKIP:new(at,txt) end
192   col = (nump(txt) and NUM or SYM):new(at,txt)
193   push(goalp(txt) and i.y or i.x, col)
194   if klassp(txt) then i.klass = col end
195   return col end
196
197 -- NUM: summarizes a stream of numbers
198 function NUM.new(k,n,s)
199   return new(k, {n=0, at=n or 0, txt=s or "", has=SOME:new(), ok=false,
200     w=lessp(s or "") and -1 or 1, lo=big, hi=-big}) end
201
202 function NUM.add(i,x)
203   if x == "?" then
204     i.n = i.n + 1
205     if i.has:add(x) then i.ok=false end
206     i.lo,i.hi = min(x,i.lo), max(x,i.hi); end end
207
208 function NUM.dist(i,x,y)
209   if x=="?" and y=="?" then return 1
210   elseif x=="?" then y=i:norm(y); x=y<0.5 and 1 or 0
211   elseif y=="?" then x=i:norm(x); y=x<0.5 and 1 or 0
212   else x,y = i:norm(x), i:norm(y) end
213   return math.abs(x-y) end
214
215 function NUM.mid(i) return per(i:sorted(), .5) end
216
217 function NUM.norm(i,x)
218   return math.abs(i.hi-i.lo)<1E-9 and 0 or (x-i.lo)/(i.hi - i.lo) end
219
220 function NUM.sorted(i)
221   if i.ok==false then table.sort(i.has.all); i.ok=true end
222   return i.has.all end
223
224 -- ROWS: manages 'rows', summarized in 'cols' (columns).
225 function ROWS.new(k,init,s, i)
226   i = new(k, {rows={}, cols=nil})
227   if type(init)=="string" then for t in csv(init) do i:add(t) end end
228   if type(init)=="table" then for t in init do i:add(t) end end
229   return i end
230
231 function ROWS.add(i,t)
232   if i.cols then push(i.rows,i.cols:add(t)) else i.cols=COLS:new(t) end end
233
234 function ROWS.clone(i, j) j = ROWS:new(); j:add(i.cols.names); return j end
235
236 function ROWS.dist(i,row1,row2, d,fun)
237   function fun(col) return col:dist(row1[col.at], row2[col.at])^the.p end
238   return (sum(i.cols.x, fun)/ #i.cols.x)^(1/the.p) end
239
240 function ROWS.furthest(i,row1,rows, fun)
241   function fun(row2) return {i:dist(row1,row2), row2} end
242   return unpack(per(sort(map(rows,fun),firsts), the.furthest)) end
243
244 function ROWS.half(i, top)
245   local some, top,c,x,y,tmp,mid,lefs,rights,_
246   some = many(i.rows, the.keep)
247   top = top or 1
248   _,x = top:furthest(any(some), some)
249   c,y = top:furthest(x, some)
250   tmp = sort(map(i.rows,function(r) return top:fastmap(r,x,y,c) end),firsts)
251   mid = #i.rows//2
252   lefs, rights = i:clone(), i:clone()
253   for at,row in pairs(tmp) do (at < mid and lefs or rights):add(row[2]) end
254   return lefs,rights,x,y,c, tmp[mid] end
255
256 function ROWS.mid(i,cols)
257   return map(cols or i.cols.all, function(col) return col:mid() end) end
258
259 function ROWS.fastmap(i, r,x,y,c, a,b)
260   a,b = idist(r,x), idist(r,y); return {(a^2 + c^2 - b^2)/(2*c), r} end
261
262 -- SKIP: summarizes things we want to ignore (so does nothing)
263 function SKIP.new(k,n,s) return new(k, {n=0, at=at or 0, txt=s or ""}) end
264 function SKIP.add(i,x) return x end
265 function SKIP.mid(i) return "?" end
266
267 -- SOME: keeps a random sample on the arriving data
268 function SOME.new(k,keep) return new(k, {n=0, all={}, keep=keep or the.keep}) end
269 function SOME.add(i,x)
270   i.n = i.n+1
271   if #i.all < i.keep then push(i.all,x) ; return i.all
272   elseif r() < i.keep/i.n then i.all[r(#i.all)]=x; return i.all end end
273
274 -- SYM: summarizes a stream of symbols
275 function SYM.new(k,n,s)
276   return new(k, {n=0, at=n or 0, txt=s or "", has={}, most=0}) end
277
278 function SYM.add(i,x,inc)
279   if x == "?" then
280     inc = inc or 1
281     i.n = i.n + inc
282     i.has[x] = inc + (i.has[x] or 0)
283     if i.has[x] > i.most then i.most,i.mode=i.has[x],x end end end
284
285 function SYM.dist(i,x,y) return (x=="?" and y=="?" and 1) or (x==y and 0 or 1) end
286 function SYM.mid(i) return i.mode end
287 function SYM.div(i, p)
288   return sum(i.has,function(k) p=-i.has[k]/i.n; return -p*math.log(p,2) end) end
289
290 function SYM.merge(i,j, k)
291   k = SYM:new(i.at,i.txt)
292   for x,n in pairs(i.has) do k:add(x,n) end
293   for x,n in pairs(j.has) do k:add(x,n) end
294   ei, ej, ek = i:div(), j:div(), k:div()
295   if i.n==0 or j.n==0 or .99*ek <= (i.n*ei + j.n*ej)/k.n then
296     return k end end

```

```

297 -- CLUSTER
298
299
300
301 -- CLUSTER: recursively divides data by clustering towards two distant points
302 function CLUSTER.new(k,sample,top)
303   local i,enough,left,right
304   top = top or sample
305   i = new(k, {here=sample})
306   enough = (#top.rows)^the.enough
307   if #sample.rows >= 2*enough then
308     left, right, i.x, i.y, i.c, i.mid = sample:half(top)
309     if #left.rows < #sample.rows then
310       i.left = CLUSTER:new(left, top)
311       i.right = CLUSTER:new(right, top) end end
312   return i end
313
314 function CLUSTER.show(i,pre, here)
315   pre = pre or ""
316   here=""
317   if not i.left and not i.right then here= o(i.here:mid(i.here.cols.y)) end
318   print(fmt("%s: %-30s %s", #i.here.rows, pre, here))
319   for _,kid in pairs(i.left, i.right) do
320     if kid then kid:show(pre .. "|." ) end end end
321
322 -- EXPLAIN
323
324
325 -- SPAN: keeps a random sample on the arriving data
326 function SPAN.new(k, col, lo, hi, has)
327   return new(k, {col=col, lo=lo, hi=hi or lo, has=has or SYM:new()}) end
328
329 function SPAN.add(i,x,y,n) i.lo,i.hi=min(x,i.lo),max(x,i.hi); i.has:add(y,n) end
330 function SPAN.merge(i,j)
331   local has = i.has:merge(j.has)
332   if now then return SPAN:new(i.col, i.lo, j.hi, has) end end
333
334 function SPAN.select(i,row, x)
335   x = row[i.col.at]
336   return (x=="") or (i.lo==i.hi and x==i.lo) or (i.lo <= x and x < i.hi) end
337
338 -- EXPLAIN
339 function EXPLAIN(k,sample,top)
340   i.here = sample
341   top = top or sample
342   enough = (#top.rows)^the.enough
343   if #top.rows >= 2*enough then
344     left, right = sample:half(top)
345     spans = {}
346     for n,col in pairs(i.cols.x) do
347       tmp = col:spans(j.cols.x[n])
348       if #tmp>1 then for _,one in pairs(tmp) do push(spans,one) end end
349       if #spans > 2 then
350         XXXX?
351       end
352     end
353   end
354   function SYM.spans(i, j)
355     local xys,all,one,last,xys,x,c,n = {},{}
356     for x,n in pairs(i.has) do push(xys, {x,"this",n}) end
357     for x,n in pairs(j.has) do push(xys, {x,"that",n}) end
358     for _,tmp in ipairs(sort(xys,firsts)) do
359       x,c,n = unpack(tmp)
360       if x ~= last then
361         last = x
362         one = push(all, Span(i,x,x)) end
363       one:add(x,y,n) end
364     return all end
365   end
366   function NUM.spans(i, j)
367     local xys,all,lo,hi,gap,xys,one,x,c,n = {},{}
368     lo,hi = min(i.lo, j.lo), max(i.hi,j.hi)
369     gap = (hi - lo) / bins
370     for x,n in pairs(i.has) do push(xys, {x,"this",1}) end
371     for x,n in pairs(j.has) do push(xys, {x,"that",1}) end
372     one = Span:new(i,lo,lo)
373     all = {one}
374     for _,tmp in ipairs(sort(xys,first)) do
375       x,c,n = unpack(tmp)
376       if one.hi - one.lo > gap then one = push(all, Span(i, one.hi, x)) end
377       one:add(x,y,n) end
378     return all end
379
380 function merge(b4, j,n,now,a,b,merged)
381   j,n,now = 0,#b4,{}
382   while j < #b4 do
383     j = j+1
384     a, b = b4[j], b4[j+1]
385     if b then
386       merged = a:merge(b)
387       if merged then a,j = merged, j+1 end end
388     push(now,a)
389     j = j+1 end
390   return #now == #b4 and b4 or merge(now) end

```

```

393 -- DEMOS
394 --
395 --
396 --
397 fails=0
398 function asserts(test, msg)
399   print(test and "PASS: " or "FAIL: ", msg or "")
400   if not test then
401     fails=fails+1
402     if the.dump then assert(test, msg) end end end
403
404 function EGS.nothing() return true end
405 function EGS.the() oo(the) end
406 function EGS.rand() print(r()) end
407 function EGS.some(s, t)
408   s=SOME:new(100)
409   for i=1,100000 do s:add(i) end
410   for j,x in pairs(sort(s.all)) do
411     --if (j % 10)==0 then print("") end
412     --io.write(fmt("%6s", x)) end end
413     fmt("%6s", x) end end
414
415 function EGS.clone( r, s)
416   r = ROWS:new(the.data)
417   s = r:clone()
418   for _,row in pairs(r.rows) do s:add(row) end
419   asserts(r.cols.x[1].lo==s.cols.x[1].lo, "clone.lo")
420   asserts(r.cols.x[1].hi==s.cols.x[1].hi, "clone.hi")
421   end
422
423 function EGS.data( r)
424   r = ROWS:new(the.data)
425   asserts(r.cols.x[1].hi == 8, "data.columns") end
426
427 function EGS.dist( r, rows, n)
428   r = ROWS:new(the.data)
429   rows = r.rows
430   n = NUM:new()
431   for _,row in pairs(rows) do n:add(r:dist(row, rows[1])) end
432   --oo(r.cols.x[2]:sorted()) end
433   o(r.cols.x[2]:sorted()) end
434
435 function EGS.many( t)
436   t={} for j=1,100 do push(t, j) end
437   --print(oo(many(t, 10))) end
438   o(many(t, 10)) end
439
440 function EGS.far( r, c, row1, row2)
441   r = ROWS:new(the.data)
442   row1 = r.rows[1]
443   c, row2 = r:far(r.rows[1], r.rows) end
444   --print(c, "\n", o(row1), "\n", o(row2)) end
445
446 function EGS.half( r, c, row1, row2)
447   local lefts, rights, x, y, x
448   r = ROWS:new(the.data)
449   r:mid(r.cols.y)
450   lefts, rights, x, y, c = r:half()
451   lefts:mid(lefts.cols.y)
452   rights:mid(rights.cols.y)
453   asserts(true, "half") end
454
455 function EGS.cluster(r)
456   r = ROWS:new(the.data)
457   --CLUSTER:new(r):show() end
458   CLUSTER:new(r) end
459
460 -- start-up
461 if arg[0] == "slua" then
462   oo(the)
463   if the.help then print(help:gsub("\nNOTES:*$", "")) else
464     local b4={} for k,v in pairs(the) do b4[k]=v end
465     for _,todo in pairs(the.todo=="all" and slots(EGS) or {the.todo}) do
466       for k,v in pairs(b4) do the[k]=v end
467       math.randomseed(the.seed)
468       if type(EGS[todo])=="function" then EGS[todo]() end end
469     end
470     for k,v in pairs(_ENV) do if not b4[k] then print("?", k, type(v)) end end
471     os.exit(fails)
472   else
473     return {CLUSTER=CLUSTER, COLS=COLS, NUM=NUM, ROWS=ROWS,
474            SKIP=SKIP, SOME=SOME, SYM=SYM, the=the, oo=oo, o=o}
475   end
476
477 -- git rid of SOME for rows
478 -- nss = NUM | SYM | SKIP
479 -- COLS = all:[nss]t, x:[nss]*, y:[nss]*, klass:col?
480 -- ROWS = cols:COLS, rows:SOME
481 -- ## References
482 -- - [Ah91]:
483 --   Aha, D.W., Kibler, D. & Albert, M.K. Instance-based
484 --   learning algorithms. Mach Learn 6, 37&M-^@M-^S66 (1991).
485 --   https://doi.org/10.1007/BF00153759
486 -- - [Boley, 1998]:
487 --   Boley, D., 1998.
488 --   [Principal directions divisive partitioning](https://www-users.cse.umn.edu/~boley/publications/papers/PDDP.pdf)
489 --   Data Mining and Knowledge Discovery, 2(4): 325-344.
490 -- - [Ch05]:
491 --   [Semi-Supervised Learning](http://www.molgen.mpg.de/3659531/MITPress--SemiSupervised-Learning)
492 --   (2005) Olivier Chapelle, Bernhard Sch&#252;l&#228;kopf, and Alexander Zien (eds).
493 --   MIT Press.
494 -- - [Ch18]:
495 --   [Sampling&M-^@M-^] as a Baseline Optimizer for Search-Based Software Engineer
496 --   ing](https://arxiv.org/pdf/1608.07617.pdf),
497 --   Jianfeng Chen; Vivek Nair; Rahul Krishna; Tim Menzies
498 --   IEEE Trans SE, (45)6, 2019
499 -- - [Ch22]:
500 --   [Can We Achieve Fairness Using Semi-Supervised Learning?](https://arxiv.org/p
501 --   df/2111.02038.pdf)
502 --   (2022), Joymallya Chakraborty, Huy Tu, Suvodeep Majumder, Tim Menzies.
503 -- - [Fal95]:
504 --   Christos Faloutsos and King-Ip Lin. 1995. FastMap: a fast algorithm for index
505 --   ing, data-mining and visualization of traditional and multimedia datasets. SIGMO
506 --   D Rec. 24, 2 (May 1995), 163&M-^@M-^S174. DOI:https://doi.org/10.1145/568271.223
507 --   812
508 -- - [Le05]:
509 --   Levina, E., Bickel, P.J.: [Maximum likelihood estimation of intrinsic dimensi
510 --   on](https://www.stat.berkeley.edu/~bickel/mldim.pdf).
511 --   In:
512 --   Advances in neural information processing systems, pp. 777&M-^@M-^S784 (2005)
513 -- - [Pl04]:
514 --   Platt, John.
515 --   [FastMap, MetricMap, and Landmark MDS are all Nystrom Algorithms](https://www
516 --   .microsoft.com/en-us/research/wp-content/uploads/2005/01/nystrom2.pdf)
517 --   AISTATS (2005).
518 -- - [Zit04]:
519 --   [Indicator-based selection in multiobjective search](https://link.springer.co
520 --   m/chapter/10.1007/978-3-540-30217-9_84)
521 --   Eckart Zitzler, Simon K&#228;nzli
522 --   Proc. 8th International Conference on Parallel Problem Solving from Nature (P
523 --   PSN VIII

```