

```

1  -- <img style="margin-left: 5px;" src=header.png width=150 align=right>
2
3  -- The next generation of AI-literature software engineers need
4  -- to know both the services within AI tools. To that end, I've
5  -- been refactoring the work of my AI graduate students (3 dozen over 20 years)
6  -- into a tool kit small enough to build in a semester, that can be refactored
7  -- many ways. So my standard "intro to AI" project is to get students to
8  -- rebuild the following code, from scratch, in any language they like
9  -- (except LUA) in six weekly homeworks. <br clear=all>
10
11 -- Standard supervised learners assume that all examples have labels.
12 -- When this is not true, then we need tools to incrementally
13 -- (a) summarize what has been seen so far; (b) find and focus
14 -- on the most interesting part of that summary, (c) collect
15 -- more data in that region, then (d) repeat.
16
17 -- <a href="div.png"></a>
18 -- To make that search manageable, it is useful to exploit a
19 -- manifold assumption; i.e.
20 -- higher-dimensional data can be approximated in a lower dimensional
21 -- manifold without loss of signal [Ch05,Le05].
22 -- Manifolds lead to _continuity
23 -- effects; i.e. if there are fewer dimensions, then there are more
24 -- similarities between examples.
25 -- Continuity simplifies _clustering_
26 -- (and any subsequent reasoning). More similarities means easier
27 -- clustering. And after clustering, reasoning just means reason about
28 -- a handful of examples (maybe even just one) from each cluster.
29
30 -- To exploit the manifold assumption, this code uses Aha's distance
31 -- measure [Aha91] (that can handle numbers and symbols) to
32 -- recursively divide data based on two distant points (these are found
33 -- in linear time using the Fastmap heuristic [Fa95]). To avoid spurious
34 -- outliers, this code use the 90% furthest points. To avoid long runtimes,
35 -- uses a subset of the data to learn where to divide data (then all the data
36 -- gets pushed down first halves).
37
38 -- **ASSIGNMENTS**
39 -- **Instance selection**: filter the data down to just a few samples per
40 -- cluster, the reason using just those.
41 -- **Anomaly detection**
42 -- **Explanation**
43 -- Discretize the numeric ranges (\*) at each level of the recursion,
44 -- then divide the data according what range best selects for one half, or the o
45 -- ther
46 -- at the data at this level of recursion.
47 -- **Multi-objective optimization**: This code
48 -- can apply Zitzler's multi-objective ranking predicate [Zit04] to prune the
49 -- worst
50 -- half of the data, then recurs on the rest [Ch18]. Assuming a large over-gener
51 -- ation
52 -- of the initial population (to say, 10,000, examples), this can be just as eff
53 -- ective
54 -- as genetic optimization [Ch18], but runs much faster.
55 -- **Semi-supervised learning**: these applications require only the _2.log(N)
56 -- _ labels at
57 -- of the pair of furthest points seen at each level of recursion.
58
59 --
60 local help = [[
61
62 l4 == a little LUA learner laboratory.
63 (c) 2022, Tim Menzies, BSD 2-clause license.
64
65 USAGE:
66 lua l4.lua [OPTIONS]
67
68 OPTIONS:
69 -Dump stack dump on assert fails = false
70 -data N data file = etc/data/auto93.csv
71 -enough F recurse until rows^enough = .5
72 -furthest F far = .9
73 -keep P max kept items = 512
74 -p P distance coefficient = 2
75 -seed P set seed = 10019
76 -todo S start up action (or 'all') = nothing
77 -help show help = false
78
79 KEY: N=fileName F=float P=posint S=string
80
81 NOTES: This code uses Aha's distance measure [Aha91] (that can
82 handle numbers and symbols) to recursively divide data based on two
83 distant points (these two are found in linear time using the Fastmap
84 heuristic [Fa95]).
85
86 To avoid spurious outliers, this code use the 90% furthest points.
87
88 To avoid long runtimes, uses a subset of the data to learn where
89 to divide data (then all the data gets pushed down first halves).
90
91 To support explanation, optionally, at each level of recursion,
92 this code reports what ranges can best distinguish sibling clusters
93 Cl,C2.. The discretizer is inspired by the ChiMerge algorithm:
94 numericize arg0, say, 16 bins. Then, while we can find
95 adjacent bins, merge them then look for other merges.
96 ]]
97
98 -- ## Namespace
99
100 -- Cache current globals, use at end to find rogue variables
101 local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
102
103 -- Defined local names.
104 local any,asserts,big,cli,csv,fails,firsts,fmt,goalp,ignorep,klassp
105 local lessp,map,main,many,max,merge,min,morep,new,nump,o,oo,per,pop,push
106 local r,rows,slots,sort,sum,thing,things,unpack
107
108 -- Classes have UPPER CASE names.
109 local CLUSTER, COLS, EGS, NUM, ROWS = {}, {}, {}, {}, {}
110 local SKIP, SOME, SPAN, SYM = {}, {}, {}, {}
111
112 -- ## Settings
113 -- Parse the help text for flags and defaults (e.g. -keep, 512).
114 -- Check for updates on those details from command line
115 -- (and and there,
116 -- some shortcuts are available;
117 -- e.g. _-k N _&rArr; 'keep=N';
118 -- and _-booleanFlag_&rArr; 'booleanFlag=not default').
119 local the={}
120 help:gsub("(%n[-](%s+)%n)*%s(%s+)",function(key,x)
121   for n,flag in ipairs(arg) do
122     if flag:sub(1,1)=="-" and key:find("^"..flag:sub(2)..".*") then
123       x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
124   if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
125     the[key] = tonumber(x) or x end end )
126
127 -----
128 -- this code reads csv files where the words on line1 define column types.
129 function ignorep(x) return x:find"$" end -- columns to ignore
130 function klassp(x) return x:find"$" end -- symbolic goals to achieve
131 function lessp(x) return x:find"$" end -- number goals to minimize
132 function morep(x) return x:find"+$" end -- numeric goals to maximize
133 function nump(x) return x:find"[A-Z]" end -- numeric columns
134 function goalp(x) return morep(x) or lessp(x) or klassp(x) end
135
136 -- strings
137 fmt = string.format

```

```

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

```

```

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

```

```

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

```