```lua
1    -- <img align=left width=150 src=head.png>
2    --
3    -- **[Repo](https://github.com/timm/lua) âM-^@¢ [Issues](https://github.com/timm
     /lua/issues) âM-^@¢ [&copy;2022](LICENSE.md)** Tim Menzies
4    --
5    -- The next generation of AI-literature software engineers need a deep
6    -- understanding of AI tools.  To that end, I've been looking back over
7    -- common themes from my
8    -- AI graduate students (30+ students, over 20 years). The result
9    -- is a
10   -- tool kit small enough to build in a semester, and which can be
11   -- refactored many ways.  The code may not be optimal, but it does
12   -- enable a tour through many of the current themes in AI. No, its
13   -- not a deep learner since I want to build succinct high-level models
14   -- that humans can read, understand, critique, and change.
15   --
16   -- Anyway, my standard "intro to AI" exercise is six
17   -- weeks of homeworks where students rebuild the following code,from
18   -- scratch, in any language they like (except LUA).  After that,
19   -- students can review all the assumptions of this code, then read the
20   -- literature looking for other tools that challenge those assumptions.
21   -- That leads to a second a 4-6 week project using these tools as a baseline aga
     inst
22   -- which they can compare other approaches.
23   --
24   -- <hr>
25   --
26   -- Standard supervised learners assume that all examples have labels.
27   -- When this is not true, then we need tools to incrementally
28   -- (a) summarize what has been seen so far; (b) find and focus
29   -- on the most interesting part of that summary, (c) collect
30   -- more data in that region, then (d) repeat.
31   --
32   -- <a href="div.png"><img align=right width=225 src="div.png"></a>
33   -- To make that search manageable, it is useful to exploit a
34   -- manifold assumption; i.e.
35   -- higher-dimensional data can be approximated in a lower dimensional
36   -- manifold without loss of signal [Ch05,Le05].
37   -- Manifolds lead to _continuity_
38   -- effects; i.e. if there are fewer dimensions, then there are more
39   -- similarities between examples.
40   -- Continuity simplifies _clustering_
41   -- (and any subsequent reasoning).  More similarities means  easier
42   -- clustering. And after clustering, reasoning just means reason about
43   -- a handful of examples (maybe even just one)  from each cluster.
44   --
45   -- **ASSIGNMENTS**
46   -- - **Instance selection**: filter the data down to just a few samples per
47   -- cluster, the reason using just those.
48   -- - **Anomaly detection**
49   -- - **Explanation**
50   -- Discretize the numeric ranges (\*) at each level of the recursion,
51   -- then divide the data according what range best selects for one half, or the o
     ther
52   -- at the data at this level of recursion.
53   -- - **Multi-objective optimization:** This code
54   -- can apply Zitzler's multi-objective rankining predicate [Zit04] to prune the
     worst
55   -- half of the data, then recurs on the rest [Ch18]. Assuming a large over-gener
     ation
56   -- of the initial population (to say, 10,000, examples), this can be just as eff
     ective
57   -- as genetic optimization [Ch18], but runs much faster.
58   -- - **Semi-supervised learning**: these applications require only the _2.log(N)
     _ labels at
59   -- of the pair of furthest points seen at each level of recursion.
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     -cohen    F   Cohen's delta               = .35
70     -data     N   data file                   = etc/data/auto93.csv
71     -Dump         stack dump on assert fails = false
72     -furthest F   far                         = .9
73     -Format   s   format string               = %5.2f
74     -keep     P   max kept items              = 512
75     -p        P   distance coefficient        = 2
76     -seed     P   set seed                    = 10019
77     -todo     S   start up action (or 'all') = nothing
78     -help         show help                   = false
79     -want     F   recurse until rows^want     = .5
80
81   KEY: N=fileName F=float P=posint S=string
82
83   NOTES: This code uses Aha's distance measure [Aha91] (that can
84   handle numbers and symbols) to recursively divide data based on two
85   distant points (these two are found in linear time using the Fastmap
86   heuristic [Fa95]).
87
88   To avoid spurious outliers, this code use the 90% furthest points.
89
90   To avoid long runtimes, uses a subset of the data to learn where
91   to divide data (then all the data gets pushed down first halves).
92
93   To support explanation, optionally, at each level of recursion,
94   this code reports what ranges can best distinguish sibling clusters
95   C1,C2.  The  discretizer is inspired by the ChiMerge algorithm:
96   numerics are divided into, say, 16 bins. Then, while we can find
97   adjacent bins with the similar distributions in C1,C2, then
98   (a) merge then (b) look for other merges.
99   ]]
100
101  -- ## Namespace
102
103  -- Cache current globals, use at end to find rogue variables
104  local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
105
106  -- Defined local names.
107  local any,asserts,big,cli,csv,fails,firsts,fmt,goalp,ignorep,klassp
108  local lessp,map,main,many,max,merge,min,morep,new,nump,o,oo,per,pop,push
109  local r,rows,rnd,rnds,slots,sort,sum,thing,things,unpack
110
111  -- Classes have UPPER CASE names.
112  local CLUSTER, COLS, EGS,  NUM, ROWS = {},{},{},{},{}
113  local SKIP,    SOME, SPAN, SYM       = {},{},{},{}
114
115  -- ## Settings
116  -- Parse the help text for flags and defaults (e.g. -keep, 512).
117  -- Check for updates on those details from command line
118  -- (and and there,
119  -- some shortcuts are available;
120  -- e.g.  _-k N_ &rArr; `keep=N`;
121  -- and  _-booleanFlag_ &rArr; `booleanFlag=not default`).
122  local the={}
123  help:gsub("\n [-]([^%s]+)[^\n]*%s([^%s]+)",function(key,x)
124    for n,flag in ipairs(arg) do
125      if flag:sub(1,1)=="-" and key:find("^"..flag:sub(2).."*") then
126        x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
127    if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
128      the[key] = tonumber(x) or x end end )
129

130  -- ---------------------------------------------------------------------
131  -- this code reads csv files where the words on line1 define column types.
132  function ignorep(x) return x:find":$" end     -- columns to ignore
133  function klassp(x)  return x:find"!$" end      -- symbolic goals to achieve
134  function lessp(x)   return x:find"-$" end      -- number goals to minimize
135  function morep(x)   return x:find"+$" end      -- numeric goals to maximize
136  function nump(x)    return x:find"^[A-Z]" end  -- numeric columns
137  function goalp(x)   return morep(x) or lessp(x) or klassp(x) end
138
139  -- strings
140  fmt = string.format
141
142  -- maths
143  big = math.huge
144  max = math.max
145  min = math.min
146  r   = math.random
147
148  function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
149  function rnd(x,f)
150    return fmt(type(x)=="number" and x~=x//1 and f or the.rnd or "%s",x) end
151
152  -- tables
153  pop = table.remove
154  unpack = table.unpack
155  function any(t)        return t[r(#t)] end
156  function firsts(a,b)   return a[1] < b[1] end
157  function many(t,n, u)  u={}; for i=1,n do push(u,any(t)) end; return u end
158  function per(t,p)      return t[ (#t*(p or .5))//1 ] end
159  function push(t,x)     table.insert(t,x); return x end
160  function sort(t,f)     table.sort(t,f); return t end
161
162  -- meta
163  function map(t,f, u)   u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
164  function sum(t,f, n)   n=0; for _,v in pairs(t) do n=n+f(v)      end; return n end
165  function slots(t, u)
166    u={}
167    for k,v in pairs(t) do k=tostring(k);if k:sub(1,1)~="_" then push(u,k) end end
168    return sort(u) end
169
170  -- print tables, recursively
171  function oo(t)  print(o(t)) end
172  function o(t)
173    if type(t)~="table" then return tostring(t) end
174    local key=function(k) return fmt(":%s %s",k,o(t[k])) end
175    local u = #t>0 and map(t,o) or map(slots(t),key)
176    return '{'..table.concat(u,"")..'}' end
177
178  -- strings to things
179  function csv(file,    x)
180    file = io.input(file)
181    return function()
182      x=io.read(); if x then return things(x) else io.close(file) end end end
183
184  function thing(x)
185    x = x:match"^%s*(-)%s*$"
186    if x=="true" then return true elseif x=="false" then return false end
187    return tonumber(x) or x end
188
189  function things(x,sep,  t)
190    t={}
191    for y in x:gmatch(sep or"([^,]+)") do push(t,thing(y)) end
192    return t end
```

```
193  --
194  -- ┌─┐ ┌   ┌─┐ ┌─┐ ┌─┐ ┌─┐
195  -- │   │   │─┤ └─┐ └─┐ ├─┘
196  --
```

```
197  function new(k,t) k.__index=k; k.__tostring=o; return setmetatable(t,k) end
198
199  -- COLS: turns list of column names into NUMs, SYMs, or SKIPs
200  function COLS.new(k,row,    i)
201    i= new(k,{all={},x={},y={},names=row})
202    for at,txt in ipairs(row) do  push(i.all, i:col(at,txt)) end
203    return i end
204
205  function COLS.add(i,t)
206    for _,col in pairs(i.all) do col:add( t[col.at] ) end
207    return t end
208
209  function COLS.col(i,at,txt,     col)
210    if ignorep(txt) then return SKIP:new(at,txt) end
211    col = (nump(txt) and NUM or SYM):new(at,txt)
212    push(goalp(txt) and i.y or i.x, col)
213    if klassp(txt) then i.klass = col end
214    return col end
215
216  -- NUM: summarizes a stream of numbers
217  function NUM.new(k,n,s)
218    return new(k,{n=0,at=n or 0,txt=s or"",has=SOME:new(),ok=false,
219                 w=lessp(s or "") and -1 or 1, lo=big, hi=-big}) end
220
221  function NUM.add(i,x)
222    if x ~= "?" then
223      i.n = i.n + 1
224      if i.has:add(x) then i.ok=false end
225      i.lo,i.hi = min(x,i.lo), max(x,i.hi); end end
226
227  function NUM.dist(i,x,y)
228    if    x=="?" and y=="?" then return 1
229    elseif x=="?" then y=i:norm(y); x=y<0.5 and 1 or 0
230    elseif y=="?" then x=i:norm(x); y=x<0.5 and 1 or 0
231    else   x,y = i:norm(x), i:norm(y) end
232    return math.abs(x-y) end
233
234  function NUM.mid(i) return per(i:sorted(), .5) end
235
236  function NUM.norm(i,x)
237    return math.abs(i.hi-i.lo)<1E-9 and 0 or (x-i.lo)/(i.hi - i.lo) end
238
239  function NUM.sorted(i)
240    if i.ok==false then table.sort(i.has.all); i.ok=true end
241    return i.has.all end
242
243  -- ROWS: manages `rows`, summarized in `cols` (columns).
244  function ROWS.new(k,inits,     i)
245    i = new(k,{rows={},cols=nil})
246    if type(inits)=="string" then for t in csv(inits) do i:add(t) end end
247    if type(inits)=="table"  then for t in inits        do i:add(t) end end
248    return i end
249
250  function ROWS.add(i,t)
251    if i.cols then push(i.rows,i.cols:add(t)) else i.cols=COLS:new(t) end end
252
253  function ROWS.clone(i,  j) j= ROWS:new(); j:add(i.cols.names);return j end
254
255  function ROWS.dist(i,row1,row2,   d,fun)
256    function fun(col) return col:dist(row1[col.at], row2[col.at])^the.p end
257    return (sum(i.cols.x, fun)/ #i.cols.x)^(1/the.p) end
258
259  function ROWS.furthest(i,row1,rows,     fun)
260    function fun(row2) return {i:dist(row1,row2), row2} end
261    return unpack(per(sort(map(rows,fun),firsts), the.furthest)) end
262
263  function ROWS.half(i, top)
264    local some, top,c,x,y,tmp,mid,lefts,rights,_
265    some= many(i.rows, the.keep)
266    top = top or i
267    _,x = top:furthest(any(some), some)
268    c,y = top:furthest(x,           some)
269    tmp = sort(map(i.rows,function(r) return top:fastmap(r,x,y,c) end),firsts)
270    mid = #i.rows//2
271    lefts, rights = i:clone(), i:clone()
272    for at,row in pairs(tmp) do (at <=mid and lefts or rights):add(row[2]) end
273    return lefts,rights,x,y,c, tmp[mid] end
274
275  function ROWS.mid(i,cols)
276    return map(cols or i.cols.all, function(col) return col:mid() end) end
277
278  function ROWS.fastmap(i, r,x,y,c,      a,b)
279    a,b = i:dist(r,x), i:dist(r,y); return {(a^2 + c^2 - b^2)/(2*c), r} end
280
281  -- SKIP: summarizes things we want to ignore (so does nothing)
282  function SKIP.new(k,n,s) return new(k,{n=0,at=at or 0,txt=s or""}) end
283  function SKIP.add(i,x)   return x end
284  function SKIP.mid(i)     return "?" end
285
286  -- SOME: keeps a random sample on the arriving data
287  function SOME.new(k,keep) return new(k,{n=0,all={}, keep=keep or the.keep}) end
288  function SOME.add(i,x)
289    i.n = i.n+1
290    if     #i.all < i.keep then push(i.all,x)        ; return i.all
291    elseif r()      < i.keep/i.n then i.all[r(#i.all)]=x; return i.all end end
292
293  -- SYM: summarizes a stream of symbols
294  function SYM.new(k,n,s)
295    return new(k,{n=0,at=n or 0,txt=s or"",has={},most=0}) end
296
297  function SYM.add(i,x,inc)
298    if x ~= "?" then
299      inc = inc or 1
300      i.n = i.n + inc
301      i.has[x] = inc + (i.has[x] or 0)
302      if i.has[x] > i.most then i.most,i.mode=i.has[x],x end end end
303
304  function SYM.dist(i,x,y) return(x=="?" and y=="?" and 1) or(x==y and 0 or 1) end
305  function SYM.mid(i)        return i.mode end
306  function SYM.div(i,    p)
307    return sum(i.has,function(k) p=-i.has[k]/i.n;return -p*math.log(p,2) end) end
308
309  function SYM.merge(i,j,    k)
310    k = SYM:new(i.at,i.txt)
311    for x,n in pairs(i.has) do k:add(x,n) end
312    for x,n in pairs(j.has) do k:add(x,n) end
313    ei, ej, ejk= i:div(), j:div(), k:div()
314    if i.n==0 or j.n==0 or .99*ek <= (i.n*ei + j.n*ej)/k.n then
315      return k end end
```

```
316  --
317  -- ┌─┐ ┌   ┌ ┌─┐ ┌─┐ ┌─┐ ┌─┐
318  -- │   │   │ │ │ └─┐  │  ├─┘
319  --
```

```
320  -- CLUSTER: recursively divides data by clustering towards two distant points
321  function CLUSTER.new(k,egs,top)
322    local i,want,left,right
323    i     = new(k, {here=egs})
324    top   = top or egs
325    want = (#top.rows)^the.want
326    if #egs.rows >= 2*want then
327      left, right, i.x, i.y, i.c, i.mid = egs:half(top)
328      if #left.rows < #egs.rows then
329        i.left = CLUSTER:new(left,  top)
330        i.right= CLUSTER:new(right, top) end end
331    return i end
332
333  function CLUSTER.show(i,pre,  here)
334    pre = pre or ""
335    here=""
336    if not i.left and not i.right then here= o(i.here:mid(i.here.cols.y)) end
337    print(fmt("%6s: %-30s %s",#i.here.rows, pre, here))
338    for _,kid in pairs{i.left, i.right} do
339      if kid then kid:show(pre .. "|.. ") end end end
340
```

```
344
345  -- SPAN: keeps a random sample on the arriving data
346  function SPAN.new(k, col, lo, hi, has)
347    return new(k,{col=col,lo=lo,hi=hi or lo,has=has or SYM:new()}) end
348
349  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
350  function SPAN.merge(i,j)
351    local has = i.has:merge(j.has)
352    if now then return SPAN:new(i.col, i.lo, j.hi, has) end end
353
354  function SPAN.select(i,row,    x)
355    x = row[i.col.at]
356    return (x=="?") or (i.lo==i.hi and x==i.lo) or (i.lo <= x and x < i.hi) end
357
358  function SPAN.score(i) return i.has.n/i.col.n, i.has:div() end
359
360  function SPAN.scores(i, ss,ds)
361    size,div = i:score()
362    size,div = ss:norm(size), ds:norm(div)
363    return ((1 - size)^2 + (0 - div)^2)^.5 end
364
365  -- EXPLAIN:
366  function EXPLAIN.new(k,egs,top)
367    local i,n,y,ds,ss,top,div,want,size,left,span,right,spans
368    i    = new(k,{here = egs})
369    top  = top or egs
370    want = (#top.rows)^the.want
371    if #top.rows >= 2*want then
372      left,right   = egs:half(top)
373      spans, ds, ss = {}, Num(), Num()
374      for n,col in pairs(i.cols.x) do
375        for _,span in pairs(col:spans(j.cols.x[n])) do
376          push(spans, one)
377          size, div = span:score()
378          ss:add(size)
379          ds:add(div) end end
380      span= sort(spans,function(x,y)return x:scores(ss,ds)<y:scores(ss,ds) end)[1]
381      y, n = egs:clone(), egs:clone()
382      for _,row in pairs(egs.rows) do (span:selects(row) and y or n):add(row) end
383      if #y.rows<#egs.rows and #y.rows>want then i.yes=EXPLAIN:new(y,top) end
384      if #n.rows<#egs.rows and #n.rows>want then i.no =EXPLAIN:new(n,top) end end
385    return i end
386
387  function EXPLAN.show(i,pre)
388    pre = pre or ""
389    if not pre then
390      tmp = i.here:mid(i.here.y)
391      print(fmt("%6s: %-30s %s", #i.here.rows, pre, o(i.here:mid(i.here.cols.y))))
392
393    for _,pair in pairs{{true,i.yes},{false,i.no}} do
394      status,kid = unpack(pair)
395      k:shpw(pre .. "|.. ") end end
396
397  function SYM.spans(i, j)
398    local xys,all,one,last,xys,x,c n = {},{}
399    for x,n in pairs(i.has) do push(xys, {x,"this",n}) end
400    for x,n in pairs(j.has) do push(xys, {x,"that",n}) end
401    for _,tmp in ipairs(sort(xys,firsts)) do
402      x,c,n = unpack(tmp)
403      if x ~= last then
404        last = x
405        one  = push(all, Span(i,x,x)) end
406      one:add(x,y,n) end
407    return all end
408
409  function NUM.spans(i, j)
410    local xys,all,lo,hi,gap,xys,one,x,c,n = {},{}
411    lo,hi = min(i.lo, j.lo), max(i.hi,j.hi)
412    gap   = (hi - lo) / (6/the.cohen)
413    for x,n in pairs(i.has) do push(xys, {x,"this",1}) end
414    for x,n in pairs(j.has) do push(xys, {x,"that",1}) end
415    one = Span:new(i,lo,lo)
416    all = {one}
417    for _,tmp in ipairs(sort(xys,first)) do
418      x,c,n = unpack(tmp)
419      if one.hi - one.lo > gap then one = push(all, Span(i, one.hi, x)) end
420      one:add(x,y) end
421    all        = merge(all)
422    all[1   ].lo = -big
423    all[#all].hi =  big
424    return all end
425
426  function merge(b4,       j,n,now,a,b,merged)
427    j,n,now = 0,#b4,{}
428    while j < #b4 do
429      j    = j+1
430      a, b = b4[j], b4[j+1]
431      if b then
432        merged = a:merge(b)
433        if merged then a,j = merged, j+1 end end
434      push(now,a)
435      j = j+1 end
436    return #now == #b4 and b4 or merge(now) end
```

```lua
--    ___   ___  ___   ___
--   |  _) |  _)| | | (__
--   |_|   |___ |_|_|(___]
--
fails=0
function asserts(test, msg)
  print(test and "PASS:"or "FAIL: ",msg or "")
  if not test then
     fails=fails+1
     if the.dump then assert(test,msg) end end end

function EGS.nothing() return true end
function EGS.the()     oo(the) end
function EGS.rand()    print(r()) end
function EGS.some(s,t)
  s=SOME:new(100)
  for i=1,100000 do s:add(i) end
  for j,x in pairs(sort(s.all)) do
    --if (j % 10)==0 then print("") end
    --io.write(fmt("%6s",x))   end end
    fmt("%6s",x)   end end

function EGS.clone( r,s)
  r = ROWS:new(the.data)
  s = r:clone()
  for _,row in pairs(r.rows) do s:add(row) end
  asserts(r.cols.x[1].lo==s.cols.x[1].lo,"clone.lo")
  asserts(r.cols.x[1].hi==s.cols.x[1].hi,"clone.hi")
  end

function EGS.data( r)
  r = ROWS:new(the.data)
  asserts(r.cols.x[1].hi == 8, "data.columns") end

function EGS.dist( r,rows,n)
  r = ROWS:new(the.data)
  rows = r.rows
  n = NUM:new()
  for _,row in pairs(rows) do n:add(r:dist(row, rows[1])) end
  --oo(r.cols.x[2]:sorted()) end
  o(r.cols.x[2]:sorted()) end

function EGS.many(   t)
  t={}; for j=1,100 do push(t,j) end
  --print(oo(many(t, 10))) end
  o(many(t, 10)) end

function EGS.far(   r,c,row1,row2)
  r = ROWS:new(the.data)
  row1   = r.rows[1]
  c,row2 = r:far(r.rows[1], r.rows) end
  --print(c,"\n",o(row1),"\n", o(row2)) end

function EGS.half(   r,c,row1,row2)
  local lefts,rights,x,y,x
  r = ROWS:new(the.data)
  r:mid(r.cols.y)
  lefts,rights,x,y,c = r:half()
  lefts:mid(lefts.cols.y )
  rights:mid(rights.cols.y)
  asserts(true,"half") end

function EGS.cluster(r)
  r = ROWS:new(the.data)
  --CLUSTER:new(r):show() end
  CLUSTER:new(r) end

-- start-up
if arg[0] == "sl.lua" then
  oo(the)
  if the.help then print(help:gsub("\nNOTES:*$","")) else
    local b4={}; for k,v in pairs(the) do b4[k]=v end
    for _,todo in pairs(the.todo=="all" and slots(EGS) or {the.todo}) do
      for k,v in pairs(b4) do the[k]=v end
      math.randomseed(the.seed)
      if type(EGS[todo])=="function" then EGS[todo]() end end
  end
  for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
  os.exit(fails)
else
  return {CLUSTER=CLUSTER, COLS=COLS, NUM=NUM, ROWS=ROWS,
          SKIP=SKIP, SOME=SOME, SYM=SYM,the=the,oo=oo,o=o}
end
-- git rid of SOME for rows
-- nss  = NUM | SYM | SKIP
-- COLS = all:[nss]+, x:[nss]*, y:[nss]*, klass;col?
-- ROWS = cols:COLS, rows:SOME
-- ## References
-- - [Ah91]:
-- Aha, D.W., Kibler, D. & Albert, M.K. Instance-based
-- learning algorithms. Mach Learn 6, 37â€“66 (1991).
-- https://doi.org/10.1007/BF00153759
-- - [Boley, 1998]:
--   Boley, D., 1998.
-- [Principal directions divisive partitioning](https://www-users.cse.umn.edu/~b
oley/publications/papers/PDDP.pdf)
--   Data Mining and Knowledge Discovery, 2(4): 325-344.
-- - [Ch05]:
-- [Semi-Supervised Learning](http://www.molgen.mpg.de/3659531/MITPress--SemiSup
ervised-Learning)
-- (2005) Olivier Chapelle,  Bernhard SchÃ¶lkopf, and Alexander Zien (eds).
-- MIT Press.
--  - [Ch18]
-- [Samplingâ€” as a Baseline Optimizer for Search-Based Software Engineer
ing](https://arxiv.org/pdf/1608.07617.pdf),
-- Jianfeng Chen; Vivek Nair; Rahul Krishna; Tim Menzies
-- IEEE Trans SE, (45)6, 2019
-- - [Ch22]:
-- [Can We Achieve Fairness Using Semi-Supervised Learning?](https://arxiv.org/p
df/2111.02038.pdf)
-- (2022), Joymallya Chakraborty, Huy Tu, Suvodeep Majumder, Tim Menzies.
-- - [Fal95]:
-- Christos Faloutsos and King-Ip Lin. 1995. FastMap: a fast algorithm for index
ing, data-mining and visualization of traditional and multimedia datasets. SIGMO
D Rec. 24, 2 (May 1995), 163â€“174. DOI:https://doi.org/10.1145/568271.223
812
-- - [Le05}
-- Levina, E., Bickel, P.J.: [Maximum likelihood estimation of intrinsic dimensi
on](https://www.stat.berkeley.edu/~bickel/mldim.pdf).
-- In:
-- Advances in neural information processing systems, pp. 777â€“784 (2005)
-- - [Pl04]:
-- Platt, John.
-- [FastMap, MetricMap, and Landmark MDS are all Nystrom Algorithms](https://www
.microsoft.com/en-us/research/wp-content/uploads/2005/01/nystrom2.pdf)
-- AISTATS (2005).
-- - [Zit04]:
-- [Indicator-based selection in multiobjective search](https://link.springer.co
m/chapter/10.1007/978-3-540-30217-9_84)
-- Eckart Zitzler , Simon KÃ¼nzli
-- Proc. 8th International Conference on Parallel Problem Solving from Nature (P
PSN VIII
```