```
1    ----------------------------------------------------------------------
2    ---    __      __                                              ,o88888
3    ---   /\ \    /\ \                                           ,o8888888'
4    ---   \ \ \   \ \ \__                          ,:o:o:oooo.      ,8O88Pd8888"
5    ---    \ \ \   \ \  _`\                      ,.::.::o:ooooOoOoO. ,oO8O8Pd888'"
6    ---     \ \ \__L\ \  \/\ \                  ,.:.::o:ooOoOoOO8O8OOo.8OOPd8O8O"
7    ---      \ \_____\ \  \ \_\                ,.::.::o:ooOoOoOOoEOOOOo.FdO8O8"
8    ---       \/_____/\ \  \/_/               , ..:.::o:ooOoOOOO8OOOOo.Fd8O8"
9    ---                \/___/                . ..:.::o:ooOoOoOO8O8OOOoEoO8"
10   ---   L5=a little LUA learning library  , ..:.::o:ooOoOooE0OOOOCOCO"
11   ---   (c) Tim Menzies 2022, BSD-2       ..:.::o:ooOoOoOO8O8OCCCC"o
12   ---   Share and enjoy.                  . ..:.::o:ooOoCoCCC"o:o
13   ---   https://menzies.us/l5             `. ..:.::o:ooooCo"oo:o:
14   ---                                       `  . . ..:.::cocoooo"'o:o::'
15   ---                                        .`   . ..::cccoc"'o:o:o::'
16   ---                                         :.:.   ,c:cccc"':.:.:.:.'
17   ---                                          ..:.::"'`:::c:"'..:.:.:.:.'
18   ---                                           ....:'.:.::::"'    . . . . '
19   ---                                            ...:::.::"/  '  . . . ''
20   ---                                             . .  ....'"/
21   ---                                              .. . .."'       -hrr-
22   ---                                               .
23   ---
24   ---
25   local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
26   local the,help={},[[
27
28   lua 15.lua [OPTIONS]
29   L5 == a very little LUA learning lab
30   (c)2022, Tim Menzies, BSD 2-clause license
31
32   OPTIONS (for changing the inference):
33
34    -cohen  -c  F  cohen's small effect size     = .35
35    -far    -F  F  look no further than "far"    = .9
36    -keep   -k     items to keep in a number     = 512
37    -leaves -l     leaf size                     = .5
38    -p      -p  P  distance calcs coefficient    = 2
39    -seed   -S  P  random number seed            = 10019
40    -some   -s     look only at "some" items     = 512
41
42   OPTIONS (for housekeeping):
43
44    -dump   -d     exit on error, with stacktrace = false
45    -file   -f  S  where to get data              = ../etc/data/auto93.csv
46    -help   -h     show help                      = false
47    -rnd    -r  S  format string                  = %5.2f
48    -todo   -t  S  start-up action                = nothing
49
50
51   KEY: S=string, P=poisint, F=float
52   ]]
53   local as,o = setmetatable
54   local function obj(  t)
55     t={__tostring=o}; t.__index=t
56     return as(t, {__call=function(_,...) return t.new(_,...) end}) end
57   ---
58   ---      _   _   _
59   ---     |_| |_| |_|
60   ---     |_| |_| |_|
61   ---
62
63   local Sym = obj() -- Where to summarize symbols
64   function Sym:new(at,s) return as({
65     is="Sym",       -- type
66     at=at or 0,     -- column index
67     name=s or "",   -- column name
68     n=0,            -- number of items summarized in this column
69     all={},         -- all[x] = n means we've seen "n" repeats of "x"
70     most=0,         -- count of the most frequently seen symbol
71     mode=nil        -- the most commonly seen letter
72     }, Sym) end
73
74   local Num = obj() -- Where to summarize numbers
75   function Num:new(at,s) return as({
76     is="Num",       -- type
77     at=at or 0,     -- column index
78     name=s or "",   -- column name
79     n=0,            -- number of items summarizes in this column
80     mu=0,           -- mean (updated incrementally)
81     m2=0,           -- second moment (updated incrementally)
82     sd=0,           -- standard deviation
83     all={},         -- a sample of items seen so far
84     lo=1E31,        -- lowest number seen; initially, big so 1st num sends it low
85     hi=-1E31,       -- highest number seen;initially, msall to 2st num sends it hi
86     w=(s or ""):find"-$" and -1 or 1 -- "-1"= minimize and "1"= maximize
87     }, Num) end
88
89   local Egs = obj() -- Where to store examples, summarized into Syms or Nums
90   function Egs:new(names,    i,col,here)  i=as({
91     is="Egs",       -- type
92     all={},         -- all the rows
93     names=names,    -- list of name
94     cols={},        -- list of all columns  (Nums or Syms)
95     x={},           -- independent columns (nothing marked as "skip")
96     y={}            -- dependent columns (nothing marked as "skip")
97     },Egs)
98     for at,name in pairs(names) do
99       col = (name:find"^[A-Z]" and Num or Sym)(at,name)
100      i.cols[1+#i.cols] = col
101      here = name:find"[-+]$" and i.y or i.x
102      if not name:find":$" then here[1 + #here] = col end end
103    return i end
104   ---
105   ---     _   _   _   _   _   _   _
106   ---    |_| |_| |\| | | |\| | | |_|
107
108   function Num.clone(i) return Num(i.at, i.name) end
109   function Sym.clone(i) return Sym(i.at, i.name) end
110
111   local data
112   function Egs.clone(i,rows,    copy)
113     copy = Egs(i.names)
114     for _,row in pairs(rows or {}) do  data(copy,row)  end
115     return copy end
116
117   --[[
118   ## Coding Conventions
119   - "i" not "self"
120   - if something holds a list of thing, name the holding variable "all"
121   - no inheritance
122   - only define a method if that is for polymorphism
123   - when you can, write functions down on one line
124   - all config items into a global "the" variable
125   - all the test cases (or demos) are "function Demo.xxx".
126   - random seed reset so carefully, just once, at the end of the code.
127   - usually, no line with just "end" on it
128   ]]
```

```
129   ---
130   ---       _   _   _   _   _
131   ---      | | | | |_  | | |_
132   ---      |_| |_|  _| |_| |_
133   ---
134   ----------------------------------------------------------------------
135   local r    = math.random
136   local fmt = string.format
137   local unpack = table.unpack
138   local function push(t,x) table.insert(t,x); return x end
139   ---
140   ---       _   _   _   _   _   _
141   ---      |   | | |_  |_| |   |_
142
143   local thing,things,file2things
144   function thing(x)
145     x = x:match"^%s*(.-)%s*$"
146     if x=="true" then return true elseif x=="false" then return false end
147     return tonumber(x) or x end
148
149   function things(x,sep,  t)
150     t={}; for y in x:gmatch(sep or"([^,]+)") do push(t,thing(y)) end
151     return t end
152
153   function file2things(file,      x)
154     file = io.input(file)
155     return function()
156       x=io.read();
157       if x then return things(x) else io.close(file) end end end
158   ---
159   ---      _   _   _   _   _   _
160   ---     |_| |_  |   _| |_  |
161   ---
162   local last,per,any,many
163   function last(a)        return a[ #a ] end
164   function per(a,p)       return a[ (p*#a)//1 ] end
165   function any(a)         return a[ math.random(#a) ] end
166   function many(a,n,  u) u={}; for j=1,n do push(u,any(a)) end; return u end
167   ---
168   ---      _   _   _
169   ---     |   | | |_
170
171   local firsts,sort,map,slots
172   function firsts(a,b)   return a[1] < b[1] end
173   function sort(t,f)     table.sort(t,f); return t end
174   function map(t,f,  u)  u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
175   function slots(t, u,s)
176     u={}
177     for k,v in pairs(t) do s=tostring(k);if s:sub(1,1)~="_" then push(u,k) end end
178     return sort(u) end
179   ---
180   ---      _   _   _   _   _
181   ---     |   | | |_| |\| |
182
183   local oo,rnd, rnds -- local o was declared above (in "new")
184   function oo(t)  print(o(t)) end
185   function o(t,seen,         key,xseen,u)
186     seen = seen or {}
187     if type(t)~="table" then return tostring(t) end
188     if seen[t]            then return "..." end
189     seen[t] = t
190     key   = function(k) return fmt(":%s %s",k,o(t[k],seen)) end
191     xseen = function(x) return o(x,seen) end
192     u = #t>0 and map(t,xseen) or map(slots(t),key)
193     return (t.is or "")..'{'..table.concat(u,"")..'}' end
194
195   function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
196   function rnd(x,f)
197     return fmt(type(x)=="number" and (x~=x//1 and f or the.rnd) or "%s",x) end
198   ---
199   ---      _   _   _   _   _   _   _
200   ---     |_| | | |_| |_| |   _| |_|
201
202   local Demo, ok = {fails=0}
203   function ok(test,msg)
204     print(test and "PASS:"or "FAIL:",msg or "")
205     if not test then
206       Demo.fails=Demo.fails+1
207       if the.dump then assert(test,msg) end end end
208
209   function Demo.main(todo,seed)
210     for k,one in pairs(todo=="all" and slots(Demo) or {todo}) do
211       if k ~= "main" and type(Demo[one]) == "function" then
212         math.randomseed(seed)
213         Demo[one]() end end
214     for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
215     return Demo.fails
216
217   local function settings(txt,  d)
218     d={}
219     txt:gsub("\n ([-]([^%s]+))[%s]+(-[^%s]+)[^\n]*%s([^%s]+)",
220       function(long,key,short,x)
221         for n,flag in ipairs(arg) do
222           if flag==short or flag==long then
223             x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
224         if x=="false" then the[key]=false elseif x=="true" then the[key]=true else
225         d[key] = tonumber(x) or x end end)
226     if d.help then print(help) end
227     return d end
```

page 2

```lua
228 ---
229 --  UPDATE COLS
230 ---
231
232 local add
233 function add(i,x, inc)
234   inc = inc or 1
235   if x ~= "?" then
236     i.n = i.n + inc
237     i:internalAdd(x,inc) end
238   return x end
239
240 function Sym.internalAdd(i,x,inc)
241   i.all[x] = inc + (i.all[x] or 0)
242   if i.all[x] > i.most then i.most, i.mode = i.all[x], x end end
243
244 function Num.internalAdd(i,x,inc,     d)
245   for j=1,inc do
246     d    = x - i.mu
247     i.mu = i.mu + d/i.n
248     i.m2 = i.m2 + d*(x - i.mu)
249     i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n-1))^0.5)
250     i.lo = math.min(x, i.lo)
251     i.hi = math.max(x, i.hi)
252     if     #i.all < the.keep       then push(i.all,x)
253     elseif r()    < they.keep/i.n then i.all[r(#i.all)]=x end end end
254 ---
255 --  MAKE DATA
256 ---
257
258 local file2Egs -- not "local data" (since defined above)
259 function data(i,row)
260   push(i.all, row)
261   for _,col in pairs(i.cols) do add(col, row[col.at]) end
262   return i end
263
264 function file2Egs(file,    i)
265   for row in file2things(file) do
266     if i then data(i,row) else i = Egs(row) end end
267   return i end
268 ---
269 --  SUMMARIZE
270 ---
271
272 local mids
273 function mids(i,rows,cols) return i:clone(rows):mid(cols) end
274
275 function Egs.mid(i,cols)
276   return map(cols or i.y,function(col) return col:mid() end) end
277
278 function Sym.mid(i) return i.mode end
279 function Num.mid(i) return i.mu end
280
281 function Num.div(i) return i.sd end
282 function Sym.div(i,   e)
283   e=0; for _,n in pairs(i.all) do e=e + n/i.n*math.log(n/i.n,2) end
284   return -e end
285 ---
286 --  DISTANCE
287 ---
288
289 local far,furthest,neighbors,dist
290 function far(        i,r1,rows,far)
291   return per(neighbors(i,r1,rows),far or the.far)[2] end
292
293 function furthest( i,r1,rows)
294   return last(neighbors(i,r1,rows))[2] end
295
296 function neighbors(i,r1,rows)
297   return sort(map(rows, function(r2) return {dist(i,r1,r2),r2} end),firsts) end
298
299 function dist(i,row1,row2,     d,n,a,b,inc)
300   d,n = 0,0
301   for _,col in pairs(i.x) do
302     a,b = row1[col.at], row2[col.at]
303     inc = a=="?" and b=="?" and 1 or col:dist1(a,b)
304     d = d + inc^the.p
305     n = n + 1 end
306   return (d/n)^(1/the.p) end
307
308 function Sym.dist1(i,a,b) return a==b and 0 or 1 end
309
310 function Num.dist1(i,a,b)
311   if     a=="?" then b=i:norm(b); a=b<.5 and 1 or 0
312   elseif b=="?" then a=i:norm(a); b=a<.5 and 1 or 0
313   else   a,b = i:norm(a), i:norm(b)   end
314   return math.abs(a - b) end
315
316 function Num.norm(i,x)
317   return i.hi - i.lo < 1E-32 and 0 or (x - i.lo)/(i.hi - i.lo) end
318 ---
319 --  CLUSTER
320 ---
321
322 local half, cluster, clusters
323 function half(i, rows,       project,row,some,left,right,lefts,rights,c,mid)
324   function project(row,a,b)
325     a= dist(i,left,row)
326     b= dist(i,right,row)
327     return {(a^2 + c^2 - b^2)/(2*c), row}
328   end ------------------------
329   some  = many(rows,        the.some)
330   left  = furthest(i,any(some), some)
331   right = furthest(i,left,       some)
332   c     = dist(i,left,right)
333   lefts,rights = {},{}
334   for n, projection in pairs(sort(map(rows,project),firsts)) do
335     if n==#rows//2 then mid=row end
336     push(n <= #rows//2 and lefts or rights, projection[2]) end
337   return lefts, rights, left, right, mid, c   end
338
339 function cluster(i,rows,   here,lefts,rights)
340   rows = rows or i.all
341   here = {all=rows}
342   if #rows >= 2 * (#i.all)^the.leaves then
343     lefts, rights, here.left, here.right, here.mid = half(i, rows)
344     if #lefts < #rows then
345       here.lefts = cluster(i,lefts)
346       here.rights= cluster(i,rights) end end
347   return here end
348
349 function clusters(i,format,t,pre,    front)
350   if t then
351     pre=pre or ""
352     front = fmt("%s%s",pre,#t.all)
353     if not t.lefts and not t.rights then
354       print(fmt("%-20s%s",front, o(rnds(mids(i,t.all),format))))
355     else
356       print(front)
357       clusters(i,format,t.lefts, "|".. pre)
358       clusters(i,format,t.rights,"|".. pre) end end end
```

```lua
359 ---
360 --  DISCRETIZE
361 ---
362
363 local merge,merged,spans,bestSpan
364 function Sym.spans(i, j)
365   local xys,all,one,last,x,y,n = {}, {}
366   for x,n in pairs(i.all) do push(xys, {x,"lefts",n}) end
367   for x,n in pairs(j.all) do push(xys, {x,"rights",n}) end
368   for _,tmp in ipairs(sort(xys,firsts)) do
369     x,y,n = unpack(tmp)
370     if x ~= last then
371       last = x
372       one  = push(all, {lo=x, hi=x, all=Sym(i.at,i.name)}) end
373     add(one.all, y, n) end
374   return all end
375
376 function Num.spans(i, j)
377   local xys,all,lo,hi,gap,one,x,y,n = {},{}
378   lo,hi = math.min(i.lo, j.lo), math.max(i.hi, j.hi)
379   gap   = (hi - lo) / (6/the.cohen)
380   for _,n in pairs(i.all) do push(xys, {n,"lefts",1}) end
381   for _,n in pairs(j.all) do push(xys, {n,"rights",1}) end
382   one = {lo=lo, hi=lo, all=Sym(i.at,i.name)}
383   all = {one}
384   for _,tmp in ipairs(sort(xys,firsts)) do
385     x,y,n = unpack(tmp)
386     if   one.hi - one.lo > gap then
387       one = push(all, {lo=one.hi, hi=x, all=one.all:clone()})
388     end
389     one.hi = x
390     add(one.all, y, n) end
391   all      = merge(all)
392   all[1   ].lo = -math.huge
393   all[#all].hi =  math.huge
394   return all end
395
396 function merge(b4,        j,n,now,a,b,both)
397   j, n, now = 0, #b4, {}
398   while j < #b4 do
399     j     = j+1
400     a, b = b4[j], b4[j+1]
401     if   b then
402       both = a.all:merge(b.all)
403       if both then
404         a = {lo=a.lo, hi=b.hi, all=both}
405         j = j + 1 end end
406     push(now,a) end
407   return #now == #b4 and b4 or merge(now) end
408
409 function Sym.merge(i,j,     k,ei,ej,ek)
410   k = i:clone()
411   for x,n in pairs(i.all) do add(k,x,n) end
412   for x,n in pairs(j.all) do add(k,x,n) end
413   ei, ej, ek= i:div(), j:div(), k:div()
414   if    ek*.99 <= (i.n*ei + j.n*ej)/k.n then
415     return k end end
416
417 function spans(egs1,egs2,        spans,tmp,col1,col2)
418   spans = {}
419   for c,col1 in pairs(egs1.x) do
420     col2 = egs2.x[c]
421     tmp = col1:spans(col2)
422     if #tmp> 1 then
423       for _,one in pairs(tmp) do push(spans,one) end end end
424   return spans end
425
426 function bestSpan(spans)
427   local divs,ns,n,div,stats,dist2heaven = Num(), Num()
428   function dist2heaven(s) return {((1 - n(s))^2 + (0 - div(s))^2)^.5,s} end
429   function div(s)         return divs:norm( s.all:div() ) end
430   function n(s)           return   ns:norm( s.all.n     ) end
431   for _,s in pairs(spans) do
432     add(divs, s.all:div())
433     add(ns,    s.all.n) end
434   return sort(map(spans, dist2heaven), firsts)[1][2]   end
435 ---
436 --  EXPLAIN
437 ---
438
439 local xplain,xplains,selects,spanShow
440 function xplain(i,rows,used,
441               stop,here,left,right,lefts0,rights0,lefts1,rights1)
442   used=used or {}
443   rows = rows or i.all
444   here = {all=rows}
445   stop = (#i.all)^the.leaves
446   if #rows >= 2*stop then
447     lefts0, rights0, here.left, here.right, here.mid, here.c  = half(i, rows)
448     if #lefts0 < #rows then
449       here.selector = bestSpan(spans(i:clone(lefts0),i:clone(rights0)))
450       push(used, {here.selector.all.name, here.selector.lo, here.selector.hi})
451       lefts1,rights1 = {},{}
452       for _,row in pairs(rows) do
453         push(selects(here.selector, row) and lefts1 or rights1, row) end
454       if #lefts1 > stop then here.lefts  = xplain(i,lefts1,used) end
455       if #rights1 > stop then here.rights = xplain(i,rights1,used) end end end
456   return here end
457
458 function xplains(i,format,t,pre,how,     sel,front)
459   pre, how = pre or "", how or ""
460   if t then
461     pre=pre or ""
462     front = fmt("%s%s%s %s",pre,how, #t.all, t.c and rnd(t.c) or "")
463     if t.lefts and t.rights then
464       print(fmt("%-35s",front))
465     else
466       print(fmt("%-35s %s",front, o(rnds(mids(i,t.all),format))))
467     end
468     sel = t.selector
469     xplains(i,format,t.lefts,  "|".. pre, spanShow(sel)..":")
470     xplains(i,format,t.rights, "|".. pre, spanShow(sel,true) ..":") end end
```

```lua
471
472  function selects(span,row,      lo,hi,at,x)
473    lo, hi, at = span.lo, span.hi, span.all.at
474    x = row[at]
475    if x=="?" then return true end
476    if lo==hi then return x==lo else return lo <= x and x < hi end end
477
478  function spanShow(span, negative,    hi,lo,x,big)
479    if not span then return "" end
480    lo, hi, x, big  = span.lo, span.hi, span.all.name, math.huge
481    if not negative then
482      if lo ==  hi  then return fmt("%s == %s",x,lo)   end
483      if hi ==  big then return fmt("%s >= %s",x,lo)   end
484      if lo == -big then return fmt("%s < %s",x,hi)   end
485      return fmt("%s <= %s < %s",lo,x,hi)
486    else
487      if lo ==  hi  then return fmt("%s != %s",x,lo)   end
488      if hi ==  big then return fmt("%s < %s",x,lo)   end
489      if lo == -big then return fmt("%s >= %s",x,hi)   end
490      return fmt("%s < %s and %s >= %s", x,lo,x,hi)   end end
491  ---
492  ---       _ _ __   _ ___  ()_ _____
493  ---      /  ` \ / / `/ _ \| `  \  ' \
494  ---     |   | | | |   | | |  | | | |  |
495  ---     |_| |_|_\__, |_| |_|_| |_| |_|
496
497  function Demo.the() oo(the) end
498
499  function Demo.many(a)
500    a={1,2,3,4,5,6,7,8,9,10}; ok("{10 2 3}" == o(many(a,3)), "manys") end
501
502  function Demo.egs()
503    ok(5140==file2Egs(the.file).y[1].hi,"reading") end
504
505  function Demo.dist(i)
506    i = file2Egs(the.file)
507    for n,row in pairs(i.all) do print(n,dist(i, i.all[1], row)) end end
508
509  function Demo.far(  i,j,row1,row2,row3,d3,d9)
510    i = file2Egs(the.file)
511    for j=1,10 do
512      row1 = any(i.all)
513      row2  = far(i,row1, i.all, .9)
514      d9    = dist(i,row1,row2)
515      row3 = far(i,row1, i.all, .3)
516      d3    = dist(i,row1,row3)
517      ok(d3 < d9, "closer far") end end
518
519  function Demo.half(  i,lefts,rights)
520    i = file2Egs(the.file)
521    lefts,rights = half(i, i.all)
522    oo(mids(i, lefts))
523    oo(mids(i, rights))
524    end
525
526  function Demo.cluster(   i)
527    i = file2Egs(the.file)
528    clusters(i,"%.0f",cluster(i)) end
529
530  function Demo.spans(    i,lefts,rights)
531    i = file2Egs(the.file)
532    lefts, rights = half(i, i.all)
533    oo(bestSpan(spans(i:clone(lefts), i:clone(rights)))) end
534
535  function Demo.xplain(    i,j,tmp,lefts,rights,used)
536    i = file2Egs(the.file)
537    used={}
538    xplains(i,"%.0f",xplain(i, i.all,used))
539    map(sort(used,function(a,b)
540           return ((a[1] < b[1]) or
541                 (a[1]==b[1] and a[2] < b[2]) or
542                 (a[1]==b[1] and a[2]==b[2] and a[3] < b[3]))end),oo) end
543
544
545  --------------------------------------------------------------------------------
546  the = settings(help)
547  Demo.main(the.todo, the.seed)
```