

```

1 -----
2 ---
3 ---
4 ---
5 ---
6 ---
7 ---
8 ---
9 ---
10 ---
11 ---
12 ---
13 ---
14 ---
15 ---
16 ---
17 ---
18 local b4={}; for k, _ in pairs(_ENV) do b4[k]=k end
19 local the, help = {}, {}
20
21 lua brknbad.lua [OPTIONS]
22 (c) 2022, T. Menzies, BSD-2-Clause
23 Divide things. Show deltas between things.
24
25 OPTIONS:
26 -cohen          -c cohen          = .35
27 -far            -F how far to seek poles = .9
28 -keep          -k items to keep   = 256
29 -minitems      -m min items in a rang e = .5
30 -p            -p euclidean coefficient =
31 -some         -S sample size for rows = 512
32
33 OPTIONS, other:
34 -dump          -d stackdump on error = false
35 -file          -f data file         = ../etc/data/auto93.csv
36 -help         -h show help         = false
37 -rnd          -r round numbers     = %5.2f
38 -seed         -s random number seed = 10019
39 -todo         -t start-up action   = nothing
40
41 ]]
42 local any, bestBin, bins, bins1, bootstrap, class, csv2egs, firsts, fmt, ish
43 local last, many, map, new, o, oo, per, push, quintiles, r, rnd, rnds, scottKnot
44 local selects, settings, slots, smallfx, sort, sum, thing, things, xplains
45 local NUM, SYM, EGS, BIN, CLUSTER, XPLAIN, GO
46
47 --[[
48
49 ## Conventions
50
51 ### Data
52
53 - First row of data are names that describe each column.
54 - Names ending with '['+-]' are dependent goals to be minimized or maximized.
55 - Names ending with '!' are dependent classes.
56 - Dependent columns are 'y' columns (the rest are independent 'x' columns).
57 - Uppercase names are numeric (so the rest are symbolic).
58 - Names ending with '.' are columns to be skipped.
59 - Data is read as rows, stored in an EGS instance.
60 - Within an EGS, row columns are summarized into NUM or SYM instances.
61
62 ### Inference
63
64 - The rows within an EGS are recursive bi-clustered into CLUSTERS
65 using random projections (Fastmap) and Aha's distance metric
66 (that can process numbers and symbols).
67 - Entropy-based discretization finds BINs that separates each pair of
68 clusters.
69 - An XPLAIN tree runs the same clustering processing, but data is divided
70 at level using the BIN that most separates the clusters.
71
72 ### Coding
73
74 - No globals (so everything is 'local').
75 - Code 80 clide indent with two spaces.
76 - Format to be read a two-pages-per-page portrait pdf.
77 - Divide code into section and subsection headings (e.g using figlet)
78 - Sections are less than 120 lines long (one column in the pdf).
79 - No lines containing only the word 'end' (unless marking the end of a
80 complex for loop or function).
81 - Usually, if an object contains a list of other objects, that sublist
82 is called 'all'.
83 - If a slot is too big to display, it is declared private (not to be printed)
84 by renaming (e.g.) 'slotx' to '_slotx' (so often, 'all' becomes '_all').
85
86 ### Classes
87
88 - Spread class code across different sections (so don't overload reader
89 with all details, at one time).
90 - Show simpler stuff before complex stuff.
91 - Reserve 'i' for 'self' (to fit more code per line).
92 - Don't use inheritance (to simplify readability).
93 - Use polymorphism (using LUA's delegation trick).
94 - Define an class of objects with 'Thing=class"Thing"' and
95 a 'function:Thing(args)' creation method.
96 - Define instances with 'new({slot1=value1, slot2=value2,...}, Thing)'.
97 - Instance methods use '.'; e.g. 'function Thing.show(i) ... end'.
98 - Class methods using ':'; e.g. 'Thing:new4strings'. Class methods
99 do things like instance creation or manage a set of instances.
100
101 ### Test suites (and demos)
102
103 - Define start-up actions as 'go' functions.
104 - In 'go' functions, check for errors with 'ok(test,mdf)'
105 (that updates an 'fails' counter when not 'ok').
106
107 ### At top of file
108
109 - Trap known globals in 'b4'.
110 - Define all locals at top-of-file (so everyone can access everything).
111 - Define options in a help string at top of file.
112 - Define command line options -h (for help); -s (for seeding random numbers)
113 -t (for startup actions, so -t all' means "run everything").
114
115 ### At end of file
116
117 - Using 'settings', parse help string to set options,
118 maybe updating from command-line.
119 - Using 'GO.main', run the actions listed on command line.
120 - 'GO.main' resets random number generator before running an action
121 - After everything else, look for 'roguess' (any global not in 'b4')
122 - Finally, return the 'fails' as the exit status of this code. --]]

```

```

232 -----
233 DATA CLASSES
234
235 NUM, SYM, EGS = class"NUM", class"SYM", class"EGS"
236
237 --- create
238
239 function SYM:new(at,name)
240     return new({at=at, name=name, most=0,n=0,all={}}, SYM) end
241
242 function NUM:new(at,name)
243     return new({at=at, name=name, _all={},
244         w=(name or ""):find"$" and -1 or 1,
245         n=0, sd=0, mu=0, m2=0, lo=math.huge, hi=-math.huge}, NUM) end
246
247 function EGS:new(names, i,col)
248     i = new({_all={}, cols={names=names, all={}, x={}, y={}}, EGS)
249     for at,name in pairs(names) do
250         col = push(i.cols.all, (name:find"^[A-Z]" and NUM or SYM) (at,name) )
251         if not name:find"$" then
252             if name:find"$" then i.cols.class = col end
253             push(name:find"[+!]"$ and i.cols.y or i.cols.x, col) end end
254     return i end
255
256 function EGS:new4file(file, i)
257     for row in things(the.file) do
258         if i then i:add(row) else i = EGS(row) end end
259     return i end
260
261 --- copy
262
263 function SYM.copy(i) return SYM(i.at, i.name) end
264
265 function NUM.copy(i) return NUM(i.at, i.name) end
266
267 function EGS.copy(i,rows, j)
268     j = EGS(i.cols.names)
269     for _,row in pairs(rows or {}) do j:add(row) end
270     return j end
271
272 --- update
273
274 function EGS.add(i,row)
275     push(i._all, row)
276     for at,col in pairs(i.cols.all) do col:add(row[col.at]) end 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.all[x] = i.all[x] + (i.all[x] or 0)
283         if i.all[x] > i.most then i.most, i.mode = i.all[x], x end end end
284
285 function SYM.sub(i,x,inc)
286     if x ~= "?" then
287         inc = inc or 1
288         i.n = i.n - inc
289         i.all[x] = i.all[x] - inc end end
290
291 function NUM.add(i,x,_, d,a)
292     if x ~= "?" then
293         i.n = i.n + 1
294         d = x - i.mu
295         i.mu = i.mu + d/i.n
296         i.m2 = i.m2 + d*(x - i.mu)
297         i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n - 1))^0.5)
298         i.lo = math.min(x, i.lo)
299         i.hi = math.max(x, i.hi)
300         a = i._all
301         if #a < the.keep then i.ok=false; push(a,x)
302         elseif r() < the.keep/i.n then i.ok=false; a[r(#a)]=x end end end
303
304 function NUM.sub(i,x,_, d)
305     if x ~= "?" then
306         i.n = i.n - 1
307         d = x - i.mu
308         i.mu = i.mu - d/i.n
309         i.m2 = i.m2 - d*(x - i.mu)
310         i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n - 1))^0.5) end end
311
312 --- quality
313
314 function EGS.better(i,row1,row2)
315     local s1,s2,n,a,b=0,0,#i.cols.y
316     for _,col in pairs(i.cols.y) do
317         a = col:norm( row1[col.at] )
318         b = col:norm( row2[col.at] )
319         s1 = s1 - 2.7183*(col.w * (a - b) / n)
320         s2 = s2 - 2.7183*(col.w * (b - a) / n) end
321     return s1 / n < s2 / n end
322
323 function EGS.bettors(i,j,k)
324     return i:better(j:mid(j.cols.all), k:mid(k.cols.all)) end
325
326 function EGS.mid(i,cols)
327     return map(cols or i.cols.y, function(col) return col:mid() end) end
328
329 function NUM.mid(i) return i.mu end
330
331 function SYM.mid(i) return i.mode end
332
333 function NUM.div(i) return i.sd end
334
335 function SYM.div(i, e)
336     e=0; for _,n in pairs(i.all) do
337         if n > 0 then e = e - n/i.n * math.log(n/i.n,2) end end
338     return math.abs(e) end
339
340 function NUM.norm(i,x)
341     return i.hi - i.lo < 1E-32 and 0 or (x - i.lo)/(i.hi - i.lo) end
342
343 function NUM.all(i)
344     if not i.ok then table.sort(i._all); i.ok=true end
345     return i._all end
346
347

```

```

353 -----
354 CLUSTER
355
356 $ lua brknbad.lua -t cluster
357
358 ---
359
360 398
361 199
362 99
363 49
364 24
365 25
366 50
367 25
368 25
369 100
370 50
371 25
372 25
373 50
374 25
375 25
376 199
377 99
378 49
379 24
380 25
381 50
382 25
383 25
384 100
385 50
386 25
387 25
388 50
389 25
390 25
391 25
392 CLUSTER=class"CLUSTER"
393 function CLUSTER:new(top,egs, i, lefts, rights)
394     egs = egs or top
395     i = new({egs=egs, top=top}, CLUSTER)
396     if #egs._all >= 2*(#top._all)^the.minItems then
397         lefts, rights, i.left, i.right, i.mid, i.c = top:half(egs._all)
398         if #lefts._all < #egs._all then
399             i.lefts = CLUSTER(top, lefts)
400             i.rights = CLUSTER(top, rights) end end
401     return i end
402
403 function CLUSTER.leaf(i) return not (i.lefts or i.rights) end
404
405 function CLUSTER.show(i, pre, front)
406     pre = pre or ""
407     local front = fmt("%s",pre,#i.egs._all)
408     if i:leaf()
409     then print(fmt("%-20s",front, o(rnds(i.egs:mid(i.egs.cols.y))))))
410     else print(front)
411         if i.lefts then i.lefts:show(" |"..pre)
412         if i.rights then i.rights:show(" |"..pre) end end end end
413
414 --- random projections
415
416 function EGS.half(i, rows)
417     local project,far,some,left,right,c,lefts,rights
418     rows = rows or i._all
419     far = function(r,t) return per(i:dist(r,t), the.far)[2] end
420     project = function(r1, a,b)
421         a,b = i:dist(left,r1), i:dist(right,r1)
422         return {(a^2 + c^2 - b^2)/(2*c), r1} end
423     some = many(rows, the.some)
424     left = far(any(some), some)
425     right = far(left, some)
426     c = i:dist(left,right)
427     lefts,rights = i:copy(), i:copy()
428     for n, projection in pairs(sort(map(rows,project),firsts)) do
429         if n==#rows//2 then mid=row end
430         (n <= #rows//2 and lefts or rights):add( projection[2] ) end
431     return lefts, rights, left, right, mid, c end
432
433 --- distances in data
434
435 function EGS.dists(i,r1,rows)
436     return sort(map(rows,function(r2) return {i:dist(r1,r2),r2} end),firsts) end
437
438 function EGS.dist(i,row1,row2, d)
439     d = sum(i.cols.x, function(c) return c:dist(row1[c.at], row2[c.at])^the.p end)
440     return (d/#i.cols.x)^(1/the.p) end
441
442 function NUM.dist(i,a,b)
443     if a=="?" and b=="?" then return 1 end
444     if a=="?" then b=i:norm(b); a=b<.5 and 1 or 0
445     elseif b=="?" then a=i:norm(a); b=a<.5 and 1 or 0
446     else a,b = i:norm(a), i:norm(b) end
447     return math.abs(a - b) end
448
449 function SYM.dist(i,a,b) return a=="?" and b=="?" and 1 or a==b and 0 or 1 end
450

```

```

453 -----
454 --- DISCRETIZE
455 ---
456 --- $ lua brknbad.lua -t bins
457 ---
458 ---
459 ---
460 ---          selects    diversity
461 ---          =====
462 ---          Clndrs    < 5      211    0.48
463 ---          Clndrs    >= 5     187    0.30  <== best overall
464 ---
465 ---          Volume    < 121    158    0.23
466 ---          121    <= Volume    < 168    63    0.84
467 ---          168    <= Volume    < 225    32    0.20
468 ---          Volume    >= 225    145    0.00  <== pretty good
469 ---
470 ---          Model     < 73     125    0.87
471 ---          73    <= Model     < 76     91    0.97
472 ---          76    <= Model     < 79     93    1.00
473 ---          Model    >= 79      89    0.47
474 ---
475 ---          origin    == 1      249    0.72  <== pretty bad
476 ---          origin    == 2       70    0.00
477 ---          origin    == 3       79    0.00
478 ---
479 BIN=class"BIN"
480 function BIN:new(col,lo,hi,n,div)
481     return new({col=col, lo=lo, hi=hi, n=n, div=div},BIN) end
482
483 function BIN.selects(i,row, x)
484     x = row[i.col.at]
485     return x=="?" or i.lo==i.hi and x==i.lo or i.lo<=x and x<i.hi end
486
487 function BIN.show(i,negative)
488     local x, lo,hi,big, s = i.col.name, i.lo, i.hi, math.huge
489     if negative then
490         if lo== hi then s=fmt("%s=%s",x,lo)
491         elseif hi== big then s=fmt("%s< %s",x,lo)
492         elseif lo==big then s=fmt("%s>= %s",x,hi)
493         else
494             s=fmt("%s< %s and %s>= %s",x,lo,x,hi) end
495     else
496         if lo== hi then s=fmt("%s== %s",x,lo)
497         elseif hi== big then s=fmt("%s>= %s",x,lo)
498         elseif lo==big then s=fmt("%s< %s",x,hi)
499         else
500             s=fmt("%s<= %s< %s",lo,x,hi) end end
501     return s end
502
503 function BIN.distance2heaven(i, divs, ns)
504     return ((1 - ns:norm(i.n))^2 + (0 - divs:norm(i.div))^2)^0.5 end
505
506 function BIN:best(bins)
507     local divs,ns, distance2heaven = NUM(), NUM()
508     function distance2heaven(bin) return (bin:distance2heaven(divs,ns),bin) end
509     for _,bin in pairs(bins) do
510         divs:add(bin.div); ns:add( bin.n)
511     end
512     return sort(map(bins, distance2heaven), firsts)[1][2] end
513
514 function EGS.bins(i,j, bins)
515     bins = {}
516     for n,col in pairs(i.cols.x) do
517         for _,bin in pairs(col:bins(j.cols.x[n])) do push(bins, bin) end end
518     return bins end
519
520 --- DISCRETIZE SYMS
521 ---
522 function SYM.bins(i,j)
523     local xys= {}
524     for x,n in pairs(i.all) do push(xys, {x=x,y="left", n=n}) end
525     for x,n in pairs(j.all) do push(xys, {x=x,y="right",n=n}) end
526     return BIN:new4SYMs(i, SYM, xys) end
527
528 function BIN:new4SYMs(col, yclass, xys)
529     local out,all={}, {}
530     for _,xy in pairs(xys) do
531         all[xy.x] = all[xy.x] or yclass()
532         all[xy.x]:add(xy.y, xy.n) end
533     for x,one in pairs(all) do push(out,BIN(col, x, x, one.n, one:div())) end
534     return out end
535
536 --- DISCRETIZE NUMS
537 ---
538 function NUM.bins(i,j)
539     local xys, all = {}, NUM()
540     for _,n in pairs(i._all) do all:add(n); push(xys,{x=n,y="left"}) end
541     for _,n in pairs(j._all) do all:add(n); push(xys,{x=n,y="right"}) end
542     return BIN:new4NUMs(i, SYM, sort(xys,function(a,b) return a.x < b.x end),
543         (#xys)^the.minItems, all.sd*the.cohen) end
544
545 function BIN:new4NUMs(col, yclass, xys, minItems, cohen)
546     local out, b4, argmin = {}, -math.huge
547     function argmin(lo,hi)
548         local lhs, rhs, cut, div, xpect, xy = yclass(), yclass()
549         for j=lo,hi do rhs:add(xys[j].y) end
550         div = rhs:div()
551         if hi-lo+1 > 2*minItems
552             then
553                 for j=lo,hi - minItems do
554                     lhs:add(xys[j].y)
555                     rhs:sub(xys[j].y)
556                     if lhs.n > minItems and -- enough items (on left)
557                         xys[j].x ~= xys[j+1].x and -- there is a break here
558                         xys[j].x - xys[lo].x > cohen and -- not trivially small (on left)
559                         xys[hi].x - xys[j].x > cohen -- not trivially small (on right)
560                     then
561                         xpect = (lhs.n*lhs:div() + rhs.n*rhs:div()) / (lhs.n+rhs.n)
562                         if xpect < div then -- cutting here simplifies things
563                             cut, div = j, xpect end end end --end for
564                 end -- end if
565             if cut
566                 then argmin(lo, cut)
567                 argmin(cut+1, hi )
568             else b4 = push(out, BIN(col, b4, xys[hi].x, hi-lo+1, div)).hi end
569         end
570         argmin(1,#xys)
571         out[#out].hi = math.huge
572     return out end
573

```

```

573 -----
574 --- XPLAIN
575 ---
576 --- % lua brknbad.lua -r xplain
577 ---
578 ---
579 ---
580 ---          Weight-  Acc+  Mpg+
581 ---          =====
582 ---
583 ---          398
584 ---          Clndrs >= 5 : 190
585 ---          Model < 73 : 50
586 ---          Volume >= 318 : 29 (4213.93 11.52 12.41)
587 ---          Volume < 318 : 21 (3412.71 14.38 18.10)
588 ---          Model >= 73 : 140
589 ---          Model >= 78 : 50 (3354.20 15.68 22.40)
590 ---          Volume >= 225 : 32 (3554.53 15.69 20.94)
591 ---          Model < 78 : 90
592 ---          Volume < 262 : 43 (3298.33 16.97 20.00)
593 ---          Model >= 75 : 28 (3401.82 17.36 20.00)
594 ---          Volume >= 262 : 47
595 ---          Model < 74 : 20 (4279.05 12.25 12.00) <== worst
596 ---          Model >= 74 : 27 (4177.30 13.40 15.93)
597 ---          Clndrs < 5 : 208
598 ---          origin == 3 : 73
599 ---          Model >= 78 : 41 (2176.20 16.37 33.66)
600 ---          Model >= 80 : 31 (2176.10 16.36 34.84) <=== best
601 ---          Model < 78 : 32 (2155.03 16.41 26.87)
602 ---          origin != 3 : 135
603 ---          origin == 2 : 63
604 ---          Model >= 75 : 36 (2363.81 16.76 30.83)
605 ---          Model < 75 : 27 (2284.96 16.67 26.30)
606 ---          origin != 2 : 72
607 ---          Model < 78 : 28 (2319.25 17.11 26.07)
608 ---          Model >= 78 : 44 (2512.20 16.16 29.77)
609 ---          Model >= 80 : 31 (2547.77 16.51 30.00)
610
611 XPLAIN=class"XPLAIN"
612 function XPLAIN:new(top,egs)
613     local i,stop,lefts,rights,yes, no
614     egs = egs or top
615     i = new({egs=egs,top=top},XPLAIN)
616     stop = (#top._all)^the.minItems
617     if #egs._all > 2*stop then
618         lefts, rights= top:half(egs._all)
619         if #lefts._all < #egs._all then
620             i.bin = BIN:best( lefts:bins(rights) )
621             yes, no = top:copy(), top:copy()
622             for _,row in pairs(egs._all) do
623                 (i.bin:selects(row) and yes or no):add(row) end
624             if #yes._all > stop then i.yes = XPLAIN(top, yes) end
625             if #no._all > stop then i.no = XPLAIN(top, no) end end end
626     return i end
627
628 function XPLAIN.show(i, pre,how)
629     pre, how = pre or "", how or ""
630     local front = fmt("%s%s", pre, how, #i.egs._all)
631     if i.yes and i.no
632     then print(fmt("%-40s",front))
633     else print(fmt("%-40s",front, o(rnds(i.egs:mid()))))
634     end
635     if i.yes then i.yes:show(" .. pre, i.bin:show() ..") end
636     if i.no then i.no:show( " .. pre, i.bin:show(true) ..") end end
637

```

