

```

1  --
2  local b4={}; for k, _ in pairs(_ENV) do b4[k]=k end
3  local l,the,help = {}, {}, {}
4  gate: guess, assess, try, expand
5  (c) 2023, Tim Menzies, BSD-2
6  Learn a little, guess a lot, try the strangest guess, learn a little more, repeat
7  USAGE:
8  lua gate.lua [OPTIONS]
9
10 --
11 OPTIONS:
12 -c --cohen      small effect size      = .35
13 -f --file       csv data file name     = ../data/diabetes.csv
14 -h --help       show help              = false
15 -k --k          low class frequency kludge = 1
16 -m --m          low attribute frequency kludge = 2
17 -s --seed       random number seed     = 31210
18 -t --todo       start up action         = help]]
19
20 -----
21 -- ## Classes
22 local function isa(x,y) return setmetatable(y,x) end
23 local function is(s, t) t={a=s}; t.__index=t; return t end
24
25 -- ## Columns
26 -- ## Numerics
27
28 -- Create
29 local NUM=is"NUM"
30 function NUM.new(s, n)
31   return isa(NUM, {txt=s or "", at=n or 0, n=0, mu=0, m2=0, hi=-1E30, lo=1E30,
32     heaven = (s or ""):find"%-$" and 0 or 1}) end
33
34 -- Update
35 function NUM:add(x, d)
36   if x ~= "" then
37     self.n = self.n+1
38     d = x - self.mu
39     self.mu = self.mu + d/self.n
40     self.m2 = self.m2 + d*(x - self.mu)
41     self.lo = math.min(x, self.lo)
42     self.hi = math.max(x, self.hi) end end
43
44 -- Query
45 function NUM:mid() return self.mu end
46
47 function NUM:div() return self.n < 2 and 0 or (self.m2/(self.n - 1))^0.5 end
48
49 function NUM:small() return the.cohen*self:div() end
50
51 function NUM:norm(x)
52   return x=="*" and x or x - self.lo / (self.hi - self.lo + 1E-30) end
53
54 -- Likelihood
55 function NUM:like(x, nom, denom)
56   local mu, sd = self:mid(), (self:div() + 1E-30)
57   nom = 2.718^(-.5*(x - mu)^2/(sd^2))
58   denom = (sd*2.5 + 1E-30)
59   return nom/denom end
60
61 -- ## Symbols
62
63 -- Create
64 local SYM=is"SYM"
65 function SYM.new(s,n)
66   return isa(SYM, {txt=s or "", at=n or 0, n=0, has={}, mode=nil, most=0}) end
67
68 -- Update
69 function SYM:add(x)
70   if x ~= "" then
71     self.n = self.n + 1
72     self.has[x] = 1 + (self.has[x] or 0)
73     if self.has[x] > self.most then
74       self.most,self.mode = self.has[x], x end end end
75
76 -- Query
77 function SYM:mid() return self.mode end
78
79 function SYM:div( e)
80   e=0; for _,v in pairs(self.has) do e=e+v/self.n*math.log(v/self.n,2) end; return e end
81
82 function SYM:small() return 0 end
83
84 -- Likelihood
85 function SYM:like(x, prior)
86   return ((self.has[x] or 0) + the.m*prior)/(self.n + the.m) end
87
88 -- ## Columns
89 -- A contrainer storing multiple 'NUM's and 'SYM's.
90
91 -- Create a set of columns from a set of strings. If uppercase
92 -- then 'NUM', else 'SYM'. 'Klass'es end in '!'. Numeric goals to
93 -- minimize of maximize end in '-', '+'. Keep all cols in 'all'.
94 -- Also add dependent columns to 'y' (anything ending in '-', '+', '!') and
95 -- independent columns in 'x' (skipping over anything ending in 'X').
96 local COLS=is"COLS"
97 function COLS.new(row)
98   local x,y,all = {}, {}, {}
99   local klass,col
100   for at,txt in pairs(row.cells) do
101     col = (txt:find"%[A-Z]" and NUM or SYM).new(txt,at)
102     all[l+fall] = col
103     if not txt:find"%S" then klass=col end
104     if txt:find"%S" then klass=col end
105     (txt:find"%[a-jS" and y or x)[at] = col end end
106   return isa(COLS, {x=x, y=y, all=all, klass=klass, names=row.cells}) end
107
108 -- Update
109 function COLS:add(row)
110   for _,cols in pairs(self.x, self.y) do
111     for _,col in pairs(cols) do
112       col:add(row.cells[col.at]) end end
113   return row end
114
115

```

```

116 -- ## ROW
117
118 -- Store cells.
119 local ROW=is"ROW"
120 function ROW:new(t) return isa(ROW, { cells = t }) end
121
122 -- Distance to best values (and _lower_ is _better_).
123 function ROW:d2h(data, d, n)
124   d, n = 0, 0
125   for _, col in pairs(data.cols.y) do
126     n = n + 1
127     d = d + math.abs(col.heaven - col:norm(self.cells[col.at])) ^ 2 end
128   return d ^ .5 / n ^ .5 end
129
130 -- Return the 'data' (from 'datas') that I like the best
131 function ROW:likes(datas, n,nHypotheses,most,tmp,out)
132   n,nHypotheses = 0,0
133   for k,data in pairs(datas) do
134     n = n + #data.rows
135     nHypotheses = 1 + nHypotheses end
136   for k,data in pairs(datas) do
137     tmp = self:like(data,n,nHypotheses)
138     if most==nil or tmp > most then most,tmp,out = tmp,k end end
139   return out,most end
140
141 -- How much does ROW like 'self'. Using logs since these
142 -- numbers are going to get very small.
143 function ROW:like(data,n,nHypotheses, prior,out,v,inc)
144   prior = (#data.rows + the.k) / (n + the.k * nHypotheses)
145   out = math.log(prior)
146   for _,col in pairs(data.cols.x) do
147     v = self.cells[col.at]
148     if v ~= "" then
149       inc = col:like(v,prior)
150       out = out + math.log(inc) end end
151   return math.exp(l)*out end
152
153 -- ## Data
154 -- Store 'rows', summarized in 'COL'umns.
155
156 local DATA=is"DATA"
157 function DATA.new(src, fun, self)
158   self = isa(DATA, {rows={}, cols=nil})
159   if type(src) == "string"
160     then for _,x in ipairs(src) do self:add(x, fun) end
161     else for _,x in pairs(src or {}) do self:add(x, fun) end end
162   return self end
163
164 -- Update. First time through, assume the row defines the columns.
165 -- Otherwise, update the columns then store the rows. If 'fun' is
166 -- defined, call it before updating anything.
167 function DATA:add(t, fun,row)
168   row = t.cells and t or ROW:new(t)
169   if self.cols
170     then if fun then fun(self,row) end
171     self.rows[l + #self.rows] = self.cols:add(row)
172     else self.cols = COLS.new(row) end end
173
174 -- Query
175 function DATA:mid(cols, u)
176   u = {}; for _, col in pairs(cols or self.cols.all) do u[l + #u] = col:mid() end
177   return ROW:new(u) end
178
179 function DATA:div(cols, u)
180   u = {}; for _, col in pairs(cols or self.cols.all) do u[l + #u] = col:div() end;
181   return ROW:new(u) end
182
183 function DATA:small( u)
184   u = {}; for _, col in pairs(self.cols.all) do u[l + #u] = col:small(); end
185   return ROW:new(u) end
186
187 function DATA:stats(cols,fun,ndivs, u)
188   u = {{"N"} = #self.rows}
189   for _,col in pairs(self.cols[cols or "y"]) do
190     v[col:txt] = 1.rnd(getmetatable(col)[fun or "mid"] (col), ndivs) end
191   return u end
192
193

```

```

194 -- Gate.
195 function DATA:gate(budget0,budget,some)
196   local rows,lite,dark
197   local stats,bests = {}, {}
198   rows = 1.shuffle(self.rows)
199   lite = 1.slice(rows,1,budget0)
200   dark = 1.slice(rows, budget0+1)
201   for i=1,budget do
202     local best, rest = self:bestRest(lite, (#lite)*some) -- assess
203     local todo, selected = self:split(best,rest,lite,dark)
204     stats[i] = selected:mid()
205     bests[i] = best.rows[1]
206     table.insert(lite, table.remove(dark,todo)) end
207   return stats,bests end
208
209 -- Find the row scoring based on our acquite function.
210 function DATA:split(best,rest,lite,dark)
211   local selected,max,out
212   selected = DATA.new(self.cols.names)
213   max = 1E30
214   out = 1
215   for i,row in pairs(dark) do
216     local b,r,tmp
217     b = row:like(best, #lite, 2)
218     r = row:like(rest, #lite, 2)
219     if b>r then selected:add(row) end
220     tmp = math.abs(b+r) / math.abs(b-r+1E-300)
221     --print(b,r,tmp)
222     if tmp > max then out,max = i,tmp end end
223   return out,selected end
224
225 -- Sort on distance to heaven, split off the first 'want' items to return
226 -- a 'best' and 'rest' data.
227 function DATA:bestRest(rows, want, rest, top)
228   table.sort(rows, function(a, b) return a:d2h(self) < b:d2h(self) end)
229   best, rest = { self.cols.names }, { self.cols.names }
230   for i, row in pairs(rows) do
231     if i <= want then best[l + #best] = row else rest[l + #rest] = row end
232   end
233   return DATA.new(best), DATA.new(rest)
234
235

```

```

236 -- ## Library Functions
237
238
239 -- ### Linting
240 function l.rogues()
241   for k,v in pairs(_ENV) do if not b4[k] then print("E",k,type(k)) end end end
242
243 -- ### Numbers
244 function l.rnd(n, ndecs)
245   if type(n) ~= "number" then return n end
246   if math.floor(n) == n then return n end
247   local mult = 10^(ndecs or 2)
248   return math.floor(n * mult + 0.5) / mult end
249
250 -- ### Lists
251
252 -- Sorted keys
253 function l.keys(t, u)
254   u = {}
255   for k,_ in pairs(t) do u[l+#u]=k end; table.sort(u); return u end
256
257 -- Deep copy
258 function l.copy(t, u)
259   if type(t) ~= "table" then return t end
260   u = {}
261   for k,v in pairs(t) do u[l.copy(k)] = l.copy(v) end
262   return u end
263
264 -- Return a new table, with old items sorted randomly.
265 function l.shuffle(t, u, j)
266   u = {}
267   for _,x in pairs(t) do u[l+#u]=x; end;
268   for i = #u,2,-1 do j=math.random(i); u[i],u[j] = u[j],u[i] end
269   return u end
270
271 -- Return 't' skipping 'go' to 'stop' in steps of 'inc'.
272 function l.slice(t, go, stop, inc, u)
273   if go and go < 0 then go=#t+go end
274   if stop and stop < 0 then stop=#t+stop end
275   u = {}
276   for j=(go or 1)//1,(stop or #t)//1,(inc or 1)//1 do u[l+#u]=t[j] end
277   return u end
278
279 -- ### String to Things
280
281 -- Coerce string to intm float, nil, true, false, or (it all else fails), a strong.
282 function l.coerce(s1, fun)
283   function fun(s2)
284     if s2=="nil" then return nil else return s2=="true" or (s2=="false" and s2) end end
285     return math.tointeger(s1) or tonumber(s1) or fun(s1:match'^(%s+)(.*)$') end
286   end
287
288 -- Parse help string to infer the settings.
289 function l.settings(s, t,pat)
290   t,pat = {}, "[|]|([%S]+)=[|]|([%S]+)"
291   for k, s1 in s:match(pat) do t[k] = l.coerce(s1) end
292   t._help = s
293   return t end
294
295 -- Return a list of comma separated values (coerced to things)
296 function l.cells(s, t)
297   t = {}
298   for sl in s:match'^(%[|]+)$' do t[l+#t]=l.coerce(s1) end; return t end
299
300 -- Return rows of a csv file.
301 function l.csv(src, _l)
302   i,src = 0,src==" " and io.stdin or io.input(src)
303   return function( s )
304     s=io.read()
305     if s then i=i+1; return i,l.cells(s) else io.close(src) end end end
306
307 -- Update a table of settings using command-line settings.
308 function l.cli(t)
309   for k, v in pairs(t) do
310     v = tostring(v)
311     for argv,i in pairs(arg) do
312       if s==" " .. (k:sub(1,1)) or s=="-" .. k then
313         v = v=="true" and "false" or v=="false" and "true" or arg[argv + 1]
314         t[k] = l.coerce(v) end end
315     if t._help then os.exit(print("u"..t._help)) end
316     return t end
317
318 -- ### Things to Strings
319
320 -- Emulate sprintf
321 l.fmt = string.format
322
323 -- Print a string of a nested structure.
324 function l.o(x) print(l.o(x)); return x end
325
326 -- Rerun a string for a nested structure.
327 function l.o(t, _n, u)
328   if type(t) == "number" then return tostring(l.rnd(t, n)) end
329   if type(t) == "table" then return tostring(t) end
330   u = {}
331   for _,k in pairs(l.keys(t)) do
332     if tostring(k:sub(1,1)) == "." then
333       u[l+#u] = #t>0 and l.o(t[k],n) or l.fmt("%s: %s", l.o(k,n), l.o(t[k],n)) end end
334   return "[" .. table.concat(u, ", ") .. "]" end
335

```

```

336 -- ## Examples
337
338 -- ### Examples support code
339
340 -- Where to store examples
341 local eg={}
342
343 local function run(k, oops,b4)
344   b4 = l.copy(the) -- set up
345   math.randomseed(the.seed) -- set up
346   oops = eg[k] or {}
347   io.stder:write(l.fmt("#%s%s\n",oops and "aM~JM~L FAIL" or "aM~JM~E PASS",k))
348   for k,v in pairs(b4) do the[k]=v end -- tear down
349   return oops end
350
351 -- Run all examples
352 function eg.all( bad)
353   bad=0
354   for _,k in pairs(l.keys(eg)) do
355     if k ~= "all" then
356       if run(k) then bad=bad+1 end end end
357   io.stder:write(l.fmt("#%s%s fail(s)\n",bad>0 and "aM~JM~L FAIL" or "aM~JM~E PASS",bad))
358   os.exit(bad) end
359
360 -- List all example names
361 function eg.egs()
362   for _,k in pairs(l.keys(eg)) do print(l.fmt("lua gate.lua -t%s",k)) end end
363
364 -- ### The actual examples
365 function eg.o()
366   return l.o(a=1,b=2,c=3,d={e=3,f=4}) == "[a:1,b:2,c:3,d:[e:3,f:4]]" end
367
368 function eg.the() l.o(the); return the.help ~= nil and the.seed and the.m and the.k end
369
370 function eg.help() print("m"..the._help) end
371
372 function eg.sym( s,mode,e)
373   e = SYM.new()
374   for _, x in pairs(1,1,1,1,2,2,3) do s:add(x) end
375   mode, e = s:mid(), s:div()
376   print(mode, e)
377   return 1.37 < e and e < 1.38 and mode == 1 end
378
379 local function norm(mu,sd, R)
380   b=math.random
381   return (mu or 0) + (sd or 1) * math.sqrt(-2 * math.log(R()))
382   * math.cos(2 * math.pi * R()) end
383
384 function eg.num( e,mu,sd)
385   e = NUM.new()
386   for _ = 1,1000 do e:add(norm(10, 2)) end
387   mu, sd = e:mid(), e:div()
388   print(l.rnd(mu,3), l.rnd(sd,3))
389   return 10 < mu and mu < 10.1 and 2 < sd and sd < 2.05 end
390
391 function eg.csv( n)
392   n=0
393   for i,t in l.csv(the.file) do
394     if i%100 == 0 then n = n + #t; print(i, l.o(t)) end end
395   return n == 63 end
396
397 function eg.data( d,n)
398   n=0
399   d = DATA.new(the.file)
400   for i, row in pairs(d.rows) do
401     if i % 100 == 0 then n = n + #row.cells; l.o(row.cells) end end
402   l.o(d.cols.x[1].cells)
403   return n == 63 end
404
405 local function learn(data,row, my,kl)
406   my.n = my.n + 1
407   kl = row.cells[data.cols.klass.at]
408   if my.n > 10 then
409     my.tries = my.tries + 1
410     my.acc = my.acc (kl == row.likes(my.datas) and 1 or 0) end
411   my.datas[kl] = my.datas[kl] or DATA.new(data.cols.names)
412   my.datas[kl]:add(row) end
413
414 function eg.bayes()
415   local wme = {acc=0,datas={},tries=0,n=0}
416   DATA.new("../data/diabetes.csv", function(data,t) learn(data,t,wme) end)
417   print(wme.acc/(wme.tries))
418   return wme.acc/(wme.tries) > .72 end
419
420 function eg.km()
421   print(l.fmt("#%4s%4s%4s", "acc", "k", "m"))
422   for k=0,3,1 do
423     for m=0,3,1 do
424       the.k = k
425       the.m = m
426       local wme = {acc=0,datas={},tries=0,n=0}
427       DATA.new("../data/soybean.csv", function(data,t) learn(data,t,wme) end)
428       print(l.fmt("#%5.2N%4s%4s", wme.acc/wme.tries, k,m)) end end end
429
430 function eg.stats()
431   return l.o(DATA.new("../data/auto93.csv"):stats()) ==
432     "[N:398, Acc+:15.57,Lbs--:2970.42, Mpg+:23.84]" end
433
434 function eg.sorted( d)
435   d=DATA.new("../data/auto93.csv")
436   table.sort(d.rows, function(a,b) return a:d2h(d) < b:d2h(d) end)
437   print("M",l.o(d.cols.names))
438   for i, row in pairs(d.rows) do
439     if i < 5 or i > #d.rows - 5 then print(i, l.o(row.cells)) end end end
440
441 function eg.gate(stats, bests, d, say,sayd)
442   local budget0,budget,some = 4,10,.5
443   print(the.seed)
444   d = DATA.new("../data/auto93.csv")
445   function sayd(row, txt) print(l.o(row.cells), txt, l.rnd(row:d2h(d))) end
446   function say( row,txt) print(l.o(row.cells), txt) end
447   print(l.o(d.cols.names), "about",d2h*)
448   print("Overall" ..
449     sayd(d:mid(), "mid")
450     say(d:div(), "div")
451     say(d:small(), "small-div"..the.cohen)
452     print("generality" ..
453       stats,bests = d:gate(budget0, budget, some)
454       for i,stat in pairs(stats) do sayd(stat,i+budget0) end
455       print("specifically" ..
456         for i,best in pairs(bests) do sayd(best,i+budget0) end
457         print("optimum" ..
458         table.sort(d.rows, function(a,b) return a:d2h(d) < b:d2h(d) end)
459         sayd(d.rows[1], #d.rows)
460         print("Random" ..
461         local rows=l.shuffle(d.rows)
462         rows = l.slice(rows,1,math.log(.05)/math.log(1-the.cohen/6))
463         table.sort(rows, function(a,b) return a:d2h(d) < b:d2h(d) end)
464         sayd(rows[1]) end
465   end
466
467 function eg.gate20( ss,bs,rs,d,stats,bests,rows,stat,best)

```

```

468   print("#average,#optimistic,#random")
469   ss,bs,rs=NUM.new(),NUM.new(),NUM.new()
470   for i=1,20 do
471     io.write(i," "); io.flush()
472     d=DATA.new(the.file)
473     d.rows = l.shuffle(d.rows)
474     stats,bests = d:gate(4, 16, .5)
475     ss:add(stats[#stats]:d2h(d))
476     bs:add(bests[#bests]:d2h(d))
477     stat,best = stats[#stats]:d2h(d), bests[#bests]:d2h(d)
478     rows=l.shuffle(d.rows)
479     rows = l.slice(rows,1,math.log(.05)/math.log(1-the.cohen/6))
480     table.sort(rows, function(a,b) return a:d2h(d) < b:d2h(d) end)
481     rs:add(rows[1]:d2h(d)) end
482   print(l.rnd(ss:mid(),2), l.rnd(bs:mid(),2), l.rnd(rs:mid(),2))
483   print(l.rnd(2*ss:div(),2), l.rnd(2*bs:div(),2), l.rnd(2*rs:div(),2))
484   end
485
486 -- ## Start-up
487
488 the = l.settings(help)
489 if not pcall(debug.getlocal,4,1) then run(l.cli(the).todo) end
490 l.rogues()
491 return {the=the, COLS=COLS, DATA=DATA, NUM=NUM, ROW=ROW, SYM=SYM}

```