

```

1  #!/usr/bin/env lua
2  -- vim : filetype=lua ts=2 sw=2 et :
3  --
4  --
5  --
6  --
7  --
8  --
9  --
10 --
11 local THE, help= {}, [[tussel [OPTIONS]
12
13 OPTIONS:
14 -Debug      on error, dump stack and exit : false
15 -dull       F small effect= stdev*dull      : .35
16 -Far        F where to find far things     : .9
17 -file       S read data from file : ../../data/auto93.csv
18 -goal       S smile,frown,xplor,doubt      : smile
19 -h          show help                      : false
20 -p          I distance coefficient          : 2
21 -Rest       F size of rest set is Rest*best : 4
22 -round      I round floats to "round" places : 2
23 -seed       I random number seed          : 10019
24 -Small      F splits at #t^small          : .5
25 -todo       S start-up action              : pass
26             -todo ALL = run all
27             -todo LS = list all
28 -verbose    show details                   : false
29 ]]
30 local function update_from_command_line(flag,x) -- maybe flipping defaults for booleans
31   for n,txt in ipairs(arg) do
32     if flag:match("^[^%s+][^%s+]"..txt:sub(2).."%") -- allow abbreviations for flags
33       then x=x=="false" and"true" or x=="true" and"false" or arg[n+1] end end
34   return x end
35
36 local function read_settings_from_2_blanks_and_1_dash()
37   help:gsub("^[^%s+][^%s+]"..txt:sub(2).."%",function(flag,x) -- flag,x = word1,last word
38     x= update_from_command_line(flag,x)
39     if x=="false" then x=false elseif x=="true" then x=true else x=tonumber(x) or x end
40     THE[flag] = x end) end
41
42

```

```

42 -----
43 --
44 --
45 --
46 --
47 local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
48 local function rogues()
49   for k,v in pairs(_ENV) do if not b4[k] then print("Rogue?",k,type(v)) end end end
50
51 local function push(t,x) table.insert(t,x); return x end
52 local function firsts(a,b) return a[1] < b[1] end
53 local function sort(t,f) table.sort(t,f); return t end
54 local function map(t,f, u) u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
55 local function keys(t, u)
56   u={}
57   for k,_ in pairs(t) do if tostring(k):sub(1,1) ~= "_" then push(u,k) end end
58   return sort(u) end
59
60 local function copy(t,u)
61   u={}; for k,v in pairs(t) do u[k] = v end; return setmetatable(u, getmetatable(t)) end
62
63 local function csv(file, x,row)
64   function row(x, t)
65     for y in x:gsub("%s+", ""):gmatch("[^,]+") do push(t,tonumber(y) or y) end; return t end
66   file = io.input(file)
67   return function() x=io.read(); if x then return row(x, {}) else io.close(file) end end end
68
69 local function green(s) return "\027[32m"..s.." \027[0m" end
70 local function yellow(s) return "\027[33m"..s.." \027[0m" end
71
72 local function rnd(x,d, n) n=10^(d or THE.round); return math.floor(x*n+0.5) / n end
73 local function rnds(t,d)
74   return map(t,function(x) return type(x)=="number" and rnd(x,d) or x end) end
75
76 local fmt = string.format
77 local function say(...) if THE.verbose then print(fmt(...)) end end
78 local function o(t, u,key)
79   function key(k) return fmt(":%s %s", yellow(k), o(t[k])) end
80   if type(t) ~= "table" then return tostring(t) end
81   u = #t>0 and map(t,o) or map(keys(t),key)
82   return green((t._is or "").."{")..table.concat(u, " ")..green("}") end
83
84 local function rand(lo,hi)
85   THE.seed = (16807 * THE.seed) % 2147483647
86   return (lo or 0) + ((hi or 1) - (lo or 0)) * THE.seed / 2147483647 end
87
88 local function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
89 local function any(t) return t[randi(1,#t)] end
90 local function many(t,n, u) u={};for j=1,n do push(u,any(t)) end; return u end
91 local function shuffle(t, j)
92   for i=#t,2,-1 do j=randi(1,i); t[i],t[j]=t[j],t[i] end; return t end
93
94 local function xpect(a,b) return (a.n*a:spread() + b.n*b:spread())/(a.n + b.n) end
95
96 local _id=0
97 local function ako(x) return getmetatable(x) end
98 local function new(mt,x) _id=_id+1; x._id=_id; return setmetatable(x,mt) end
99 local function klass(s, klass)
100   klass = {_is=s, __tostring=o}
101   klass.__index = klass
102   return new({__call=function(_, ...) return klass.new(...) end},klass) end
103
104

```

```

104 -----
105 --
106 -- NUM
107 --
108 --
109 local NUM=class"NUM"
110 function NUM.new(n,s)
111   return new(NUM, {txt=s or"", at=n or 0,lo=math.huge, hi=-math.huge,
112     _has={},
113     n=0,mu=0,m2=0,w=(s or ""):find"-" and -1 or 1}) end
114
115 function NUM.mid(i)   return i.mu end
116 function NUM.spread(i) return i.n<2 and 0 or (i.m2/(i.n-1))^0.5 end
117
118 function NUM.add(i,x, d)
119   if x ~= "?" then
120     push(i._has,x)
121     i.n = i.n+1; d=x-i.mu; i.mu=i.mu+d/i.n; i.m2=i.m2+d*(x-i.mu)
122     i.hi= math.max(i.hi,x)
123     i.lo= math.min(i.lo,x) end
124   return x end
125
126 function NUM.norm(i,x)
127   return math.abs(i.lo - i.hi) < 1E-32 and 0 or (x - i.lo) / (i.hi - i.lo) end
128
129 function NUM.merge(i,j, k)
130   k=NUM(i.at, i.txt)
131   for _,x in pairs(j._has) do k:add(x) end
132   return k end
133
134 local _bins
135 function NUM.bins(i,j, x,y,s,xstats)
136   xys = {}
137   for _,x in pairs(i._has) do push(xys, {x=x, y="best"}) end
138   for _,x in pairs(j._has) do push(xys, {x=x, y="rest"}) end
139   return _bins(xys, xpect(i,j)^THE.dull, (#xys)^THE.Small, i, SYM) end
140
141 function _bins(xys,dull,small,col,yklass, bin,bins,merge,span,spans)
142   function merge(b4, j,tmp,maybe,now,after)
143     j, tmp = 0, {}
144     while j < #b4 do
145       j = j + 1
146       now, after = b4[j], b4[j+1]
147       if after then
148         maybe = now.has:merge(after.has)
149         if maybe:spread()*1.01 <= xpect(now.has, after.has) then
150           now = {col=col, lo=now.lo, hi= after.hi, has=maybe}
151           j = j + 1 end end
152       push(tmp,now) end
153       return #tmp==#b4 and b4 or merge(tmp) end
154
155   bin = {col=col, lo=xys[1].x, hi=xys[1].x, has=yklass()}
156   bins = {bin}
157   for j,xy in pairs(sort(xys, function(a,b) return a.x < b.x end)) do
158     if j < #xys - small and -- enough items remaining after split
159       xy.x == xys[j+1].x and -- next item is different (so can split here)
160       bin.has.n > small and -- bin has enough items
161       bin.hi - bin.lo > dull -- bin is not trivially small
162     then bin = push(bins, {col=col, lo=bin.hi, hi=xy.x, has=yklass()}) end
163     bin.hi = xy.x
164     bin.has:add(xy.y) end
165   bins[1].lo = -math.huge
166   bins[#bins].hi = math.huge
167   return merge(bins) end

```

```

168 --
169 -- SKIP
170 --
171 --
172 local SKIP=class"SKIP"
173 function SKIP.new(n,s) return new(SKIP, {txt=s or"", at=n or 0}) end
174 function SKIP.add(i,x) return x end
175 function SKIP.mid() return "?" end
176 function SKIP.bins(...) return {} end
177
178 --
179 -- SYM
180 --
181 --
182 local SYM=class"SYM"
183
184 function SYM.new(n,s)
185   return new(SYM, {n=0,has={},txt=s or"", at=n or 0,mode=nil,most=0}) end
186 function SYM.add(i,x,n)
187   if x ~= "?" then
188     n = n or 1
189     i.n = i.n+ n
190     i.has[x] = n+(i.has[x] or 0)
191     if i.has[x] > i.most then i.most = i.has[x], x end end
192   return x end
193
194 function SYM.mid(i)   return i.mode end
195 function SYM.spread(i, e)
196   e=0; for _,n in pairs(i.has) do e = e - n/i.n*math.log(n/i.n,2) end; return e end
197
198 function SYM.merge(i,j, k)
199   k = SYM(i.at,i.txt)
200   for x,n in pairs(i.has) do k:add(x,n) end
201   for x,n in pairs(j.has) do k:add(x,n) end
202   return k end
203
204 function SYM.bins(i,j, bins,t)
205   t,bins = {},{}
206   for x,n in pairs(i.has) do t[x] = t[x] or SYM(); t[x]:add("best",n) end
207   for x,n in pairs(j.has) do t[x] = t[x] or SYM(); t[x]:add("rest",n) end
208   for x,stats in pairs(t) do
209     push(bins, {col=1, lo=x,hi=x, has=stats}) end
210   return bins end
211
212 function SYM.score(i,goal,tmp)
213   local goals={}
214   function goals.smile(b,r) return r>b and 0 or b*b/(b+r +1E-31) end
215   function goals.frown(b,r) return b<r and 0 or r*r/(b+r +1E-31) end
216   function goals.xplor(b,r) return 1/(b+r +1E-31) end
217   function goals.doubt(b,r) return 1/(math.abs(b-r) +1E-31) end
218   local best, rest = 0, 0
219   for x,n in pairs(i.has) do
220     if x==goal then best = best+n/i.n else rest = rest+n/i.n end end
221   if best + rest < 0.01 then return 0
222   return goals[THE.goal](best,rest) end
223
224 --
225 -- EG
226 --
227 --
228 local EG=class"EG"
229 function EG.new(t) return new(EG, {klass=0,has=t}) end
230
231 function EG.cols(i,cols) return map(cols, function(x) return i.has[x.at] end) end
232 function EG.dist(i,j,smpl, a,b,d,n,inc,dist1)
233   function dist1(num,a,b)
234     if num
235       then if a=="?" then b=num:norm(b); a=b>.5 and 0 or 1
236       elseif b=="?" then a=num:norm(a); b=a>.5 and 0 or 1
237       else a,b = num:norm(a), num:norm(b) end
238       return math.abs(a-b)
239     else return a==b and 0 or 1 end end
240
241   d,n = 0,1E-31
242   for col,_ in pairs(smpl.xs) do
243     n = n+1
244     a,b = i.has[col], j.has[col]
245     inc = a=="?" and b=="?" and 1 or dist1(smpl.num[col],a,b)
246     d = d + inc^THE.p end
247   return (d/n)^(1/THE.p) end
248
249 function EG.better(eg1,eg2,smpl, e,n,a,b,s1,s2)
250   s1,s2,e,n = 0,0,10,#smpl.y
251   for _,col in pairs(smpl.y) do
252     a = col:norm(eg1.has[col.at])
253     b = col:norm(eg2.has[col.at])
254     s1 = s1 - e^(col.w * (a-b)/n)
255     s2 = s2 - e^(col.w * (b-a)/n) end
256   return s1/n < s2/n end

```

```

257 --
258 -- Sample
259 --
260 --
261 local SAMPLE=class"SAMPLE"
262 function SAMPLE.new(inits, i)
263   l= new(SAMPLE, {head=nil, eggs={}, all={}, num={}, sym={}, xs={}, ys={}})
264   if type(inits)=="table" then for _,eg in pairs(inits) do i:add(eg) end end
265   if type(inits)=="string" then for eg in csv(inits) do i:add(eg) end end
266   return i end
267
268 function SAMPLE.skip(i, x) return x:find"." end
269 function SAMPLE.nump(i, x) return x:find"[A-Z]" end
270 function SAMPLE.goalp(i, x) return x:find"-" or x:find"+" end
271
272 function SAMPLE.add(i, eg, now)
273   eg = eg.has and eg.has or eg
274   if not i.head then
275     i.head = eg
276     for n,s in pairs(eg) do
277       now = (i:skip(s) and SKIP or i:nump(s) and NUM or SYM)(n,s)
278       push(i.all, now)
279       if not i:skip(s) then
280         push(i.goalp(s) and i.ys or i.xs, now) end end
281   else
282     push(i.egs, EG(eg))
283     for n,one in pairs(i.all) do one:add(eg[one.at]) end end
284   return i end
285
286 function SAMPLE.clone(i, inits, j)
287   j= SAMPLE()
288   j:add(copy(i.head))
289   for x in pairs(inits or {}) do j:add(x) end
290   return j end
291
292 function SAMPLE.stats(i, cols)
293   return map(cols or i.all, function(x) return x:mid() end) end
294
295 function SAMPLE.far(i, eg1, eggs, gap, tmp)
296   gap = function(eg2) return {eg2, eg1:dist(eg2,i)} end
297   tmp = sort(map(egs, gap), function(a,b) return a[2] < b[2] end)
298   return table.unpack(tmp[#tmp*THE.Far//1] ) end
299
300 local evals=0
301 function SAMPLE.split(i, eggs, here)
302   local a,b,c,there,best,rest,tmp
303   eggs = eggs or i.egs
304   evals = evals + (here and 1 or 2)
305   here = here or i:far(any(egs),egs)
306   there,c = i:far(here, eggs)
307   tmp = {}
308   for _,eg in pairs(egs) do
309     a = eg:dist(here, i)
310     b = eg:dist(there,i)
311     push(tmp, ((a^2 + c^2 - b^2) / (2*c), eg)) end
312   best,rest = {},{}
313   for n,eg in pairs(sort(tmp, firsts)) do
314     push(n <= .5*#egs and best or rest, eg[2]) end
315   if there:better(here,i) then rest,best = best,rest end
316   return i:clone(best), i:clone(rest),there end
317
318 function SAMPLE.twain(i,min,lv1,here, there)
319   lv1 = lv1 or 0
320   min = min or 2*(#i.egs)^THE.Small
321   if #i.egs < min then return i end
322   local best,rest,there = i:split(i.egs,here)
323   local bins = {}
324   for n,bestx in pairs(best.xs) do
325     for _,bin in pairs(bestx:bins(rest.xs[n])) do push(bins, bin) end end
326   local score = function(a,b) return a.has:score("best") > b.has:score("best") end
327   local bin = sort(bins, score)[1]
328   print(fmt("%s %s %s %s=(%s,%s)", o(rnds(i:stats(i.ys),0 )),
329     string.rep(".",lv1),
330     #i.egs, bin.col.txt, bin.lo, bin.hi ))
331   local left, right = i:clone(), i:clone()
332   for _,eg in pairs(i.egs) do
333     local x = eg.has[ bin.col.at ]
334     if x=="-" then left:add(eg); right:add(eg)
335     elseif bin.lo<=x and x<bin.hi then left:add(eg)
336     else right:add(eg) end end
337   if #left.egs < #i.egs then left:twain(min, lv1+1, there) end
338   if #right.egs < #i.egs then right:twain(min, lv1+1, there) end
339   end

```

```

340 --
341 -- Main
342 --
343 --
344 local go, nogo = {},{} -- places to store tests
345 local fails = 0 -- counter for failure
346
347 local function azzert(test,msg) -- update failure count before calling the real assert
348   msg=msg or ""
349   if test then print(" PASS:".msg)
350   else fails=fails+1
351     print(" FAIL:".msg)
352     if THE.Debug then assert(test,msg) end end end
353
354 local function main()
355   read_settings_from_2_blanks_and_1_dash() -- set up system
356   if THE.h then print(help); os.exit() end -- maybe show help
357   go[THE.todo]() -- run something, maybe changing failure count
358   rogues() -- report any stray globals
359   os.exit(fails) end
360
361 function go.ALL() -- run all tests, resetting the system before each test
362   for _,k in pairs(keys(go)) do
363     if k:match"^[a-z]" then
364       read_settings_from_2_blanks_and_1_dash()
365       print("\n"..k)
366       go[k]() end end end
367
368 function go.LS() -- list all tests
369   for _,k in pairs(keys(go)) do
370     if k:match"^[a-z]" then print(" -t"..k) end end end
371
372 --
373 -- Evals
374 --
375 --
376 function go.the(s) say(o(THE)) end -- to disable, change "go" to "nogo"
377 function nogo.fail(s) azzert(false, "can you handle failure?") end
378 function go.pass(s) azzert(true, "can you handle success?") end
379 function go.sample(s, eggs)
380   s=SAMPLE(THE.file)
381   azzert(398==#s.egs, "got enough rows?")
382   azzert(s.ys[1].w==1, "minimizing goals are -1?") end
383
384 function go.clone(s, t,s1,s2)
385   s=SAMPLE(THE.file)
386   s1=o(s.ys)
387   t=s:clone(s.egs)
388   s2=o(t.ys)
389   azzert(s1==s2, "cloning works?") end
390
391 function go.dominate(s, eggs)
392   s=SAMPLE(THE.file)
393   eggs = sort(s.egs, function(a,b) return a:better(b,s) end)
394   for i=1,5 do say(o(egs[i]:cols(s.ys))) end; say("")
395   for i=#egs-5,#egs do say(o(egs[i]:cols(s.ys))) end
396   azzert(egs[1]:better(egs[#egs],s), "y-sort working?") end
397
398 function go.distance( s, eg1,dist,tmp,j1,j2,d1,d2,one)
399   s=SAMPLE(THE.file)
400   eg1=s.egs[1]
401   dist = function(eg2) return {eg2, eg1:dist(eg2,s)} end
402   tmp = sort(map(s.egs, dist), function(a,b) return a[2] < b[2] end)
403   one = tmp[1][1]
404   for j=1,30 do
405     j1=randi(1,#tmp)
406     j2=randi(1,#tmp)
407     if j1>j2 then j1,j2=j2,j1 end
408     d1 = tmp[j1][1]:dist(one,s)
409     d2 = tmp[j2][1]:dist(one,s)
410     azzert(d1 <= d2, "distance?") end end
411
412 function go.num( m,n)
413   m=NUM()
414   for i=1,10 do m:add(i) end
415   n = copy(m)
416   for i=1,10 do n:add(i) end
417   azzert(2.95 == rnd(n:spread()), "sd ok?") end
418
419 -- bring stats back
420 function go.label( s,x)
421   s = SAMPLE(THE.file)
422   x= s:twain()
423   print("evals",evals)
424   end
425   -- cuts={}
426   -- for n,i in pairs(best.xs) do
427   --   j=rests.xs[n]
428   --   for _,cut in pairs(i:bins(j)) do push(cuts,cut) end end
429   --   for _,cut in pairs(sort(cuts,function(a,b)
430   --     return a.has:score("best") > b.has:score("best") end)) do
431   --     print(rnd(cut.has:score("best")), cut.col.txt, cut.lo, cut.hi) end end
432
433 main()

```