

```

1  #!/usr/bin/env lua
2  -- vim : filetype=lua ts=2 sw=2 et :
3  local THE, help= {}, [[de OPTIONS
4
5  -Debug      on error, dump stack and exit : false
6  -dull       F small effect= stdev*dull      : .35
7  -Far        F where to find far things      : .9
8  -file       S read data from file : ../data/autoc93.csv
9  -goal       S smile, frown, xplor, doubt    : smile
10 -h          show help                      : false
11 -p          I distance coefficient           : 2
12 -Rest       F size of rest set is Rest*best : 4
13 -round      I round floats to "round" places : 2
14 -seed       I random number seed           : 10019
15 -Small      F splits at #t^small            : .5
16 -todo       S start-up action               : pass
17             -todo ALL = run all
18             -todo LS  = list all
19 -verbose    show details                   : false
20 ]]
21 local function update_from_command_line(flag,x) -- maybe flipping defaults for booleans
22 for n,t in ipairs(arg) do
23   if flag:match("^"..txt:sub(2).."*") -- allow abbreviations for flags
24   then x=x=="false" and true or x=="true" and false or arg[n+1] end end
25 return x end
26
27 local function read_settings_from_2_blanks_and_1_dash()
28 help:gsub("^\\n\\n[^(%s)+][^\\n]*%s([^(%s)+]",function(flag,x) -- flag,x = word1,last word
29   x= update_from_command_line(flag,x)
30   if x=="false" then x=false elseif x=="true" then x=true else x=tonumber(x) or x end
31   THE[flag] = x end end
32
33

```

```

33 -----
34 local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
35 local function rogues()
36   for k,v in pairs(_ENV) do if not b4[k] then print("Rogue?",k,type(v)) end end end
37
38 local function push(t,x) table.insert(t,x); return x end
39 local function firsts(a,b) return a[1] < b[1] end
40 local function sort(t,f) table.sort(t,f); return t end
41 local function map(t,f, u) u={};for k,v in pairs(t) do push(u,f(v)) end; return u end
42 local function keys(t, u)
43   u={};
44   for k,_ in pairs(t) do if tostring(k):sub(1,1) ~= "_" then push(u,k) end end
45   return sort(u) end
46
47 local function copy(t,u)
48   u={}; for k,v in pairs(t) do u[k] = v end; return setmetatable(u, getmetatable(t)) end
49
50 local function csv(file, x,row)
51   function row(x, t)
52     for y in x:gsub("(%s+",""):gmatch("([^\r]+)" do push(t,tonumber(y) or y) end; return t end
53     file = io.input(file)
54     return function() x=io.read(); if x then return row(x, {}) else io.close(file) end end end
55
56 local function green(s) return "\027[32m"..s.."027[0m" end
57 local function yellow(s) return "\027[33m"..s.."027[0m" end
58
59 local function rnd(x,d, n) n=10^(d or THE.round); return math.floor(x*n+0.5) / n end
60 local function rnds(t,d)
61   return map(t,function(x) return type(x)=="number" and rnd(x,d) or x end) end
62
63 local fmt = string.format
64 local function say(...) if THE.verbose then print(fmt(...)) end end
65 local function o(t, u,key)
66   function key(k) return fmt(":%s%s", yellow(k), o(t[k])) end
67   if type(t) ~= "table" then return tostring(t) end
68   u = #t>0 and map(t,o) or map(keys(t),key)
69   return green((t._is or "").."{")..table.concat(u, " ")..green("}") end
70
71 local function rand(lo,hi)
72   THE.seed = (16807 * THE.seed) % 2147483647
73   return (lo or 0) + ((hi or 1) - (lo or 0)) * THE.seed / 2147483647 end
74
75 local function randi(lo,hi) return math.floor(0.5 + rand(lo,hi)) end
76 local function any(t) return t[randi(1,#t)] end
77 local function many(t,n, u) u={};for j=1,n do push(u,any(t)) end; return u end
78 local function shuffle(t, j)
79   for i=#t,2,-1 do j=randi(1,i); t[i],t[j]=t[j],t[i] end; return t end
80
81 local function xpect(a,b) return (a.n*a:spread() + b.n*b:spread())/(a.n + b.n) end
82
83 local _id=0
84 local function ako(x) return getmetatable(x) end
85 local function new(mt,x) _id=_id+1; x._id=_id; return setmetatable(x,mt) end
86 local function klass(s, klass)
87   klass = {__is=s, __tostring=o}
88   klass.__index = klass
89   return new({__call=function(_,...) return klass.new(...) end},klass) end
90
91

```

```

91 -----
92 local NUM=class"NUM"
93 function NUM.new(n,s)
94     return new(NUM, {txt=s or"", at=n or 0,lo=math.huge, hi=-math.huge,
95         _has={},
96         n=0,mu=0,m2=0,w=(s or ""):find"- " and -1 or 1}) end
97
98 function NUM.mid(i) return i.mu end
99 function NUM.spread(i) return i.n<2 and 0 or (i.m2/(i.n-1))^0.5 end
100
101 function NUM.add(i,x, d)
102     if x ~= "?" then
103         push(i._has,x)
104         i.n = i.n+1; d=x-i.mu; i.mu=i.mu+d/i.n; i.m2=i.m2+d*(x-i.mu)
105         i.hi= math.max(i.hi,x)
106         i.lo= math.min(i.lo,x) end
107     return x end
108
109 function NUM.norm(i,x)
110     return math.abs(i.lo - i.hi) < 1E-32 and 0 or (x - i.lo) / (i.hi - i.lo) end
111
112 function NUM.merge(i,j, k)
113     k=NUM(i.at, i.txt)
114     for _,x in pairs(j._has) do k:add(x) end
115     return k end
116
117 local _bins
118 function NUM.bins(i,j, x,y,xstats)
119     xys = {}
120     for _,x in pairs(i._has) do push(xys, {x=x, y="best"}) end
121     for _,x in pairs(j._has) do push(xys, {x=x, y="rest"}) end
122     return _bins(xys, xpect(i,j)*THE.dull, (#xys)^THE.Small, i, SYM) end
123
124 function _bins(xys,dull,small,col,yklass, bin,bins,merge,span,spans)
125     function merge(b4, j,tmp,maybe,now,after)
126         j, tmp = 0, {}
127         while j < #b4 do
128             j = j + 1
129             now, after = b4[j], b4[j+1]
130             if after then
131                 maybe = now.has:merge(after.has)
132                 if maybe:spread()*1.01 <= xpect(now.has, after.has) then
133                     now = {col=col, lo=now.lo, hi= after.hi, has=maybe}
134                     j = j + 1 end end
135                 push(tmp,now) end
136             return #tmp==#b4 and b4 or merge(tmp) end
137
138     bin = {col=col, lo=xys[1].x, hi=xys[1].x, has=yklass()}
139     bins = {bin}
140     for j,xy in pairs(sort(xys, function(a,b) return a.x < b.x end)) do
141         if j < #xys - small and -- enough items remaining after split
142             xy.x ~= xys[j+1].x and -- next item is different (so can split here)
143             bin.has.n > small and -- bin has enough items
144             bin.hi - bin.lo > dull -- bin is not trivially small
145             then bin = push(bins, {col=col, lo=bin.hi, hi=xy.x, has=yklass()}) end
146             bin.hi = xy.x
147             bin.has:add(xy.y) end
148     bins[1].lo = -math.huge
149     bins[#bins].hi = math.huge
150     return merge(bins) end
151
152

```

```

152 -----
153 local SKIP=class"SKIP"
154 function SKIP.new(n,s) return new(SKIP, {txt=s or"", at=n or 0}) end
155 function SKIP.add(i,x) return x end
156 function SKIP.mid() return "?" end
157 function SKIP.bins(...) return {} end
158
159 local SYM=class"SYM"
160
161 function SYM.new(n,s)
162     return new(SYM, {n=0,has={},txt=s or"", at=n or 0,mode=nil,most=0}) end
163 function SYM.add(i,x,n)
164     if x ~= "?" then
165         n = n or 1
166         i.n = i.n+ n
167         i.has[x] = n+(i.has[x] or 0)
168         if i.has[x] > i.most then i.most, i.mode = i.has[x], x end end
169     return x end
170
171 function SYM.mid(i) return i.mode end
172 function SYM.spread(i, e)
173     e=0; for _,n in pairs(i.has) do e = e - n/i.n*math.log(n/i.n,2) end; return e end
174
175 function SYM.merge(i,j, k)
176     k = SYM(i.at,i.txt)
177     for x,n in pairs(i.has) do k:add(x,n) end
178     for x,n in pairs(j.has) do k:add(x,n) end
179     return k end
180
181 function SYM.bins(i,j, bins,t)
182     t,bins = {},{}
183     for x,n in pairs(i.has) do t[x] = t[x] or SYM(); t[x]:add("best",n) end
184     for x,n in pairs(j.has) do t[x] = t[x] or SYM(); t[x]:add("rest",n) end
185     for x,stats in pairs(t) do
186         push(bins, {col=i, lo=x,hi=x, has=stats}) end
187     return bins end
188
189 function SYM.score(i,goal,tmp)
190     local goals={}
191     function goals.smile(b,r) return r>b and 0 or b*b/(b+r +1E-31) end
192     function goals.frown(b,r) return b<r and 0 or r*r/(b+r +1E-31) end
193     function goals.xplor(b,r) return 1/(b+r +1E-31) end
194     function goals.doubt(b,r) return 1/(math.abs(b-r) +1E-31) end
195     local best, rest = 0, 0
196     for x,n in pairs(i.has) do
197         if x==goal then best = best+n/i.n else rest = rest+n/i.n end end
198     if best + rest < 0.01 then return 0 end
199     return goals[THE.goal](best,rest) end
200
201
202 -----
203 local EG=class"EG"
204 function EG.new(t) return new(EG, {klass=0,has=t}) end
205
206 function EG.cols(i,cols) return map(cols, function(x) return i.has[x.at] end) end
207 function EG.dist(i,j,smpl, a,b,d,n,inc,dist1)
208     function dist1(num,a,b)
209         if num
210             then if a=="?" then b=num:norm(b); a=b>.5 and 0 or 1
211                 elseif b=="?" then a=num:norm(a); b=a>.5 and 0 or 1
212                 else a,b = num:norm(a), num:norm(b) end
213                 return math.abs(a-b)
214             else return a==b and 0 or 1 end end
215
216     d,n = 0,1E-31
217     for col, in pairs(smpl.xs) do
218         n = n+1
219         a,b = i.has[col], j.has[col]
220         inc = a=="?" and b=="?" and 1 or dist1(smpl.num[col],a,b)
221         d = d + inc^THE.p end
222     return (d/n)^(1/THE.p) end
223
224 function EG.better(eg1,eg2,smpl, e,n,a,b,s1,s2)
225     s1,s2,e,n = 0,0,10,#smpl.ys
226     for _,col in pairs(smpl.ys) do
227         a = col:norm(eg1.has[col.at])
228         b = col:norm(eg2.has[col.at])
229         s1 = s1 - e^(col.w * (a-b)/n)
230         s2 = s2 - e^(col.w * (b-a)/n) end
231     return s1/n < s2/n end
232
233

```

```

233 -----
234 local SAMPLE=class"SAMPLE"
235 function SAMPLE.new(inits, i)
236   i= new(SAMPLE, {head=nil,egs={},all={},num={},sym={},xs={},ys={})
237   if type(inits)=="table" then for _,eg in pairs(inits) do i:add(eg) end end
238   if type(inits)=="string" then for eg in csv(inits) do i:add(eg) end end
239   return i end
240
241 function SAMPLE.skip(i, x) return x:find"." end
242 function SAMPLE.nump(i, x) return x:find"^[A-Z]" end
243 function SAMPLE.goalp(i, x) return x:find"-" or x:find"+" end
244
245 function SAMPLE.add(i,eg, now)
246   eg = eg.has and eg.has or eg
247   if not i.head then
248     i.head = eg
249     for n,s in pairs(eg) do
250       now = (i:skip(s) and SKIP or i:nump(s) and NUM or SYM)(n,s)
251       push(i.all, now)
252       if not i:skip(s) then
253         push(i:goalp(s) and i.ys or i.xs, now) end end
254     else
255       push(i.egs, EG(eg))
256       for n,one in pairs(i.all) do one:add(eg[one.at]) end end
257     return i end
258
259 function SAMPLE.clone(i,inits, j)
260   j= SAMPLE()
261   j:add(copy(i.head))
262   for _,x in pairs(inits or {}) do j:add(x) end
263   return j end
264
265 function SAMPLE.stats(i, cols)
266   return map(cols or i.all, function(x) return x:mid() end) end
267
268 function SAMPLE.far(i,egl,egs, gap,tmp)
269   gap = function(eg2) return {eg2, egl:dist(eg2,i)} end
270   tmp = sort(map(egs, gap), function(a,b) return a[2] < b[2] end)
271   return table.unpack(tmp[#tmp*THE.Far//1] ) end
272
273 local evals=0
274 function SAMPLE.split(i,egs, here)
275   local a,b,c,there,best,rest,tmp
276   egs = egs or i.egs
277   evals = evals + (here and 1 or 2)
278   here = here or i:far(any(egs),egs)
279   there,c = i:far(here, egs)
280   tmp = {}
281   for _,eg in pairs(egs) do
282     a = eg:dist(here, i)
283     b = eg:dist(there,i)
284     push(tmp, {(a^2 + c^2 - b^2) / (2*c), eg}) end
285   best,rest = {},{}
286   for n,eg in pairs(sort(tmp, firsts)) do
287     push(n <= 5*#egs and best or rest, eg[2]) end
288   if there:better(here,i) then rest,best = best,rest end
289   return i:clone(best), i:clone(rest),there end
290
291 function SAMPLE.twain(i,min,lvl,here, there)
292   lvl = lvl or 0
293   min = min or 2*(#i.egs)^THE.Small
294   if #i.egs < min then return i end
295   local best,rest,there = i:split(i.egs,here)
296   local bins = {}
297   for n,bestx in pairs(best.xs) do
298     for _,bin in pairs(bestx:bins(rest.xs[n])) do push(bins, bin) end end
299   local score = function(a,b) return a.has:score("best") > b.has:score("best") end
300   local bin = sort(bins, score)[1]
301   print(fmt("%s %s %s %s=(%s,%s)", o(rnds(i:stats(i.ys),0)),
302     string.rep(".",lvl),
303     #i.egs, bin.col.txt, bin.lo, bin.hi ))
304   local left, right = i:clone(), i:clone()
305   for _,eg in pairs(i.egs) do
306     local x = eg.has[ bin.col.at ]
307     if x=="?" then left:add(eg); right:add(eg)
308     elseif bin.lo<=x and x<bin.hi then left:add(eg)
309     else right:add(eg) end end
310   if #left.egs < #i.egs then left:twain(min, lvl+1, there) end
311   if #right.egs < #i.egs then right:twain(min, lvl+1, there) end
312   end
313
314

```

```

314 -----
315 local go, nogo = {},{} -- places to store tests
316 local fails = 0 -- counter for failure
317
318 local function azzert(test,msg) -- update failure count before calling the real assert
319   msg=msg or ""
320   if test then print(" PASS:"..msg)
321   else fails=fails+1
322     print(" FAIL:"..msg)
323     if THE.Debug then assert(test,msg) end end end
324
325 local function main()
326   read_settings_from_2_blanks_and_1_dash() -- set up system
327   if THE.h then print(help) end -- maybe show help
328   go[THE.todo]() -- run something, maybe changing failure count
329   rogues() -- report any stray globals
330   os.exit(fails) end -- exit, reporting the failure counts
331
332 function go.ALL() -- run all tests, resetting the system before each test
333   for _,k in pairs(keys(go)) do
334     if k:match"^[a-z]" then
335       read_settings_from_2_blanks_and_1_dash()
336       print("\n"..k)
337       go[k]() end end end
338
339 function go.LS() -- list all tests
340   for _,k in pairs(keys(go)) do
341     if k:match"^[a-z]" then print(" -t"..k) end end end
342
343 function go.the(s) say(o(THE)) end -- to disable, change "go" to "nogo"
344 function nogo.fail(s) azzert(false,"can you handle failure?") end
345 function go.pass(s) azzert(true, "can you handle success?") end
346 function go.sample(s, egs)
347   s=SAMPLE(THE.file)
348   azzert(398==#s.egs, "got enough rows?")
349   azzert(s.ys[1].w==1, "minimizing goals are-1?") end
350
351 function go.clone(s, t,s1,s2)
352   s=SAMPLE(THE.file)
353   s1=o(s.ys)
354   t=s:clone(s.egs)
355   s2=o(t.ys)
356   azzert(s1==s2, "cloning works?") end
357
358 function go.dominate(s, egs)
359   s=SAMPLE(THE.file)
360   egs = sort(s.egs, function(a,b) return a:better(b,s) end)
361   for i=1,5 do say(o(egs[i]:cols(s.ys))) end; say("")
362   for i=#egs-5,#egs do say(o(egs[i]:cols(s.ys))) end
363   azzert(egs[1]:better(egs[#egs],s), "y-sort working?") end
364
365 function go.distance( s,egl,dist,tmp,j1,j2,d1,d2,one)
366   s=SAMPLE(THE.file)
367   egl=s.egs[1]
368   dist = function(eg2) return {eg2,egl:dist(eg2,s)} end
369   tmp = sort(map(s.egs, dist), function(a,b) return a[2] < b[2] end)
370   one = tmp[1][1]
371   for j=1,30 do
372     j1=randi(1,#tmp)
373     j2=randi(1,#tmp)
374     if j1>j2 then j1,j2=j2,j1 end
375     d1 = tmp[j1][1]:dist(one,s)
376     d2 = tmp[j2][1]:dist(one,s)
377     azzert(d1 <= d2,"distance?") end end
378
379 function go.num( m,n)
380   m=NUM()
381   for i=1,10 do m:add(i) end
382   n = copy(m)
383   for i=1,10 do n:add(i) end
384   azzert(2.95 == rnd(n:spread()),"sd ok?") end
385
386 -- bring stats back
387 function go.label( s,x)
388   s = SAMPLE(THE.file)
389   x= s:twain()
390   print("cvals",evals)
391   end
392   -- cuts={}
393   -- for n,i in pairs(best.xs) do
394   --   j=rests.xs[n]
395   --   for _,cut in pairs(i:bins(j)) do push(cuts,cut) end end
396   --   for _,cut in pairs(sort(cuts,function(a,b)
397   --     return a.has:score("best") > b.has:score("best") end)) do
398   --     print(rnd(cut.has:score("best")), cut.col.txt, cut.lo, cut.hi) end end
399
400 main()

```