

```

1 local help=[[
2 TINYXAI.lua: dont fuse on non-numerics
3
4 USAGE:
5   lua tinyxai.lua [OPTIONS]
6
7 OPTIONS:
8   -b bins          max number of bins      = 16
9   -c cohen          difference in nums      = .35
10  -f file           source                  = ../data/auto93.csv
11  -F Far            how far is far          = .95
12  -g go             action                  = help
13  -h help           show help               = false
14  -m min            size of small           = .5
15  -p                distance coefficient    = 2
16  -r rests          number of rests to use  = 4
17  -s Samples        samples                 = 64
18  -seed            random number seed      = 10019]]
19
20 ----- Names
21 ----- Locals
22 -- Store old names (so, on last line, we can check for rogue locals)
23 local b4={}; for k, _ in pairs(_ENV) do b4[k]=k end
24 -- Define new names from library code (so we can mention them before defining them).
25 local any,big,cat,chat,cli,coerce,csv,data,fmt,lt
26 local map,max,min,per,push,rand,shuffle,sort
27 -- Place to store test suites (to disable a test, move it 'go' to 'no').
28 local go,no={},{}
29 -- Place to store the settings (and this variable is parsed from help text; e.g. 'the.
30 bins=16').
31 local the = {}
32
33 ----- Objects
34 -- obj(str,fun): class -- 'Fun' is a constructor for instances of class 'str'.
35 -- BTW, polymorphism, encapsulation, classes, instance, constructors, all in 3 lines.
36 --[[
37 local function obj(txt,fun, t,i)
38   local function new(k,...) i=metatable({},k); fun(i,...); return i end
39   t[_tostring = cat]; t.__index = t;return setmetatable(t,{__call=new}) end
40 --]]
41
42 --- ROW(tab) -- Stores one record. ROWs are stored in a contained called ROWS.
43 --- Implementation note: ROWs are created when data is read from CSV files.
44 --- After this, if a table is added to more than one ROWs object then the same
45 --- ROW will be held in different ROWSs. This makes certain labelling and
46 --- record keep tasks easier (e.g. tracking how many rows we have evaluated).
47 local ROW=obj("ROW", function(self,cells)
48   self.cells = {} -- place to hold one record
49   self.label = false -- true if we have decided this ROW is "best"?
50   self.evaluated = false end -- have we accessed this row's y-values?
51
52 --- SYM(?num=0, ?str="") -- Summarizes streams of symbols in ROWs
53 local SYM=obj("SYM", function(self,at,txt)
54   self.n = 0 -- number of items seen
55   self.at = at or 0 -- column number
56   self.txt = txt or "" -- column name
57   self.kept= {} end -- counters for symbols
58
59 --- NUM(?num=0, ?str="") -- Summarize streams of numbers in ROWs
60 local NUM=obj("NUM", function(self,at,txt)
61   self.n = 0 -- number of items seen
62   self.at = at or 0 -- column number
63   self.txt = txt -- column name
64   self.w = txt:find"$" and -1 or 1 -- If minimizing, then -1. Else 1
65   self.kept= {} -- some sample of the seen items
66   self.ok = false end -- true if sorted, set to false by each add
67
68 --- COLS({str}+) -- Factory for making NUMs or SYMs from list of col names.
69 --- Column names starting with '+' are NUMs (others are SYMs).
70 --- Anything adding with '+-+' is a dependent goal column.
71 --- Column names ending with '!' are "skipped"; i.e.
72 --- not added to the list of independent or dependent columns.
73 local COLS=obj("COLS", function(self, names)
74   self.names=names -- list of column names
75   self.all = {} -- [NUM|SYM] all names, converted to NUMs or SYMs
76   self.x = {} -- [NUM|SYM] just the independent columns
77   self.y = {} -- [NUM|SYM] just the dependent columns
78   self.klass= nil -- SYM the klass column (if it exists)
79   for k,v in pairs(names) do
80     col= push(self.all, (v:find"^[A-Z]" and NUM or SYM) (at,txt))
81     if not v:find"$" then
82       if v:find"$" then self.klass = col end
83       push(v:find"[+-$]" and self.y or self.x, col) end end end
84
85 --- ROWS() -- Stores 'rows' and their summaries in 'cols'.
86 local ROWS=obj("ROWS", function(self)
87   self.rows = {} -- [ROW] records, stored as ROW
88   self.cols = nil end -- a COLS instance (if nil, no data read yet)
89
90 --- BIN|NUM|SYM, num, ?num=lo, ?SYM) -- Values from same rows in 2 columns
91 local BIN=obj("BIN", function(self,col, lo, hi, has)
92   self.col = col -- What column does this bin handle?
93   self.lo = lo -- Lowest value of column1
94   self.hi = hi or lo -- Highest value of column1
95   self.has = has or SYM() end -- Symbol counts of column2 values.
96
97 ----- Columns
98 ----- Sym
99 -- Create
100 -- SYM:merge(SYM): SYM -- Create a new SYM by merging two others.
101 function SYM:merge(other, k)
102   k= SYM(self.at, self.txt)
103   for x,n in pairs(self.kept) do k:add(x,n) end
104   for x,n in pairs(other.kept) do k:add(x,n) end
105   return k end
106
107 ----- Update
108 -- SYM:add(any,?num=1) -- Add a symbols 'x'. Do it 'n' times.
109 function SYM:add(x,n)
110   n = n or 1
111   if x=="?" then self.n=self.n+n; self.kept[x]=n + (self.kept[x]+0) end end
112   return mode end
113
114 ----- Query
115 -- SYM:div(i):num -- Diversity. Return entropy.
116 function SYM:div(i)
117   return sum(self.kept, function(n) return -n/i.n*math.log(n/i.n,2) end) end
118
119 -- SYM:mid():num -- Return 'mid'dle (mode) symbol.
120 function SYM:mid()
121   local most,mode = -1,nil
122   for x,n in pairs(self.kept) do if n>most then most,mode=n,x end end
123   return mode end
124
125 ----- Distance
126 -- SYM:dist(i,tm,atom):num -- Identical symbols have distance 0. Otherwise, 1.
127 -- If any unknowns, assume max distance.
128 function SYM:dist(x,y) return (x=="?" or y=="?" and 1 or x==y and 0 or 1 end
129
130 ----- Discretization
131 -- SYM:bin(any):any -- Discretize a symbol (do nothing)
132 function SYM:bin(x) return x end
133 -- SYM:merges(t,...):SYM -- Merge adjacent bins (do nothing: SYM ranges don't merge)
134 function SYM:merges(t,...) return t end
135
136 ----- Num
137 -- NUM:add(num) -- Add 'x'. If no space, at prob 'some/n', replace any old number.
138 function NUM:add(x)
139   if x=="?" then
140     self.n = self.n + 1
141     local pos
142     if #self.kept < the.some then pos= #i.kept+1
143     elseif rand() < the.some/self.n then pos= rand(#i.kept) end
144     if pos then
145       self.ok=false -- the 'kept' list is no longer in sorted order
146       self.kept[pos]=x end end end
147
148 ----- Query
149 -- NUM:has() -- Return 'kept', ensuring it is sorted.
150 function NUM:has()
151   self.kept = self.ok and self.kept or sort(self.kept)
152   self.ok = true
153   return self.kept end
154
155 -- NUM:mid():num -- Return 'mid'dle (median) number.
156 function NUM:mid() return per(self.has(),.5) end
157
158 -- NUM:norm(x):num -- Normalize x,y to 0..1.
159 function NUM:norm(x)
160   local a = self:has()
161   local lo,hi = a[1], a[#a]
162   return x=="?" and x or math.abs(hi-lo)<1E-9 and 0 or (x-lo)/(hi-lo+1/big) end
163
164 ----- Distance
165 -- NUM:dist(x,y):num -- Normalize x,y to 0..1, report their difference.
166 -- If any unknowns, assume max distance.
167 function NUM:dist(x,y)
168   if x=="?" and y=="?" then return 1
169   elseif x=="?" then y=self:norm(y); x=y*.5 and 1 or 0
170   elseif y=="?" then x=self:norm(x); y=x*.5 and 1 or 0
171   else x,y = self:norm(x), self:norm(y) end
172   return math.abs(x-y) end
173
174 ----- Discretization
175 -- NUM:bin(any):any -- Discretize a num to one of 'the.bins'.
176 function NUM:bin(x)
177   local s = self:has()
178   local lo,hi = a[1], a[#a]
179   local b = (hi - lo)/the.bins
180   return hi==lo and 1 or math.floor((x-b+.5)*b) end
181
182 --- NUM:merges([BIN],...) :[BIN] -- Prune superfluous bins.
183 function NUM:merges(b4, min)
184   local n,now = 1,{}
185   while n <= #b4 do
186     local merged = n<#b4 and b4[n]:merged(b4[n+1],min) -- defined in BIN
187     now[#now+1] = merged or b4[n] -- merged or b4[n]
188     n = n + (merged and 2 or 1) -- if merged, skip over merged bin
189   end -- end while
190   if #now < #b4 then return self:merges(now,min) end -- seek others to merge
191   bins[1].lo,bins[#bins].hi = -big,big -- grow to plus/minus infinity
192   return bins end
193
194 ----- Data
195 ----- COLS
196 ----- Update
197 -- COLS:add(ROW) -- update the non-skipped columns with values from ROW
198 function COLS:add(row)
199   for _,col in pairs(self.x, self.y) do
200     for _,col in pairs(cols) do col:add(row.cells[col.at]) end end end
201
202 ----- Distance
203 -- COLS:dist(ROW,ROW):num -- Using 'x' columns, compute distance.
204 function COLS:dist(r1,r2)
205   local d,x1,x2 = 0
206   for _,col in pairs(self.x) do
207     x1 = r1.cells[col.at]
208     x2 = r2.cells[col.at]
209     d = d+(col:dist(x1,x2))*the.p end
210   return (d/#self.x)^(1/the.p) end
211
212 -- COLS:half(ROW) -- Divide 'rows' by their distance to two distant points A,B
213 -- Find two distant points by 'the.samples' times, look at random pairs.
214 function COLS:half(rows, b4)
215   local function Abc(A,B) return (A=A, B=B, As={}, Bs={}, c=self:dist(A,B)) end
216   local ABCs={}
217   for n=1,the.Samples do push(ABCs, Abc(b4 or any(rows), -- if b4, use it for one pole
218     any(rows))) end
219   local i = per(sort(ABCs,lt"C"), the.Far) -- avoid outliers: don't go right to the ed
220   ge
221   local function xCs(C)
222     return {i=self:dist(C,i.A)^2+i.c^2-self:dist(C,i.B)^2/(2*i.c),
223       C = C} end
224   for j,xCol in pairs(sort(map(rows,xCs,lt"x"))) do
225     push(j<#rows/2 and i.As or i.Bs, xC.C) end
226     return i end
227
228 ----- Optimize
229 -- COLS:best(ROW,ROW):bool -- Multi-objective comparisons. True if moving to self 1
230 -- cases least than moving to other.
231 function COLS:best(r1,r2)
232   r1.evaluated,r2.evaluated = true,true
233   local s1, s2, ys, e = 0, 0, self.y, math.exp(1)
234   for _,col in pairs(ys) do
235     local x = col:norm(r1.cells[col.at])
236     local y = col:norm(r2.cells[col.at])
237     s1 = s1+e*(col.w * (x-y)/ys)
238     s2 = s2 - e*(col.w * (y-x)/ys) end
239   return s1/#ys < s2/#ys end
240
241 --- COLS:bests([ROW]):{bests:[ROW],rests:[ROW]} -- Recursively apply 'best' to findReturn
242 n most preferred rows, and the rest.
243 function COLS:bests(rows, b4,stop,rests)
244   rests = rests or {}
245   stop = stop or (#rows)^the.min
246   if #rows < stop then return rows,rests end -- return best=[ROW],rests=[ROW]
247   local two = self:half(rows,b4) -- if b4 supplied, then half will use it as one pole.
248   local bests, rests, restal
249   if self:best(two.A, two.B)
250   then best, rests, restal = two.A, two.As, two.Bs
251   for i=#restal,1,-1 do push(rests, restal[i]) end -- rests sorted, L to R, worst
252   to better
253   else best, rests, restal = two.B, two.Bs, two.As
254   for i=1,#restal, 1 do push(rests, restal[i]) end -- rests sorted, L to R, worst
255   to better
256   end
257   return self:best(bests, best, stop,rests) end
258
259 ----- COLS
260 ----- Create
261 -- BIN:merged(BIN, num) -- Combine two bins if we should or can do so.
262 function BIN:merged(i, min)
263   local a, b, c = self:has, j.has, self:has:merge(j.has)
264   local should = a.n < min or b.n < min -- "should" if either too small
265   local can = cidiv() <= (a.n*a:div() + b.n*b:div())/c.n -- "can" if whole simpler
266   than parts.
267   if should or can then return BIN(a.col,self.lo, j.hi, c) end end
268
269 ----- Update
270 -- BIN:add(num,sym) -- extend 'lo,hi' to cover 'x'; remember we saw 'y'.
271 function BIN:add(x,y)
272   self.lo = min(x,self.lo)
273   self.hi = max(x,self.hi)
274   self:has(y) end
275
276 ----- Query
277 -- BIN:show() -- pretty print the range
278 function BIN:show(i)
279   local x,lo,hi = self.y.s.txt, self.lo, self.hi
280   if lo == hi then return fmt("%s==%s", x, lo)
281   elseif hi == big then return fmt("%s> %s", x, lo)
282   elseif lo == -big then return fmt("%s<= %s", x, hi)
283   else return fmt("%s< %s<= %s", lo,x,hi) end end
284
285 --- BIN:selects([ROW]):[ROW] -- Returns the subset of rows that fall within this BIN.
286 -- Returns nil if the subset is same size as original sets.
287 function BIN:selects(rows, select,tmp)
288   v= row.cells[self.col.at]
289   if v=="?" or self.lo==self.hi or self.lo<v and v <self.hu then return row end end
290   tmp= map(rows,select)
291   if #tmp < #rows then return rows end end
292

```

```

288 -----
289 function ROWS:clone(t) return ROWS():add(self.cols.names):adds(t or {}) end
290
291 function ROWS:file(x) for t in csv(x) do self:add(t) end; return self end
292
293 function ROWS:adds(t) for _,t1 in pairs(t) do self:add(t1) end; return self end
294
295 function ROWS:add(t)
296   if self.cols
297   then self.cols:add(push(i.rows, t.cells and t or ROW(t)))
298   else self.cols=COLS(t) end end
299
300 -- ROWS:mids(?int=2,?[COL]=self.cols.y):[key=num] -- Return 'mid' of columns rounded t
301 o 'p' places
302 function ROWS:mids(p,cols)
303   local t=(n=#self.rows)
304   for _,col in pairs(cols or self.cols.y) do t[col.txt]=col:mid(p) end
305   return rnds(t,p or 2) end
306
307 function ROWS:splitter(rows)
308   function split(col)
309     for _,row in pairs(rows) do
310       local v = row.cells[col.at]
311       if v ~= "?" then
312         local pos = col:bin(v)
313         dict[pos] = dict[pos] or push(list, BIN(col,v))
314         dict[pos]:add(v,row.label) end end
315   list = col:merges(sort(list,lt"lo", n"the.min)
316   return (bins = list,
317         div = sum(list,function(z) return z.has:div() * z.ys.n/n end))
318 end -----
319 return sort(map(self.cols.x, split),lt"div")[1].bin end
320
321 function ROWS:tree()
322   local bests,rests= self.cols:bests(self.rows)
323   local rows={}
324   for _,best in pairs(bests) do push(rows, best).label=1 end
325   for i = 1,the.rests#bests do push(rows,rests[i]).label=0 end
326   self:grow(rows, (#rows)"the.min)
327   return self end
328
329 function ROWS:grow(rows,stop,when)
330   local function kid(bin)
331     local t = bin:selects(rows)
332     if t then return self:clone(t):grow(t,stop,bin) end end
333   self.when=when
334   if #rows < stop then return self end
335   self.kids = map(self:splitter(rows), kid) end
336

```

```

336 -----
337 -- ## Lib
338 big=math.huge
339 min=math.min
340 max=math.max
341 fmt=string.format
342 rand=math.random
343
344 function any(a) return a[rand(#a)] end
345 function per(t,p) p=p*#t//1; return t[math.max(1,math.min(#t,p))] end
346
347 function push(t,x) t[1+#t]=x; return x end
348 function map(t,f, u) u={};for _,x in pairs(t) do u[1+#u]=f(x)end;return u end
349 function sum(t,f, u) u=0;for _,x in pairs(t) do u=u+f(x) end;return u end
350
351 --- rnd(num, places:int):num -- Return 'x' rounded to some number of 'place' 6#9312;
352 . <
353 function rnd(x, places) -- 6#9312;
354   local mult = 10^(places or 2)
355   return math.floor(x * mult + 0.5) / mult end
356 --- rnds(t:num, places:?int=2):num -- Return items in 't' rounds to 'places'.
357 function rnds(t, places)
358   local u={};for k,x in pairs(t) do u[k]=rnd(x,places or 2)end;return u end
359
360 function sort(t,f) table.sort(t,f); return t end
361 function lt(x) return function(a,b) return a[x] < b[x] end end
362
363 function shuffle(t, j)
364   for i=#t,2,-1 do j=rand(i); t[i],t[j]=t[j],t[i] end; return t end
365
366 function coerce(x)
367   x = x:match("(%s*(-)%s*$")
368   if x=="true" then return true elseif x=="false" then return false
369   else return math.tointeger(x) or tonumber(x) or x end end
370
371 function cli(t)
372   for k,v in pairs(t) do
373     v = tostring(v)
374     for n,x in pairs(arg) do if x=="-"..(k:sub(1,1)) then
375       x = v=="false" and "true" or v=="true" and "false" or arg[n+1] end end
376     t[k] = coerce(v) end end
377
378 function chat(t) print(cat(t)) return t end
379 function cat(t, show,u)
380   function show(k,v) return #t==0 and ("%s%s"):format(k,v) or tostring(v) end
381   u={}; for k,v in pairs(t) do u[1+#u]=show(k,v) end
382   return (t._is or "").."["..table.concat(#t==0 and sort(u) or u," ").."]" end
383
384 function csv(file,fun)
385   function lines(file, fun)
386     local file = io.input(file)
387     while true do
388       local line = io.read()
389       if not line then return io.close(file) else fun(line) end end
390   end
391   function words(s,sep,fun, t)
392     fun = fun or same
393     t={};for x in s:gmatch(fmt("(%s%s)",sep)) do t[1+#t]=fun(x) end; return t
394   end
395   lines(file, function(line) fun(words(line, ",", coerce)) end) end
396

```

```

396 -----
397 -- ## Start-up
398 function go.the() chat(the) end
399
400 help:gsub("\n[-]%%$[%s]+([%S]+)(^n)+=([%S]+)",function(k,x) the[k]=coerce(x)end)
401
402 cli(the)
403 if the.help then os.exit(print(help:gsub("(%u[%u%d]+","%27[1:32m%1:27[0m)","")) end
404 math.randomseed(the.seed)
405 if go[the.go] then go[the.go]() end
406 for k,v in pairs(_ENV) do if not b4[k] then print("?*",k,type(v)) end end

```