

```

1 local b4={}; for k,v in pairs(_ENV) do b4[k]=k end
2 local the,help = {},{}
3
4 TODO: dont fuse on non-numerics
5
6 -c cohen difference in nums      = .35
7 -f file source                  = ../data/auto93.csv
8 -g go action                    = help
9 -m min size of small            = .5
10 -s seed random number seed     = 10019]]
11
12 local cat,chat,cli,csv,fmt,fuse, fume1,kap,lines,lt,map,new
13 local obj,order,push,roques,same,sort,thing,trim,words
14
15 function lt(x) return function(a,b) return a[x] < b[x] end end
16
17 function trim(x) return x:match("^%s*(-)%s*$" end
18
19 function thing(x)
20   if x=="true" then return true else if x=="false" then return false
21   else return math.tointeger(x) or tonumber(x) or x end end
22
23 function lines(file, fun)
24   local file = io.input(file)
25   while true do
26     local line = io.read()
27     if not line then return io.close(file) else fun(line) end end end
28
29 function words(s,sep,fun, t)
30   fun = fun or same
31   t={};for x in s:gmatch(fmt("([%s]+)",sep)) do t[#t+1]=fun(x) end; return t end
32
33 function csv(file,fun)
34   lines(file, function(line) fun(words(line, ",", thing)) end end
35
36 function cli(t)
37   for key,x in pairs(t) do
38     x = tostring(x)
39     for n,flag in ipairs(arg) do
40       if flag=="-".key:sub(1,1)
41       then x = x=="false" and "true" or x=="true" and "false" or arg[n+1] end end
42     t[key] = thing(x) end
43   return t end
44
45 function roques()
46   for k,v in pairs(_ENV) do if not b4[k] then print("'",k,type(v)) end end end
47
48 fmt=string.format
49 function same(x) return x end
50 function map(t,f, u) u={};;for _,x in pairs(t)do u[#u+1]=f(x) end;return u end
51 function kap(t,f, u) u={};;for k,x in pairs(t)do u[#u+1]=f(k,x)end;return u end
52 function sort(t,f) table.sort(t,f); return t end
53 function push(t,x) t[#t+1]=x; return x end
54
55 function chat(t) print(cat(t)); return t end
56 function cat(t)
57   local function pub(k,v) return (tostring(k)):sub(1,1)=="-" end
58   local function key(k,v) if pub(k) then return fmt("%s%s",k,v) end end
59   local u= {}+1 and map(t,f or tostring) or sort(kap(t,key))
60   return (t._is or "")..{"..table.concat(u,"").."}" end
61
62 local _id = 0
63 function new(kl,...)
64   local x
65   _id=_id+1; x=setmetatable({_id=_id,kl};kl.new(x,...); return x end
66
67 function obj(name)
68   local t = {_tostring=cat,_is=name}; t.__index=t
69   return setmetatable(t, {_call=new}) end
70
71
72
73 local function col(self, at,txt)
74   self.at = at or 0
75   self.txt = txt or ""
76   self.n = 0
77   self.all = {} end
78
79 local SYM = obj"SYM"
80 function SYM:new(at,txt); self.kept={}; col(self,at,txt) end
81
82 function SYM:add(x)
83   if x ~= "" then
84     self.n=self.n+1
85     self.kept[x] = 1+(self.kept[x] or 0) end end
86
87 function SYM:bins(of,rows, x)
88   for _,row in pairs(rows) do
89     x = row.raw[self.at]
90     if x ~= "" then
91       self.all[x] = self.all[x] or (at=self,lo=x,hi=x,has=of:clone()) end
92     self.all[x].has:add(row) end end
93
94 function SYM:ent( e)
95   local function z(p) return p*math.log(p,2) end
96   e=0;for _,n in pairs(self.kept)do if n>0 then e=e-z(n/self.n)end end;return e end
97
98 local NUM = obj"NUM"
99 function NUM:new(at,txt)
100   col(self,at,txt)
101   self.lo = math.huge; self.hi=self.lo
102   self.mu, self.m2, self.sd = 0,0,0
103   self._bins= {}
104   self.w = self.txt:find"$" and -1 or 1 end
105
106 function NUM:add(x)
107   if x ~= "" then
108     self.n = self.n+1
109     local d = x - self.mu
110     self.mu = self.mu + d/self.n
111     self.m2 = self.m2 + d*(x - self.mu)
112     self.sd = (self.n < 2 or self.m2<0) and 0 or ((self.m2)/(self.n-1)).5
113     if x > self.hi then self.hi = x end
114     if x < self.lo then self.lo = x end end end
115
116 function NUM:bins(of, rows)
117   local function lt(x,y)
118     return (x=="?" and -math.huge or x) < (y=="?" and -math.huge or y) end
119   local order = function(a,b) return lt(a.raw[self.at],b.raw[self.at]) end
120   local x = function(k) return rows[k].raw[self.at] end
121   local xis = function(k,x) rows[k].cooked[self.at] = x; end
122   local n,b4 = 1,nil
123   local min = self.n,the.min
124   local epsilon= self.sd*the.cohen
125   local rows = sort(rows,order)
126   local one = push(self.all, (at=self, lo=x(1), hi=x(1), has=of:clone()))
127   for ,row in pairs(rows) do
128     if x(j) ~= "" then
129       if #rows>j+min and #one.has.rows>min and x(j)-=x(j+1) and x(j)-one.lo>epsilon
130       then one = push(self.all, (at=self,lo=one.hi,hi=x(j),has=of:clone())) end
131       one.hi = x(j)
132       one.has:add(row) end end
133   end
134
135 local is={}
136 is.skip= function(x) return x:find"$" end -- what to ignore
137 is.klass= function(x) return x:find"$" end -- single goal
138 is.goal= function(x) return x:find"[!-]$" end -- dependent column
139 is.num= function(x) return x:find"[A-Z]" end -- NUM or SYM?
140
141 local COLS = obj"COLS"
142 function COLS:new(names)
143   self.names = names -- {[str] list of known columns names
144   self.all = {} -- {[NUM|SYM] all the columns
145   self.x = {} -- {[NUM|SYM] list of pointers to just the independent column
146   s
147   self.y = {} -- {[NUM|SYM] list of pointers to just the dependent columns
148   self.klass = nil -- :?[NUM|SYM] pointer to the klass column, may be nil.
149   for at,txt in pairs(names) do
150     local col = (is.num(txt) and NUM or SYM)(at,txt)
151     push(self.all, col)
152     if not is.skip(txt) then
153       push(is.goal(txt) and self.y or self.x, col)
154       if is.klass(txt) then self.klass = col end end end end
155
156 function COLS:add(row)
157   for _,cols in pairs(self.x,self.y) do
158     for _,col in pairs(cols) do col:add(row.raw[col.at]) end end
159   return row end
160
161 local ROW = obj"ROW"
162 function ROW:new(of,cells)
163   self.raw = cells
164   self.cooked= cells
165   self._of = of
166   self.evald = false end
167
168 local ROWS = obj"ROWS"
169 function ROWS:new() self.rows={}; self.cols=nil end
170
171 function ROWS:clone(src)
172   return ROWS():add(self.cols.names):adds(src) end
173
174 function ROWS:adds(src)
175   if type(src) == "string"
176   then csv(src, function(row) self:add(row) end)
177   else for _,row in pairs(src or {}) do self:add(row) end end
178   return self end
179
180 function ROWS:add(row)
181   if self.cols
182   then push(self.rows, self.cols:add( row.raw and row or ROW(self,row)))
183   else self.cols = COLS(row) end
184   return self end
185
186 local go={}
187 function go.all()
188   local want = function(k,_if k=="all" then return k end end
189   for _,x in pairs(sort(kap(go,want))) do
190     math.randomseed(the.seed)
191     go[x]() end end
192
193 function go.help() print(help) end
194 function go.the() chat(the) end
195 function go.csv() csv(the.file, function(x) chat(x) end) end
196 function go.clone(t,s)
197   print(the.file)
198   r=ROWS():adds(the.file)
199   s=rc:clone()
200   print(s) end
201   --chat(s.cols.all[1]) end
202
203 local function ents(rows)
204   local e = 0
205   for _,col in pairs(rows.cols.x) do
206     local s = SYM()
207     for _,row in pairs(rows.rows) do
208       local x = row.cooked[col.at]
209       if x ~= "" then s:add(x) end
210     end
211     e = e + s:ent() end
212   return e end
213
214 local function fusel(i,j,epsilon)
215   if math.abs(i.e - j.e) <= epsilon then
216     return (at=i.at, lo=i.lo, hi=j.hi, n=i.n+j.n,
217     e= (i.n*i.e + j.n*j.e)/(i.n + j.n)) end end
218
219 function fuse(epsilon, b4)
220   local n,now = 1,{}
221   while n<=#b4 do
222     local merged = n#b4 and fusel(b4[n],b4[n+1],epsilon)
223     now[#now+1] = merged or b4[n]
224     n = n + (merged and 2 or 1)
225   end
226   if #now<#b4 then return fuse(epsilon,now) end
227   now[1].lo = -math.huge
228   now[#now].hi = math.huge
229   return now end
230
231 function go.rows(r, es,all)
232   es = NUM()
233   all = {}
234   r=ROWS():adds(the.file)
235   print(ents(r))
236   es:adds(ents(r))
237   for _,col in pairs(r.cols.x) do col:bins(r,r.rows) end
238   for _,col in pairs(r.cols.x) do
239     print(col.txt, #col.all)
240     local here={}
241     all[col.at]=here
242     for _,range in pairs(col.all) do
243       local e = ents(range.has)
244       es:add(e)
245       print("[".. range.lo, range.hi, #range.has.rows,e)
246       push(here, (at=range.at,lo=range.lo, hi=range.hi, n=#range.has.rows,e=e)) end e
247   nd
248   for x,one in pairs(all) do
249     all[x] = fuse(es.sd * the.cohen, sort(one,lt"lo")) end
250     print(".....")
251     for x,one in pairs(all) do
252       print(one[1].at.txt)
253       for _,range in pairs(one) do
254         print("[".. range.lo, range.hi, range.n, range.e) end end end
255
256 -----ddp-----
257 help:gs:sub("n[~]%S[%s]k[%S]o")"n"=+(%S)k",
258   function(k,x) the[k]=thing(trim(x)) end)
259 the=cli(the)
260 go(the.go)()
261 roques()
262 -----ddp-----
263 -- 12.348722672624
264 -- Clndrs
265 -- 3 4 207 10.186135496684
266 -- 3 6 87 8.2438268409432
267 -- 6 8 104 7.0495719622724
268 -- Volume
269 -- 68 105 106 9.1032053779412
270 -- 105 141 88 9.1314349037882
271 -- 141 173 31 8.1942729502156
272 -- 173 225 27 5.831712755027
273 -- 225 262 46 5.7889030186077
274 -- 262 305 23 5.294944385008
275 -- 305 350 39 4.7986838713161
276 -- 350 455 38 5.4934036234341
277 -- Model
278 -- 70 72 84 9.265662209105
279 -- 72 74 67 8.8229185085156
280 -- 74 76 64 8.6789383050277
281 -- 76 78 64 9.0203804206288
282 -- 78 80 58 9.159537586695
283 -- 80 82 61 7.8732012623659
284 -- origin
285 -- 1 1 249 10.186516126735
286 -- 2 2 70 8.5363112369701
287 -- 3 3 79 8.4061141489268
288 -- .....
289 -- Clndrs
290 -- -inf 4 207 10.186135496684
291 -- 4 6 87 8.2438268409432
292 -- 6 inf 104 7.0495719622724
293 -- Volume
294 -- -inf 141 194 9.1160105236862
295 -- 141 173 31 8.1942729502156
296 -- 173 inf 173 5.441775380802
297 -- Model
298 -- -inf 80 337 9.001367139192
299 -- 80 inf 61 7.8732012623659
300 -- origin
301 -- -inf 1 249 10.186516126735
302 -- 2 inf 149 8.4672805661283

```