

```

1 local help=[
2
3 XPLAN.lua: semi-supervised multi-objective explanation
4 (c)2022, Tim Menzies ttimm@ieee.org, BSD-2 license
5
6 SYNOPSIS:
7 Data, with multiple dependent goals, is recursive
8 bi-clustered on independent variables by partitioning
9 on the distance to two distant points (found after a
10 few dozen random projections). A decision tree reports
11 the difference between the "best" and "worst" clusters
12 (defined using a multi-objective domination predicate).
13 This process only makes log2(N) queries to y-values
14 (while clustering, just on the pairs of distance objects).
15
16 USAGE:
17 lua xplan.lua [OPTIONS]
18
19 OPTIONS:
20 -b bins          number of bins          = 16
21 -c cohen         difference in nums       = .35
22 -f file         source                   = ./../data/auto93.csv
23 -F Far          how far is far          = .95
24 -g go           action                  = nothing
25 -h help         show help              = false
26 -m min         size of small           = 5
27 -p p           distance coefficient     = 2
28 -r rests       number of rests to use  = 4
29 -p Projections number of random projections = 64
30 -s seed        random number seed      = 10019 ]]
31
32 ---- ---- ---- Names
33 ---- ---- Locals
34 -- Store old names (so, on last line, we can check for rogue locals)
35 local b4={}; for k, _ in pairs(_ENV) do b4[k]=k end
36 -- Define new names from library code (so we can mention them before defining them).
37 local any,big,cat,chi,coerce,chat,data,fft,kap,lt
38 local map,max,min,on,per,push,rand,rnd,ends,shuffle,sort,sum
39 -- Place for test suites (to disable a test, move it 'go' to 'no').
40 local go,no = {},{}
41 -- Place for settings (parsed from help text; e.g. 'the.bins=16').
42 local the = {}
43
44 ---- ---- Objects
45 ---- obj(str,fun): class -- 'Fun' is a constructor for instances of class 'str'.
46 -- Polymorphism, encapsulation, classes, instance, constructors: all in 3 lines. :-)
47 local function obj(txt,fun, t,t1)
48 local function new(k,...) i=metatable(t,{}); fun(i,...); return i end
49 t.__tostring = cat; t.__index = t:return setmetatable(t, __call=new) end
50
51 ---- ROW(tab) -- Stores one record. ROWs are stored in a contained called ROWS.
52 -- Implementation note: ROWs are created when data is read from CSV files.
53 -- After that, if a ROW is added to more than one ROWS object then the same
54 -- ROW will be held in different ROWS. This makes certain labelling and
55 -- record keep tasks easier (e.g. tracking how many rows we have evaluated).
56 local ROW=obj("ROW", function(self,cells)
57 self.cells = {} -- place to hold one record
58 self.label = false -- true if we have decided this ROW is "best"?
59 self.evald = false end) -- have we accessed this row's y-values?
60
61 ---- SYM(?num=0, ?str="") -- Summarize streams of symbols in ROWs
62 local SYM=obj("SYM", function(self,at,txt)
63 self.n = 0 -- number of items seen
64 self.at = at or 0 -- column number
65 self.txt = txt or "" -- column name
66 self.kept = {} end) -- counters for symbols
67
68 ---- NUM(?num=0, ?str="") -- Summarize streams of numbers in ROWs
69 local NUM=obj("NUM", function(self,at,txt)
70 self.n = 0 -- number of items seen
71 self.at = at or 0 -- column number
72 txt=txt or "" -- column name
73 self.txt = txt -- If minimizing, then -1. Else 1
74 self.v = txt:find("$" and -1 or 1) -- some sample of the seen items
75 self.kept = {} -- true if sorted, set to false by each add
76 self.ok = false end)
77
78 ---- COLS([str]+) -- Factory for making NUMs or SYMs from list of col names.
79 -- Column names starting with upper case are NUMs (others are SYMs).
80 -- Anything adding with "!"+" is a dependent goal column.
81 -- Column names ending with ".*" are "skipped"; i.e.
82 -- not added to the list of independent or dependent columns.
83 local COLS=obj("COLS", function(self, names)
84 self.names=names -- list of column names
85 self.all = {} -- [NUM|SYM] all names, converted to NUMs or SYMs
86 self.x = {} -- [NUM|SYM] just the independent columns
87 self.y = {} -- [NUM|SYM] just the dependent columns
88 self.klass= nil -- SYM the class column (if it exists)
89 for k,v in pairs(names) do
90 col= push(self.all, (v:find("[A-Z]" and NUM or SYM) (at,txt)))
91 if not v:find("$" then
92 if v:find("$" then self.klass = col end
93 push(v:find("[!-]*" and self.y or self.x, col) end end end)
94
95 ---- ROWS() -- Stores 'rows' and their summaries in 'cols'.
96 local ROWS=obj("ROWS", function(self)
97 self.rows = {} -- [ROW] records, stored as ROW
98 self.cols = nil end) -- a COLS instance (if nil, no data read yet)
99
100 ---- BIN( NUM[SYM, num, ?num=lo, ?SYM) -- Values from same rows in 2 columns
101 local BIN=obj("BIN",function(self,col, lo, hi, has)
102 self.col = col -- What column does this bin handle?
103 self.lo = lo -- Lowest value of column1.
104 self.hi = hi or lo -- Highest value of column1
105 self.has = has or SYM() end) -- Symbol counts of column2 values.
106
107
108 ---- ---- ---- Columns
109 ---- ---- ---- Sym
110 ---- ---- Create
111 ---- SYM:merge(SYM): SYM -- Create a new SYM by merging two others.
112 function SYM:merge(other, k)
113 k= SYM(self.at, self.txt)
114 for x,n in pairs(self.kept) do k:add(x,n) end
115 for x,n in pairs(other.kept) do k:add(x,n) end
116 return k end
117
118 ---- ---- Update
119 ---- SYM:add(any,?num=1) -- Add a symbols 'x'. Do it 'n' times.
120 function SYM:add(x,n)
121 n = n or 1
122 if x=="?" then self.n=self.n+n; self.kept[x]=n + (self.kept[x] or 0) end end
123
124 ---- ---- Query
125 ---- SYM:div(i):num -- Diversity. Return entropy.
126 function SYM:div(i)
127 return sum(self.kept, function(n) return -n/i.n*math.log(n/i.n,2) end) end
128
129 ---- SYM:mid():num -- Return 'mid'dle (mode) symbol.
130 function SYM:mid()
131 local most,mode = -1,nil
132 for x,n in pairs(self.kept) do if n>most then most,mode=n,x end end
133 return mode end
134
135 ---- ---- Distance
136 ---- SYM:dist(stom,atom):num -- Identical symbols have distance 0. Otherwise, 1.
137 -- If any unknowns, assume max distance.
138 function SYM:dist(x,y) return (x=="?" or y=="?" and 1 or x==y and 0 or 1 end)
139
140 ---- ---- Discretization
141 ---- SYM:bin(any):any -- Discretize a symbol (do nothing)
142 function SYM:bin(x) return x end
143 ---- SYM:merges(t,...):SYM -- Merge adjacent bins (do nothing: SYMs don't merge)
144 function SYM:merges(t,...) return t end
145
146 ---- ---- ---- Num
147 ---- ---- Update
148 ---- NUM:add(num) -- Add 'x'. If no space, at prob 'some/n', replace any old number.
149 function NUM:add(x)
150 if x=="?" then
151 self.n = self.n + 1
152 local pos
153 if #self.kept < the.some then pos= #i.kept+1
154 elseif rand() < the.some/self.n then pos= rand(#i.kept) end
155 if pos then
156 self.ok=false -- the 'kept' list is no longer in sorted order
157 self.kept[pos]=x end end end
158
159 ---- ---- Query
160 ---- NUM:has() -- Return 'kept', ensuring it is sorted.
161 function NUM:has()
162 self.kept = self.ok and self.kept or sort(self.kept)
163 self.ok = true
164 return self.kept end
165
166 ---- NUM:mid():num -- Return 'mid'dle (median) number.
167 function NUM:mid() return per(self.has(),.5) end
168
169 ---- NUM:norm(x):num -- Normalize x,y to 0..1.
170 function NUM:norm(x)
171 local a = self:has()
172 local lo,hi = a[1], a[#a]
173 return x=="?" and x or math.abs(hi-lo)<1E-9 and 0 or (x-lo)/(hi-lo+1/big) end
174
175 ---- ---- Distance
176 ---- NUM:dist(x,y):num -- Normalize x,y to 0..1, report their difference.
177 -- If any unknowns, assume max distance.
178 function NUM:dist(x,y)
179 if x=="?" and y=="?" then return 1
180 elseif x=="?" then y=self:norm(y); x=y<.5 and 1 or 0
181 elseif y=="?" then x=self:norm(x); y=x<.5 and 1 or 0
182 else x,y = self:norm(x), self:norm(y) end
183 return math.abs(x-y) end
184
185 ---- ---- Discretization
186 ---- NUM:bin(any):any -- Discretize a num to one of 'the.bins'.
187 function NUM:bin(x)
188 local a = self:has()
189 local lo,hi = a[1], a[#a]
190 local b = (hi - lo)/the.bins
191 return hi==lo and 1 or math.floor((x+b+.5)*b) end
192
193 ---- NUM:merges([BIN],...):[BIN] -- Prune superfluous bins.
194 function NUM:merges(b4, min)
195 local n,now = 1,{}
196 while n <= #b4 do
197 local merged = n<#b4 and b4[n]:merged(b4[n+1],min) -- defined in BIN
198 now[#now+1] = merged or b4[n]
199 n = n + (merged and 2 or 1) -- if merged, skip over merged bin
200 end -- end while
201 if #now < #b4 then return self:merges(now,min) end -- seek others to merge
202 bins[1].lo,bins[#bins].hi = -big,big -- grow to plus/minus infinity
203 return bins end
204
205
206 ---- ---- ---- Data
207 ---- ---- ---- COLS
208 ---- ---- Update
209 ---- COLS:add(ROW) -- update the non-skipped columns with values from ROW
210 function COLS:add(row)
211 for _,cols in pairs(self.x, self.y) do
212 for _,col in pairs(cols) do col:add(row.cells[col.at]) end end end
213
214 ---- ---- Distance
215 ---- COLS:dist(ROW,ROW):num -- Using 'x' columns, compute distance.
216 function COLS:dist(r1,r2)
217 local d,x1,x2 = 0
218 for _,col in pairs(self.x) do
219 x1 = r1.cells[col.at]
220 x2 = r2.cells[col.at]
221 d = d+(col:dist(x1,x2))*the.p end
222 return (d/#self.x)^(1/the.p) end
223
224 ---- COLS:half([ROW]) -- Divide 'rows' by their distance to two distant points A,B
225 -- Find two distant points A,B using a few dozen random projections.
226 function COLS:half(rows, b4)
227 local function ABC(A,B) return (A=A, B=B, As={}, Bs={}, c=self:dist(A,B)) end
228 local ABCs={}
229 for n=1,the.Projections do
230 push(ABCs, ABC(b4 or any(rows), the.Far) -- avoid outliers: only go so Far
231 local i = per(sort(ABCs,lt"), the.Far)
232 local function xCs(C)
233 return {x= (self:dist(C,i,A)^2+i.c^2-self:dist(C,i,B)^2)/(2*i.c),
234 C = C) end
235 for j,xC in pairs(sort(map(rows,xCs,lt"x"))) do
236 push(j<#rows/2 and i.As or i.Bs, xC.C) end
237 return i end
238
239 ---- ---- Optimize
240 ---- COLS:best(ROW,ROW):bool -- True if better on multi-objectives
241 function COLS:best(r1,r2)
242 r1.evald,r2.evald = true,true
243 local s1, s2, ys, e = 0, 0, self.y, math.exp(1)
244 for _,col in pairs(ys) do
245 local x = col:norm(r1.cells[col.at])
246 local y = col:norm(r2.cells[col.at])
247 s1 = s1 - e^(col.w * (x-y)/#ys)
248 s2 = s2 - e^(col.w * (y-x)/#ys) end
249 return s1/#ys < s2/#ys end -- i.e. we lose less going to r2->r1 than r2->r1
250
251 ---- COLS:bests([ROW]):{bests=[ROW],rests=[ROW]} -- Recursively apply 'best'.
252 function COLS:bests(rows, b4,stop,rests)
253 rests = stop or {}
254 stop = stop or (#rows)>the.min
255 if #rows < stop then return rows,rests end -- return best=[ROW],rests=[ROW]
256 local two = self:half(rows,b4) -- if b4 supplied, then half will use it as one pole.
257 local best, bests, restal
258 if self:best(two.A, two.B)
259 then best, bests, restal = two.A, two.As, two.Bs
260 for i=restal,1,-1 do push(rests,restal[i]) end --sort L to R, worst to better
261 else best, bests, restal = two.B, two.Bs, two.As
262 for i=1,restal,1 do push(rests,restal[i]) end --sort L to R, worst to better
263 return self:best(bests, best, stop,rests) end
264
265 ---- ---- ---- COLS
266 ---- ---- ---- BINS
267 ---- ---- Create
268 ---- BIN:merged(BIN, num) -- Combine two bins if we should or can do so.
269 -- "Should" is true if either is too small.
270 -- "Can" is true of the whole is simpler than the parts.
271 function BIN:merged(j, min)
272 local a, b, c = self:has, j:has, self:has:merge(j:has)
273 local should = a.n < min or b.n < min
274 local can = c:div() <= (a.n*a:div() + b.n*b:div())/c.n
275 if should or can then return BIN(a.col,self.lo, j.hi, c) end end
276
277 ---- ---- Update
278 ---- BIN:add(num,sym) -- extend 'lo,hi' to cover 'x'; remember we saw 'y'.
279 function BIN:add(x,y)
280 self.lo = min(x,self.lo)
281 self.hi = max(x,self.hi)
282 self:has(y) end
283
284 ---- ---- Query
285 ---- BIN:show() -- pretty print the range
286 function BIN:show(i)
287 local x,lo,hi = self:ys.txt, self.lo, self.hi
288 if lo == hi then return fmt("%%<=%s", x, lo)
289 elseif hi == big then return fmt("%%>=%s", x, lo)
290 elseif lo == -big then return fmt("%%<=%s", x, hi)
291 else return fmt("%%<=%s<=%s", lo,x,hi) end end
292
293 ---- BIN:selects([ROW]):[ROW] -- Returns rows that fall within this BIN.
294 -- Returns nil if the subset is same size as original sets.
295 function BIN:selects(rows, select,tmp)
296 function select(row, v)
297 v=row.cells[self.col.at]
298 if v=="?" or v==self.lo==self.hi or self.lo<v and v <=self.hu then return row end end
299 tmp= map(rows,select)
300 if #tmp < #rows then return rows end end
301
302

```

```

303 ---- ROWS
304 function ROWS:clone(t) return ROWS():add(self.cols.names):adds(t or {}) end
305
306 function ROWS:file(x) for t in csv(x) do self:add(t) end; return self end
307
308 function ROWS:adds(t) for _,t1 in pairs(t) do self:add(t1) end; return self end
309
310 function ROWS:add(t)
311   if self.cols
312   then self.cols:add(push(i.rows, t.cells and t or ROW(t)))
313   else self.cols=COLS(t) end end
314
315 -- ROWS:mids(?int=2,?[COL]=self.cols.y):[key=num] -- Return 'mid' of columns
316 -- rounded to 'p' places.
317 function ROWS:mids(p,cols)
318   local t=(n=#self.rows)
319   for _,col in pairs(cols or self.cols.y) do t[col.txt]=col:mid(p) end
320   return rnds(t,p or 2) end
321
322 function ROWS:splitter(rows)
323   function split(col)
324     for _,row in pairs(rows) do
325       local v = row.cells[col.at]
326       if v ~= "?" then
327         n=n+1
328         local pos = col:bin(v)
329         dict[pos] = dict[pos] or push(list, BIN(col,v))
330         dict[pos]:add(v,row.label) end end
331   list = col:merges(sort(list,lt"lo", n"the.min)
332   return (bins = list,
333         div = sum(list,function(z) return z.has:div()*z.ys.n/n end))
334   end -----
335   return sort(map(self.cols.x, split),lt"div")[1].bin end
336
337 function ROWS:tree()
338   local bests,rests= self.cols:bests(self.rows)
339   local rows={}
340   for _,best in pairs(bests) do push(rows, best).label=1 end
341   for i = 1,the.rests#bests do push(rows,rests[i]).label=0 end
342   self:grow(rows, (#rows)"the.min)
343   return self end
344
345 function ROWS:grow(rows,stop,when)
346   local function kid(bin)
347     local t = bin:selects(rows)
348     if t then return self:clone(t):grow(t,stop,bin) end end
349   self.when=when
350   if #rows < stop then return self end
351   self.kids = map(self:splitter(rows), kid) end
352

```

```

353 ---- Lib
354 big=math.huge
355 min=math.min
356 max=math.max
357 fmt=string.format
358 rand=math.random
359
360 function any(a) return a[rand(#a)] end
361 function per(t,p) p=p*#t//1; return t[math.max(1,math.min(#t,p))] end
362
363 function push(t,x) t[1+#t]=x; return x end
364 function map(t,f, u) u={};for _,x in pairs(t) do u[1+#u]=f(x);return u end
365 function kap(t,f, u) u={};for k,x in pairs(t) do u[1+#u]=f(k,x);return u end
366 function sum(t,f, u) u=0; for _,x in pairs(t) do u=u+f(x) end;return u end
367
368 ---- rnd(num, places:int):num -- Return 'x' rounded to some number of 'place'.
369 function rnd(x, places) -- 49312;
370   local mult = 10^(places or 2)
371   return math.floor(x * mult + 0.5) / mult end
372 ---- rnds(t:num, places:?int=2):num -- Return items in 't' rounds to 'places'.
373 function rnds(t, places)
374   local u={};for k,x in pairs(t) do u[k]=rnd(x,places or 2);return u end
375
376 function sort(t,f) table.sort(t,f); return t end
377 function lt(x) return function(a,b) return a[x] < b[x] end end
378
379 function shuffle(t, j)
380   for i=#t,2,-1 do j=rand(i); t[i],t[j]=t[j],t[i] end; return t end
381
382 function coerce(x)
383   x = x:match("^%s*(-)%s*$")
384   if x=="true" then return true elseif x=="false" then return false
385   else return math.tointeger(x) or tonumber(x) or x end end
386
387 function cli(t)
388   for k,v in pairs(t) do
389     v = tostring(v)
390     for n,x in pairs(arg) do if x=="-." (k:sub(1,1)) then
391       v = v=="false" and "true" or v=="true" and "false" or arg[n+1] end end
392     t[k] = coerce(v) end
393   if t.help then os.exit(go.help()) end
394   return t end
395
396 function chat(t) print(cat(t)) return t end
397 function cat(t, show,u)
398   function show(k,v) return #t==0 and ("%s%s"):format(k,v) or tostring(v) end
399   u={}; for k,v in pairs(t) do u[1+#u]=show(k,v) end
400   return (t._is or "").."["..table.concat(#t==0 and sort(u) or u," ").."]" end
401
402 function csv(file,fun)
403   function lines(file, fun)
404     local file = io.input(file)
405     while true do
406       local line = io.read()
407       if not line then return io.close(file) else fun(line) end end
408   end
409   function words(s,sep,fun, t)
410     fun = fun or same
411     t={};for x in s:gmatch(fmt("(%s)",sep)) do t[1+#t]=fun(x) end; return t
412   end
413   lines(file, function(line) fun(words(line, ",", coerce)) end) end
414
415 ---- on(tab,tab) -- Runs some (or all) of the demos. Return number of failures.
416 -- Resets 'the' and the random number seed before each demo.
417 function on(the,go)
418   local the, fails, defaults=cli(the), 0, {}
419   for k,v in pairs(the) do defaults[k]=v end
420   local todos = sort(kap(go,function(k,_) return k end))
421   for _,todo in pairs(the.go=="all" and todos or (the.go)) do
422     if type(go[todo])=="function" then
423       for k,v in pairs(defaults) do the[k]=v end
424       math.randomseed(the.seed)
425       if true == go[todo]() then
426         print("FAIL",todo)
427         fails=fails+1 end end end
428   for k,v in pairs(_ENV) do if not b4[k] then print("?",k,type(v)) end end
429   os.exit(fails) end

```

```

429 ---- Start-up
430 function go.the() chat(the) ; return true end
431 function go.help()
432   print(help:gsub("[\u]%\u%d]+","%27[1:32m%127[0m)",""); return true end
433
434 help:gsub("\n[-]S[%s]t[%S]+)(\"n)+= ([%S]+)",function(k,x) the[k]=coerce(x)end)
435 if pcall(debug.getlocal, 4, 1)
436 then return (ROWS=ROWS, the=the)
437 else on(the,go) end

```