```
1  --------------------------------------------------------------------------------
2  ---      _/\          _/\                                           _/\
3  ---     /\ \         /\ \                                         _/\ \
4  ---    /\_\ \     ___\ \ \___      _/\___       _/\___     _/\___ \ \_\ \
5  ---    \/_/\ \   /\___\ \___\    /\ ___\     /\ ___\   /\ ___\ \/_/\ \
6  ---       \ \ \  \/_/\ \/__/    /\ \__/      /\ \__/   /\ \__/     \ \ \
7  ---        \ \_\    \ \_\       \ \____\      \ \____\  \ \____\     \ \_\
8  ---         \/_/     \/_/        \/____/       \/____/   \/____/      \/_/
9  ---
10 ---        .------.
11 ---        |  Ba  | Bad <----.   planning= (better - bad)
12 ---        |   56 |          |   monitor = (bad - better)
13 ---        .------.------.    |
14 ---               |  B   |    v
15 ---               |   5  | Better
16 ---               .------.
17 ---
18 --------------------------------------------------------------------------------
19
20 local b4={}; for k,_ in pairs(_ENV) do b4[k]=k end
21 local help=[[
22
23   -bins   -b    number of bins             = 16
24   -cohen  -c    cohen                      = .35
25   -file   -f    file name                  = ../etc/data/breastcancer.csv
26   -goal   -g    goal                       = recurrence-events
27   -K      -K    manage low class counts    = 1
28   -M      -M    manage low evidence counts = 2
29   -seed   -S    seed                       = 10019
30   -todo   -t    start up action            = nothing
31   -wait   -w    wait                       = 10
32 ]]
33
34 local max,min,ent,per
35 local push,map,sort,up1,upx,down1,slots,up1,down1
36 local words,thing, things, lines
37 local cli
38 local fmt,o,oo
39 local inc,inc2,inc3,has,has2,has3
40 local rogues
41 local classify,test,train,score,nb1,nb2,abcd
42 local bins,nb3
43 local eg,the,ako={},{},{}
44
45 ---      _  _  _ _ ___ ___   _  _ ___  _
46 ---     (_ (_) ||_|| | | |   |_|/ (_  (_
47 ---                        / |
48
49 local ako={}
50 ako.num    = function(x) return x:find"^[A-Z]" end
51 ako.goal   = function(x) return x:find"[-+!]"  end
52 ako.klass  = function(x) return x:find"!$"     end
53 ako.ignore = function(x) return x:find":$"     end
54 ako.less   = function(x) return x:find"-$"     end
55
56 --------------------------------------------------------------------------------
57 -- BSD 2-Clause License
58 -- Copyright (c) 2022, Tim Menzies
59 --
60 -- Redistribution and use in source and binary forms, with or without
61 -- modification, are permitted provided that the following conditions are met:
62 --
63 -- 1. Redistributions of source code must retain the above copyright notice,this
64 --    list of conditions and the following disclaimer.
65 --
66 -- 2. Redistributions in binary form must reproduce the above copyright notice,
67 --    this list of conditions and the following disclaimer in the documentation
68 --    and/or other materials provided with the distribution.
69 --
70 -- THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS"
71 -- AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
72 -- IMPLIED WARRANTIES OF MERCHANTABILITY & FITNESS FOR A PARTICULAR PURPOSE ARE
73 -- DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR CONTRIBUTORS BE LIABLE
74 -- FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL
75 -- DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR
76 -- SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER
77 -- CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY,
78 -- OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE
79 -- OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
80
```

```
80 --------------------------------------------------------------------------------
81 ---      ___  ___  ___  ___
82 ---     |_ | |_  |   |  |
83 ---     |_] |__ _|_  |  |__
84
85 function classify(i,t)
86   local hi,out = -1
87   for h,_ in pairs(i.h) do
88     local prior = ((i.h[h] or 0) + the.K)/(i.n + the.K*i.nh)
89     local l = prior
90     for col,x in pairs(t) do
91       if x ~= "?" and col ~= #t then
92         l=l*(has3(i.e,col,x,h) + the.M*prior)/((i.h[h] or 0) + the.M) end end
93     if l>hi then hi,out=l,h end end
94   return out end
95
96 function test(i,t)
97   if i.n > i.wait then push(i.log,{want=t[#t], got=classify(i,t)}) end  end
98
99 function train(i,t)
100  local more, kl = false, t[#t]
101  for col,x in pairs(t) do
102    if x ~="?" then
103      more = true
104      inc3(i.e, col, x, kl)
105      if col ~= #t then
106        inc2(kl==the.goal and i.best or i.rest, col,x) end end end
107  if more then
108    i.n = i.n + 1
109    if not i.h[kl] then i.nh = i.nh + 1 end
110    inc(i.h, kl)
111    if kl==the.goal then i.bests=i.bests+1 else i.rests=i.rests+1 end end end
112
113 function score(i)
114  local acc,out=0,{}
115  for _,x in pairs(i.log) do if x.want==x.got then acc=acc+1/#i.log end end
116  for col,xns in pairs(i.best) do
117    for x,b in pairs(xns) do
118      local r1 = has2(i.rest,col,x)/i.rests
119      local b1 = b/i.bests
120      push(out, {100*(b1^2/(b1+r1))//1, col,x,b}) end end
121  return acc, sort(out,down1) end
122
123 function nb1(file, log)
124  local i = {h={}, nh=0,e={}, names=nil, n=0, wait=the.wait,
125       bests=0,rests=0,best={}, rest={},log=log or {}}
126  for row in lines(file) do
127    if not i.names then i.names=row else
128      test(i,row); train(i,row) end end
129  return i end
130
131 ---      _        _   _   _      _
132 ---     VV |_ |_ |_|    (/_ VV (_|
133
134 function nb2(file,  log)
135  local tmp, i, create, update, discretize = {}
136  i = {h={}, nh=0,e={}, names=nil, n=0, wait=the.wait,
137       bests=0,rests=0,best={}, rest={},log=log or {},
138       hi={},lo={}, nums={}}
139
140  function create(t)
141    for j,txt in pairs(t) do
142      if ako.num(txt) then i.nums[j] = {lo=1E32, hi=-1E32} end end; return t end
143
144  function update(t,   x)
145    for j,n in pairs(i.nums) do
146      x=t[j]
147      if x~="?" then n.lo=min(x,n.lo); n.hi=max(x,n.hi) end end; return t end
148
149  function discretize(t, x)
150    for j,n in pairs(i.nums) do
151      x=t[j]
152      t[j]=x=="?" and x or (x - n.lo) // ((n.hi - n.lo+1E-32) / the.bins) end end
153
154  tmp={}
155  for row in lines(file) do
156    if not i.names then i.names = create(row) else push(tmp,update(row)) end end
157  for _,row in pairs(tmp) do
158    discretize(row); test(i,row); train(i,row) end
159  return i end
160
161 ---      _    _   _  _
162 ---     |_| |_) (/_|  |  |__|
163
164 function abcd(gotwants, show)
165  local i, exists, add, report, pretty = {
166    data=data or "data", rx= rx or "rx",known={},a={},b={},c={},d={},yes=0,no=0}
167
168  function exists(x,   new)
169    new = not i.known[x]
170    inc(i.known,x)
171    if new then
172      i.a[x]=i.yes + i.no; i.b[x]=0; i.c[x]=0; i.d[x]=0 end end
173
174  function report(   p,out,a,b,c,d,pd,pf,pn,f,acc,g,prec)
175    p = function (z) return math.floor(100*z + 0.5) end
176    out= {}
177    for x,_ in pairs( i.known ) do
178      pd,pf,pn,prec,g,f,acc = 0,0,0,0,0,0,0
179      a= (i.a[x] or 0); b= (i.b[x] or 0); c= (i.c[x] or 0); d= (i.d[x] or 0);
180      if b+d > 0    then pd   = d    / (b+d)          end
181      if a+c > 0    then pf   = c    / (a+c)          end
182      if a+c > 0    then pn   = (b+d) / (a+c)         end
183      if c+d > 0    then prec = d    / (c+d)          end
184      if 1-pf+pd > 0 then g=2*(1-pf) * pd / (1-pf+pd) end
185      if prec+pd > 0 then f=2*prec*pd / (prec + pd)   end
186      if i.yes + i.no > 0 then
187        acc= i.yes / (i.yes + i.no) end
188      out[x] = {data=i.data,rx=i.rx,num=i.yes+i.no,a=a,b=b,c=c,d=d,acc=p(acc),
189             prec=p(prec), pd=p(pd), pf=p(pf),f=p(f), g=p(g), class=x} end
190    return out end
191
192  function pretty(t)
193    print""
194    local s1  = "%10s|%10s|%4s|%4s|%4s|%4s "
195    local s2  = "|%3s|%3s|%3s|%4s|%3s|%3s|"
196    local d,s = "---", (s1 .. s2)
197    print(fmt(s,"db","rx","a","b","c","d","acc","pd","pf","prec","f","g"))
198    print(fmt(s,d,d,d,d,d,d,d,d,d,d,d,d))
199    for _,x in pairs(slots(t)) do
200      local u = t[x]
201      print(fmt(s.."%s", u.data,u.rx,u.a, u.b, u.c, u.d,
202               u.acc, u.pd, u.pf, u.prec, u.f, u.g, x)) end end
203
204  for _,one in pairs(gotwants) do
205    exists(one.want)
206    exists(one.got)
207    if one.want == one.got then i.yes=i.yes+1 else i.no=i.no+1 end
208    for x,_ in pairs(i.known) do
209      if   one.want == x
210      then inc(one.want == one.got and i.d or i.b, x)
211      else inc(one.got  == x        and i.c or i.a, x) end end end
212  return show and pretty(report()) or report() end
```

```
213   -------------------------------------------------------------------
214   ---
215   ---   ⎯⎵⎿⎼⎵⎾  ⏐⎺⎵⎾⏋⎵⎾⎺
216   ---
217
218   function nb3(file,  log)
219     local tmp, i, create, update, discretize1, discretize = {}
220     i = {h={}, nh=0,e={}, names=nil, n=0, wait=the.wait,
221          bests=0,rests=0,best={}, rest={},log=log or {},
222          nums={}}
223
224     function create(t)
225       for j,txt in pairs(t) do
226         if ako.num(txt) then i.nums[j] = {} end end; return t end
227
228     function update(t,     x)
229       for j,n in pairs(i.nums) do
230         x=t[j]
231         if x~="?" then push(n, {x=x, y= t[#t]}) end end; return t end
232
233     function discretize1(t,x)
234       if x == "?" then return x end
235       for j,b in pairs(t) do if b.lo <= x and x < b.hi then return j end end end
236
237     function discretize(t, x)
238       for j,bins in pairs(i.nums) do t[j] = discretize1(bins,t[j]) end end
239
240     tmp={}
241     for row in lines(file) do
242       if not i.names then i.names = create(row) else push(tmp,update(row)) end end
243     for j,xys in pairs(i.nums) do i.nums[j] = bins(xys) end
244     for _,row in pairs(tmp) do
245       discretize(row);
246       test(i,row); train(i,row) end
247     return i end
248
249   ---    �títⅽⅾ  ⎿⅀⏐⎾⎵⎵
250   ---    ~⎿ⅈ⎺⏐⅃  ⎾⅀ⅈ⎾⎵⎵
251
252   function bins(xys)
253     xys  = sort(xys, upx)
254     local cohen    = the.cohen * (per(xys,.9).x - per(xys, .1).x) / 2.54
255     local minItems = #xys / the.bins
256     local out, b4  = {}, -math.huge
257     local function add(f,z)  f[z] = (f[z] or 0) + 1 end
258     local function sub(f,z)  f[z] =  f[z] - 1        end
259     local function argmin(lo,hi)
260       local lhs, rhs, cut, div, xpect, xy = {},{}
261       for j=lo,hi do add(rhs, xys[j].y) end
262       div = ent(rhs)
263       if hi-lo+1 > 2*minItems
264       then
265         for j=lo,hi - minItems do
266           add(lhs, xys[j].y)
267           sub(rhs, xys[j].y)
268           local n1,n2 = j - lo +1, hi-j
269           if   n1          > minItems and        -- enough items (on left)
270                xys[j].x ~= xys[j+1].x and         -- there is a break here
271                xys[j].x  - xys[lo].x > cohen and  -- not trivially small (on left)
272                xys[hi].x - xys[j].x  > cohen      -- not trivially small (on right)
273           then xpect = (n1*ent(lhs) + n2*ent(rhs)) / (n1+n2)
274                if xpect < div then               -- cutting here simplifies things
275                   cut, div = j, xpect end end end --end for
276       end -- end if
277       if   cut
278       then argmin(lo,     cut)
279            argmin(cut+1, hi )
280       else b4 = push(out, {lo=b4, hi=xys[hi].x, n=hi-lo+1, div=div}).hi end
281     end ---------------------------------------
282     argmin(1,#xys)
283     out[#out].hi =  math.huge
284     return out end
```

```
285   --------------------------------------------------------------------
286   ---
287   ---   ⎧⎾⎰⎿⎾
288   ---
289
290   ---
291   ---   ⎺⎾⎺⎿⎿⎵
292
293   min = math.min
294   max = math.max
295
296   function per(t,p) return t[ (p or .5)*#t//1 ] end
297
298   function ent(t)
299     local n=0; for _,m in pairs(t) do n = n+m end
300     local e=0; for _,m in pairs(t) do if m>0 then e= e+m/n*math.log(m/n,2) end end
301     return -e end
302
303   ---
304   ---   ⎵ ⎾⎰⎴ ⎵⎵ ⎸<
305
306   function rogues()
307     for k,v in pairs(_ENV) do if not b4[k] then print("??",k,type(v)) end end end
308
309   ---
310   ---   ⎵⎵⎵⎸⎵⎾⎾⎺⎺
311
312   function inc(f,a,n)       f=f or{};f[a]=(f[a] or 0) + (n or 1) return f end
313   function inc2(f,a,b,n)    f=f or{};f[a]=inc( f[a] or {},b,n);  return f end
314   function inc3(f,a,b,c,n)  f=f or{};f[a]=inc2(f[a] or {},b,c,n);return f end
315
316   function has(f,a)        return  f[a]                       or 0 end
317   function has2(f,a,b)     return  f[a] and has( f[a],b)      or 0 end
318   function has3(f,a,b,c)   return  f[a] and has2(f[a],b,c)  or 0 end
319
320   ---
321   ---   ⎸⎵⎺⎾⎵⎺
322
323   function push(t,x) t[1 + #t] = x; return x end
324
325   function map(t,f,  u) u={};for k,v in pairs(t) do u[1+#u]=f(v) end;return u end
326
327   function sort(t,f) table.sort(t,f); return t end
328
329   function upx(a,b)     return a.x < b.x end
330   function up1(a,b)     return a[1] < b[1] end
331   function down1(a,b) return a[1] > b[1] end
332
333
334   function slots(t, u)
335     local function public(k)  return tostring(k):sub(1,1) ~= "_" end
336     u={};for k,v in pairs(t) do if public(k) then u[1+#u]=k end end
337     return sort(u) end
338
339   ---     '~)
340   ---   ⎵⎺⎾⎾⎺⎾⎵  '⎱  ⎾⎾⎺⎾⎾⎵⎵
341   ---
342
343   function words(s,sep,   t)
344     sep="([^" .. (sep or ",")  .. "]+)"
345     t={}; for y in s:gmatch(sep) do t[1+#t] = y end; return t end
346
347   function things(s)  return map(words(s), thing) end
348
349   function thing(x)
350     x = x:match"^%s*(.-)%s*$"
351     if x=="true" then return true elseif x=="false" then return false end
352     return tonumber(x) or x end
353
354   function lines(file,f,      x)
355     file = io.input(file)
356     f    = f or things
357     return function() x=io.read(); if x then return f(x) else io.close(file) end end end
358
359   ---     '~)
360   ---   ⎾⎸⎺⎵⎾⎾⎵⎵⎵  '⎱  ⎵⎾⎾⎾⎾⎺
361   ---
362
363   fmt = string.format
364
365   function oo(t) print(o(t)) end
366
367   function o(t,  seen, u)
368     if type(t)~="table" then return tostring(t) end
369     seen = seen or {}
370     if seen[t] then return "..." end
371     seen[t] = t
372     local function show1(x)  return o(x, seen) end
373     local function show2(k)  return fmt(":%s %s",k, o(t[k],seen)) end
374     u = #t>0 and map(t,show1)  or map(slots(t),show2)
375     return (t.s or "").."{"..table.concat(u," ").."}" end
376
377   ---
378   ---   ⎵⎸⎸
379
380   function cli(help)
381     local d,used = {},{}
382     help:gsub("\n ([-]([^%s]+))[%s]+(-[^%s]+)[^\n]*%s([^%s]+)",
383       function(long,key,short,x)
384         assert(not used[short], "repeated short flag ["..short.."]")
385         used[short]=short
386         for n,flag in ipairs(arg) do
387           if flag==short or flag==long then
388             x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
389           d[key] = x=="true" and true or thing(x) end)
390     if d.help then os.exit(print(help)) end
391     return d end
392
```

```
393  ---
394  ---  DEMOS
395  ---       []
396
397  function eg.ent()
398    print(ent{a=9,b=7}) end
399
400  function eg.nb1()
401    local i = nb1(the.file);
402    local acc, out = score(i); print(acc); map(out,oo) end
403
404  function eg.nb2()
405    local i = nb2(the.file);
406    local acc, out = score(i); print(acc); map(out,oo) end
407
408  function eg.nb2a()
409    local i = nb2(the.file);
410    local acc, out = score(i)
411    abcd(i.log, true)
412    map(out,oo) end
413
414  function eg.bins(   t)
415    local t,n = {},30
416    for j=1,n do push(t, {x=j, y=j<.6*n and 1 or j<.8*n and 2 or 3}) end
417    map(bins(t),oo)
418  end
419
420  function eg.nb3(  i)
421    print(20)
422    i=nb3("../etc/data/diabetes.csv")
423    for n,bins in pairs(i.nums) do
424      print(n,#bins) end
425    local acc, out = score(i)  -- XXX
426    print(#out)
427    print(acc)
428    map(out,oo)
429    end
430
```

```
431  ---
432  ---  [] START
433  ---    []
434
435  the=cli(help)
436  math.randomseed( the.seed or 10019 )
437  if eg[the.todo] then eg[the.todo]() end
438  rogues()
```