

```

1 #!/usr/bin/env python3
2 # vim: ts=2 sw=2 sts=2 et :
3
4 #
5 #
6 #
7 #
8 #
9 #
10 #
11 #
12 #
13 #
14 #
15 #
16 #
17 #
18 #
19 #
20 #
21 #
22 #
23 #
24 #
25 #
26 #
27 #
28 #
29 #
30 #
31 #sublime.py [OPTIONS]
32 (c)2022 Tim Menzies, unlicense.org.
33 Look, a little, before you leap, a lot. Random projections
34 for bi-clustering. Iterative dichotomization using ranges
35 that most distinguish sibling clusters. Repeat, recursively.
36 Use results for various knowledge-level tasks.
37
38 OPTIONS:
39
40 -Max max numbers to keep :512
41 -Some find 'far' in this many egs:512
42 -data data file :./data/aut093.csv
43 -help show help :False
44 -far how far to look in 'Some' :9
45 -p distance coefficient :2
46 -seed random number seed :10019
47 -todo start up task :nothing
48 -xsmall Cohen's small effect :.35
49
50 """
51 import re,sys,random
52
53 # This is free and unencumbered software released into the public domain.
54 #
55 # Anyone is free to copy, modify, publish, use, compile, sell, or
56 # distribute this software, either in source code form or as a compiled
57 # binary, for any purpose, commercial or non-commercial, and by any
58 # means.
59 #
60 # In jurisdictions that recognize copyright laws, the author or authors
61 # of this software dedicate any and all copyright interest in the
62 # software to the public domain. We make this dedication for the benefit
63 # of the public at large and to the detriment of our heirs and
64 # successors. We intend this dedication to be an overt act of
65 # relinquishment in perpetuity of all present and future rights to this
66 # software under copyright law.
67 #
68 # THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND,
69 # EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF
70 # MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT.
71 # IN NO EVENT SHALL THE AUTHORS BE LIABLE FOR ANY CLAIM, DAMAGES OR
72 # OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE,
73 # ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR
74 # OTHER DEALINGS IN THE SOFTWARE.
75 #
76 # For more information, please refer to <http://unlicense.org/>
77
78
79

```

```

78 #
79 #
80 #
81 #
82 #
83 #
84 #
85 #
86 #
87 # randoms stuff
88 r = random.random
89 anywhere = lambda a: random.randint(0, len(a)-1)
90
91 # useful constants
92 big = sys.maxsize
93
94 # list membership
95 first = lambda a: a[0]
96 second = lambda a: a[1]
97
98 def atom(x):
99     "Return a number or trimmed string."
100     x=x.strip()
101     if x=="True" : return True
102     elif x=="False": return False
103     else:
104         try: return int(x)
105         except:
106             try: return float(x)
107             except: return x.strip()
108
109 def demo(want,all):
110     "Maybe run a demo, if we want it, resetting random seed first."
111     for one in dir(all):
112         if (not want or (want and one.startswith(want))):
113             random.seed(the.seed)
114             all.__dict__[one]()
115
116 def file(f):
117     "Iterator. Returns one row at a time, as cells."
118     with open(f) as fp:
119         for line in fp:
120             line = re.sub(r'([\n\r\t`])#.*', '', line)
121             if line:
122                 yield [atom(cell.strip()) for cell in line.split(",")]
123
124 class o(object):
125     "Class that can pretty print its slots, with fast init."
126     def __init__(i, **d): i.__dict__.update(**d)
127     def __repr__(i):
128         pre = i.__class__.__name__ if isinstance(i,o) else ""
129         return pre+str(
130             {k: v for k, v in sorted(i.__dict__.items()) if str(k)[0] != "_"})
131
132 def options(doc):
133     "Convert __doc__ string to options directory."
134     d={}
135     for line in doc.splitlines():
136         if line and line.startswith(" -"):
137             key, *, x = line.strip()[1:].split(" ") # get 1st,last word on each line
138             for j,flag in enumerate(sys.argv):
139                 if flag and flag[0]=="-" and key.startswith(flag[1:]):
140                     x= "True" if x=="False" else ("False" if x=="True" else sys.argv[j+1])
141             d[key] = atom(x)
142     if d["help"]: exit(print(doc))
143     return o(**d)
144
145 the = options(__doc__)
146

```

```

146 #
147 #
148 #
149 #
150 #
151 #
152 #
153 #-----
154
155 class Range(o):
156     "Track the 'y' symbols seen in the range 'lo' to 'hi'."
157     def __init__(i,col=None, lo=None, hi=None, ys=None):
158         i.col, i.xlo, i.xhi, i.yhas = col, lo, hi, ys or Sym()
159
160     def __add__(i,x,y):
161         if x != "y":
162             i.lo = min(x,i.lo)
163             i.hi = max(x,i.hi)
164             i.ys += y
165             return x
166
167     def merge(i,j):
168         return Range(col=i.col,ys=i.ys.merge(j.ys),lo=min(i.lo,j.lo),hi=max(i.hi,j.h
169 i))
170
171     def __lt__(i,j): return i.val() < j.val()
172
173     def __repr__(i):
174         if i.lo == i.hi: return f"[{i.col.txt}]==[{i.lo}]"
175         if i.lo == -big: return f"[{i.col.txt}]==[{i.hi}]"
176         if i.hi == big: return f"[{i.col.txt}]>=[{i.lo}]"
177         return f"[{i.lo} <= [{i.col.txt} < [{i.hi}]"
178
179     def val(i):
180         z=1E-31; B,R = i.B+z, i.R+z; return (i.b/B)**2/(i.b/B + i.r/R)
181
182     def selects(i,row):
183         x = row[col.at]; return x=="?" or i.lo<=x and x<i.hi
184
185 class Col(o):
186     "Summarize columns."
187     def __init__(i,at=0,txt=""):
188         i.n,i.at,i.txt,i.w=0,at,txt,(-1 if "<" in txt else 1)
189
190     def __add__(i,x,inc=1):
191         if x!="?": i.n += inc; i.add(x,inc)
192         return x
193     def dist(i,x,y): return 1 if x=="?" and y=="?" else i.dist1(x,y)
194
195 class Num(Col):
196     "Summarize numeric columns."
197     def __init__(i,**kw):
198         super().__init__(**kw)
199         i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, the.Max, False
200
201     def add(i,x,_):
202         i.lo = min(x,i.lo)
203         i.hi = max(x,i.hi)
204         if len(i._all) < i.max : i.ok=False; i._all += [x]
205         elif r() < i.max/i.n: i.ok=False; i._all[anywhere(i._all)] = x
206
207     def all(i):
208         if not i.ok: i.ok=True; i._all.sort()
209         return i._all
210
211     def per(i,p=.5):
212         a = i.all(); return a[ int(p*len(a)) ]
213
214     def mid(i): return i.per(.5)
215     def div(i): return (i.per(.9) - i.per(.1)) / 2.56
216
217     def norm(i,x):
218         return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)
219
220     def dist1(i,x,y):
221         if x=="?": y=i.norm(y); x=(1 if y<.5 else 0)
222         elif y=="?": x=i.norm(x); y=(1 if x<.5 else 0)
223         else : x,y = i.norm(x), i.norm(y)
224         return abs(x-y)
225
226     def ranges(i,j, all):
227         # def merge(b4):
228         #     j,n = -1,len(b4)
229         #     while j < n:
230         #         j += 1
231         #         a = b4[j]
232         #         if j < n-1:
233         #             b=b4[j+1]
234         lo = min(i.lo, j.lo)
235         hi = max(i.hi, j.hi)
236         gap = (hi-lo) / (6/the.xsmall)
237         at = lambda z: lo + int((z-lo)/gap)*gap
238         all = {}
239         for x in map(at, i._all): s=all[x]=(all[x] if x in all else Sym()); s.add(1)
240         for x in map(at, j._all): s=all[x]=(all[x] if x in all else Sym()); s.add(0)
241         all = merge(sorted(all.items()),key=first)
242
243 class Sym(Col):
244     "Summarize symbolic columns."
245     def __init__(i,**kw):
246         super().__init__(**kw)
247         i.has, i.mode, i.most = {}, None, 0
248
249     def add(i,x,inc):
250         tmp = i.has[x] = inc + i.has.get(x,0)
251         if tmp > i.most: i.most, i.mode = tmp, x
252
253     def dist(i,x,y): return 0 if x==y else 1
254
255     def div(i):
256         p=lambda x: x/i.n
257         return sum( -p(x)*math.log(p(x),2) for x in i.has.values() )
258
259     def mid(i): return i.mode
260
261     def merge(i,j, reckless=True):
262         k = Sym(at=i.at, txt=i.txt)
263         for k,n in i.has.items(): k.add(x,n)
264         for k,n in j.has.items(): k.add(x,n)
265         if reckless:
266             return k
267         else:
268             if k.div()*1.99 <= (i.n*i.div() + j.n*j.div()/(i.n + j.n)): return k
269
270     def ranges(i,j, all):
271         for x,b in i.has.items(): all += [Range(i,x,x, b,i.n, j.has.get(x,0), j.n)]
272         for x,b in j.has.items(): all += [Range(j,x,x, b,j.n, i.has.get(x,0), i.n)]
273
274
275
276
277
278
279
280 #-----
281
282 class Sample(o):
283     "Load, then manage, a set of examples."
284     def __init__(i,init=[]):
285         i.rows, i.cols, i.x, i.y = [], [], [], []
286         if str==type(inits): [i + row for row in file(inits)]
287         if list==type(inits): [i + row for row in inits]
288
289     def __add__(i,a):
290         def col(at,txt):
291             what = Num if txt[0].isupper() else Sym
292             now = what(at=at, txt=txt)
293             where = i.y if "a" in txt or "-" in txt or "!" in txt else i.x
294             if txt[-1] != ".": where += [now]
295             return now
296
297         #-----
298         if i.cols: i.rows += [[col + a[col.at] for col in i.cols]]
299         else: i.cols = [col(at,txt) for at,txt in enumerate(a)]
300
301     def mid(i,cols=None): return [col.mid() for col in (cols or i.all)]
302     def div(i,cols=None): return [col.div() for col in (cols or i.all)]
303
304     def clone(i,init=[]):
305         out = Sample()
306         out + [col.txt for col in i.cols]
307         [out + x for x in inits]
308         return out
309
310     def dist(i,x,y):
311         d = sum( col.dist(x[col.at], y[col.at])**the.p for col in i.x )
312         return (d/len(i.x)) ** (1/the.p)
313
314     def far(i, x, rows=None):
315         tmp= sorted([(i.dist(x,y),y) for y in (rows or i.rows)],key=first)
316         return tmp[ int(len(tmp)*the.far) ]
317
318     def proj(i,row,x,y,c):
319         a = i.dist(row,x)
320         b = i.dist(row,y)
321         return (a**2 + c**2 - b**2) / (2*c) , row
322
323     def half(i, top=None):
324         top = top or i
325         some = random.choices(i.rows, k=the.Some)
326         w = some[0]
327         _,x = top.far(w, some)
328         c,y = top.far(x, some)
329         left, right = i.clone(), i.clone()
330         for n,(_r) in enumerate(
331             sorted([top.proj(r,x,y,c) for r in i.rows],key=first)):
332             (left if n <= len(i.rows)//2 else right).__add__(r)
333         return left,right
334

```

```

334 #
335 #
336 #
337 #
338 #
339 #
340 #
341 #
342 -----
343 class Demos:
344     "Possible start-up actions."
345     def num():
346         n=Num()
347         for i in range(10000): n + i
348         print(sorted(n._all),n)
349
350     def sym():
351         s=Sym()
352         for i in range(10000): s + int(r()*20)
353         print(s)
354
355     def rows():
356         for row in file(the.data): print(row)
357
358     def sample(): s=Sample(the.data); print(len(s.rows))
359
360     def done(): s=Sample(the.data); s.dist(s.rows[1], s.rows[2])
361
362     def dist():
363         s=Sample(the.data)
364         for row in s.rows: print(s.dist(s.rows[0], row))
365
366     def far():
367         s=Sample(the.data)
368         for row in s.rows: print(row,s.far(row))
369
370     def clone():
371         s=Sample(the.data); s1 = s.clone(s.rows)
372         print(s.x[0])
373         print(s1.x[0])
374
375     def half():
376         s=Sample(the.data); s1,s2 = s.half()
377         print(s1.mid(s1.y))
378         print(s2.mid(s2.y))
379
380 if __name__ == "__main__":
381     demo(the.todo,Demos)

```