

```

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\`'])#.*', '', 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):
158         i.col, i.xlo, i.xhi, i.yhas = col, lo, hi, 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.yhas += y
165         return x
166
167     def merge(i,j):
168         lo = math.min(i.lo, j.lo)
169         hi = math.max(i.hi, j.hi)
170         z = 1E-31
171         B,R = i.B+z, i.R+z
172         k = Range(i.col, lo, hi, i.b+j.b, i.B, i.r+j.r, j.R)
173         if k.b/B < .01 or k.r/R < .01 : return k
174         if k.val() > i.val() and k.val() > j.val(): return k
175
176     def __lt__(i,j): return i.val() < j.val()
177
178     def __repr__(i):
179         if i.lo == i.hi: return f"#{i.col.txt}#{i.lo}"
180         if i.lo == -big: return f"#{i.col.txt}#{i.hi}"
181         if i.hi == big: return f"#{i.col.txt}#{i.lo}"
182         return f"#{i.lo} <= #{i.col.txt} < #{i.hi}"
183
184     def val(i):
185         z=1E-31; B,R = i.B+z, i.R+z; return (i.b/B)**2/( i.b/B + i.r/R)
186
187     def selects(i,row):
188         x = row[col.at]; return x=="?" or i.lo<=x and x<i.hi
189
190 class Col(o):
191     "Summarize columns."
192     def __init__(i,at=0,txt=""):
193         i.n,i.at,i.txt,i.w=0,at,txt,(-1 if "<" in txt else 1)
194
195     def __add__(i,x,inc=1):
196         if x != "y": i.n += inc; i.add(x,inc)
197         return x
198     def dist(i,x,y): return 1 if x=="?" and y=="?" else i.dist1(x,y)
199
200 class Num(Col):
201     "Summarize numeric columns."
202     def __init__(i,**kw):
203         super().__init__(**kw)
204         i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, the.Max, False
205
206     def add(i,x,_):
207         i.lo = min(x,i.lo)
208         i.hi = max(x,i.hi)
209         if len(i._all) < i.max : i.ok=False; i._all += [x]
210         elif r() < i.max/i.n: i.ok=False; i._all[anywhere(i._all)] = x
211
212     def all(i):
213         if not i.ok: i.ok=True; i._all.sort()
214         return i._all
215
216     def per(i,p=.5):
217         a = i.all(); return a[ int(p*len(a)) ]
218
219     def mid(i): return i.per(.5)
220     def div(i): return (i.per(.9) - i.per(.1)) / 2.56
221
222     def norm(i,x):
223         return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)
224
225     def dist1(i,x,y):
226         if x=="?": y=i.norm(y); x=(1 if y<.5 else 0)
227         elif y=="?": x=i.norm(x); y=(1 if x<.5 else 0)
228         else : x,y = i.norm(x), i.norm(y)
229         return abs(x-y)
230
231     def ranges(i,j, all):
232         # def merge(b4):
233         #     j,n = -1,len(b4)
234         #     while j < n:
235         #         j += 1
236         #         a = b4[j]
237         #         if j<n-1:
238         #             b=b4[j+1]
239         lo = min(i.lo, j.lo)
240         hi = max(i.hi, j.hi)
241         gap = (hi-lo) / (6/the.xsmall)
242         at = lambda z: lo + int((z-lo)/gap)*gap
243         all = {}
244         for x in map(at, i._all): s=all[x]=(all[x] if x in all else Sym()); s.add(1)
245         for x in map(at, j._all): s=all[x]=(all[x] if x in all else Sym()); s.add(0)
246         all = merge(sorted(all.items(),key=first))
247
248 class Sym(Col):
249     "Summarize symbolic columns."
250     def __init__(i,**kw):
251         super().__init__(**kw)
252         i.has, i.mode, i.most = {}, None, 0
253
254     def add(i,x,inc):
255         tmp = i.has[x] = inc + i.has.get(x,0)
256         if tmp > i.most: i.most, i.mode = tmp, x
257
258     def dist(i,x,y): return 0 if x==y else 1
259
260     def mid(i): return i.mode
261     def div(i):
262         p=lambda x: x/i.n
263         return sum( -p(x)*math.log(p(x),2) for x in i.has.values() )
264
265     def ranges(i,j, all):
266         for x,b in i.has.items(): all += [Range(i,x,x, b,i.n, j.has.get(x,0), j.n)]
267         for x,b in j.has.items(): all += [Range(j,x,x, b,j.n, i.has.get(x,0), i.n)]
268
269
270
271
272
273
274
275 #-----
276
277 class Sample(o):
278     "Load, then manage, a set of examples."
279     def __init__(i,init=[]):
280         i.rows, i.cols, i.x, i.y = [], [], [], []
281         if str == type(inits): [i + row for row in file(inits)]
282         if list == type(inits): [i + row for row in inits]
283
284     def __add__(i,a):
285         def col(at,txt):
286             what = Num if txt[0].isupper() else Sym
287             now = what(at=at, txt=txt)
288             where = i.y if "a" in txt or "-" in txt or "!" in txt else i.x
289             if txt[-1] != ".": where += [now]
290             return now
291
292         #-----
293         if i.cols: i.rows += [[col + a[col.at] for col in i.cols]]
294         else: i.cols = [col(at,txt) for at,txt in enumerate(a)]
295
296     def mid(i,cols=None): return [col.mid() for col in (cols or i.all)]
297     def div(i,cols=None): return [col.div() for col in (cols or i.all)]
298
299     def clone(i,init=[]):
300         out = Sample()
301         out + [col.txt for col in i.cols]
302         [out + x for x in inits]
303         return out
304
305     def dist(i,x,y):
306         d = sum( col.dist(x[col.at], y[col.at])**the.p for col in i.x )
307         return (d/len(i.x)) ** (1/the.p)
308
309     def far(i, x, rows=None):
310         tmp= sorted([(i.dist(x,y),y) for y in (rows or i.rows)],key=first)
311         return tmp[ int(len(tmp)*the.far) ]
312
313     def proj(i,row,x,y,c):
314         a = i.dist(row,x)
315         b = i.dist(row,y)
316         return (a**2 + c**2 - b**2) / (2*c) , row
317
318     def half(i, top=None):
319         top = top or i
320         some = random.choices(i.rows, k=the.Some)
321         w = some[0]
322         _,x = top.far(w, some)
323         c,y = top.far(x, some)
324         left, right = i.clone(), i.clone()
325         for n,(_r) in enumerate(
326             (left if n <= len(i.rows)//2 else right).__add__(r)
327         ):
328             (left if n <= len(i.rows)//2 else right).__add__(r)
329             return left,right
330

```

```

329 #
330 #
331 #
332 #
333 #
334 #
335 #
336 #
337
338 class Demos:
339     "Possible start-up actions."
340     def num():
341         n=Num()
342         for i in range(10000): n + i
343         print(sorted(n._all),n)
344
345     def sym():
346         s=Sym()
347         for i in range(10000): s + int(r()*20)
348         print(s)
349
350     def rows():
351         for row in file(the.data): print(row)
352
353     def sample(): s=Sample(the.data); print(len(s.rows))
354
355     def done(): s=Sample(the.data); s.dist(s.rows[1], s.rows[2])
356
357     def dist():
358         s=Sample(the.data)
359         for row in s.rows: print(s.dist(s.rows[0], row))
360
361     def far():
362         s=Sample(the.data)
363         for row in s.rows: print(row,s.far(row))
364
365     def clone():
366         s=Sample(the.data); s1 = s.clone(s.rows)
367         print(s.x[0])
368         print(s1.x[0])
369
370     def half():
371         s=Sample(the.data); s1,s2 = s.half()
372         print(s1.mid(s1.y))
373         print(s2.mid(s2.y))
374
375 if __name__ == "__main__":
376     demo(the.todo,Demos)

```