

```

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 <timm@ieee.org> unlicense.org.
33 Sublime's unsupervised bifurcation: let's infer minimal explanations.
34
35 OPTIONS:
36
37 -Max max numbers to keep : 512
38 -Some find 'far' in this many egs : 512
39 -data data file : ./data/aut093.csv
40 -enough min leaf size : 5
41 -help show help : False
42 -far how far to look in 'Some' : 9
43 -p distance coefficient : 2
44 -seed random number seed : 10019
45 -todo start up task : nothing
46 -xsmall Cohen's small effect : .35
47
48 ## See Also
49
50 [issues](issues) &M-^@& [repo](github)
51
52 ## Algorithm
53
54 Stochastic clustering to generate tiny models. Uses random projections
55 then unsupervised iterative dichotomization using ranges that
56 most distinguish sibling clusters.
57
58 ## License
59
60 This is free and unencumbered software released into the public
61 domain.
62
63 Anyone is free to copy, modify, publish, use, compile, sell, or
64 distribute this software, either in source code form or as a compiled
65 binary, for any purpose, commercial or non-commercial, and by any
66 means.
67
68 In jurisdictions that recognize copyright laws, the author or authors
69 of this software dedicate any and all copyright interest in the
70 software to the public domain. We make this dedication for the
71 benefit of the public at large and to the detriment of our heirs
72 and successors. We intend this dedication to be an overt act of
73 relinquishment in perpetuity of all present and future rights to
74 this software under copyright law.
75
76 THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND,
77 EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF
78 MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT.
79 IN NO EVENT SHALL THE AUTHORS BE LIABLE FOR ANY CLAIM, DAMAGES OR
80 OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE,
81 ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE
82 OR OTHER DEALINGS IN THE SOFTWARE.
83
84 For more information, please refer to <http://unlicense.org/>
85 """
86
87 import random, sys, re
88 from random import random as r
89 from typing import Any

```

```

90 #
91 #
92 #
93 #
94 #
95 #
96 #
97 #
98 #
99 #
100 #
101 #
102 #
103 #
104 #
105 #
106 #
107 #
108 #
109 #
110 #
111 #
112 #
113 #
114 #
115 #
116 #
117 #
118 #
119 #
120 #
121 #
122 #
123 #
124 #
125 #
126 #
127 #
128 #
129 #
130 #
131 #
132 #
133 #
134 #
135 #
136 #
137 #
138 #
139 #
140 #
141 #
142 #
143 #
144 #
145 #
146 #
147 #
148 #
149 #
150 #
151 #
152 #
153 #
154 #
155 #
156 #
157 #
158 #
159 #
160 #
161 #
162 #
163 #
164 #
165 #
166 #
167 #
168 #
169 #
170 #
171 #
172 #
173 #
174 #
175 #
176 #
177 #
178 #

```

```

179 #
180 #
181 #
182 #
183 #
184 #
185 #
186 #
187 #
188 #
189 #
190 #
191 #
192 class Span(o):
193     """Given two 'Sample's and some 'x' range 'lo..hi'.
194     a 'Span' holds often that range appears in each 'Sample'."""
195     def __init__(i,col, lo, hi, ys=None):
196         i.col, i.lo, i.hi, i.ys = col, lo, hi, ys or Sym()
197
198     def add(i, x:float, y:Any, inc=1) -> None:
199         "y" is a label identifying one 'Sample' or another."
200         i.lo = min(x,i.lo)
201         i.hi = max(x,i.hi)
202         i.ys.add(y,inc)
203
204     def merge(i, j)-> o|None:
205         "If the merged span is simpler, return that merge."
206         a, b, c = i.ys, j.ys, i.ys.merge(j.ys)
207         if c.div()*.99 <= (a.n*a.div() + b.n*b.div())/(a.n + b.n):
208             return Span(i.col, min(i.lo,j.lo),max(i.hi,j.hi), ys=c)
209
210     def selects(i,row:list) -> bool:
211         "True if the range accepts the row."
212         x = row[col.at]; return x=="?" or i.lo<=x and x<i.hi
213
214     def show(i, positive=True) -> None:
215         "Show the range."
216         txt = i.col.txt
217         if positive:
218             if i.lo == i.hi: return f"[txt] = {i.lo}"
219             elif i.lo == -big: return f"[txt] < {i.hi}"
220             elif i.hi == big: return f"[txt] >= {i.lo}"
221             else: return f"[i.lo] <= [txt] < {i.hi}"
222         else:
223             if i.lo == i.hi: return f"[txt] != {i.lo}"
224             elif i.lo == -big: return f"[txt] >= {i.hi}"
225             elif i.hi == big: return f"[txt] < {i.lo}"
226             else: return f"[txt] < {i.lo} or [txt] >= {i.hi}"
227
228     def support(i) -> float:
229         "Returns 0..1."
230         return i.ys.n / i.col.n
231
232     @staticmethod
233     def sort(spans : list) -> list:
234         "Good spans have large support and low diversity."
235         divs, supports = Num(), Num()
236         sn = lambda s: supports.norm( s.support())
237         dn = lambda s: divs.norm( s.ys.div())
238         f = lambda s: (1 - sn(s))**2 + dn(s)**2**.5
239         for s in spans:
240             divs.add( s.ys.div())
241             supports.add(s.support())
242         return sorted(spans, key=f)
243
244 #
245 #
246 #
247 #
248 #
249 class Col(o):
250     "Summarize columns."
251     def __init__(i,at=0,txt="",):
252         i.n,i.at,i.txt,i.w=0,at,txt,(-1 if "<" in txt else 1)
253
254     def dist(i,x:Any, y:Any) -> float:
255         return 1 if x=="?" and y=="?" else i.dist1(x,y)
256
257 #
258 #
259 #
260 #
261 class Num(Col):
262     "Summarize numeric columns."
263     def __init__(i,**kw):
264         super().__init__(**kw)
265         i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, the.Max, False
266
267     def add(i,x: float ,_):
268         if x != "?":
269             i.n += 1
270             i.lo = min(x,i.lo)
271             i.hi = max(x,i.hi)
272             if len(i._all) < i.max : i.ok=False; i._all += [x]
273             elif r() < i.max/i.n: i.ok=False; i._all[anywhere(i._all)] = x
274         return x
275
276     def all(i):
277         if not i.ok: i.ok=True; i._all.sort()
278         return i._all
279
280     def per(i,p=.5):
281         a = i.all(); return a[ int(p*len(a)) ]
282
283     def mid(i): return i.per(.5)
284     def div(i): return (i.per(.9) - i.per(.1)) / 2.56
285
286     def merge(i,j):
287         k = Num(at=i.at, txt=i.txt)
288         for x in i._all: k.add(x)
289         for x in j._all: k.add(x)
290         return k
291
292     def norm(i,x):
293         return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)
294
295     def dist1(i,x,y):
296         if x=="?": y=i.norm(y); x=(1 if y<.5 else 0)
297         elif y=="?": x=i.norm(x); y=(1 if x<.5 else 0)
298         else: x,y = i.norm(x), i.norm(y)
299         return abs(x-y)
300
301     def spans(i,j, all):
302         lo = min(i.lo, j.lo)
303         hi = max(i.hi, j.hi)
304         gap = (hi-lo) / (6/the.xsmall)
305         at = lambda z: lo + int((z-lo)/gap)*gap
306         tmp = {}
307         for x in map(at, i._all):
308             s = tmp[x] = tmp[x] if x in tmp else Span(i,x,x+gap)
309             s.add(x,0)
310         for x in map(at, j._all):
311             s = tmp[x] = tmp[x] if x in tmp else Span(i,x,x+gap)
312             s.add(x,1)
313         tmp = merge([x for _,x in sorted(tmp.items(),key=first)])
314         if len(tmp) > 1 : all + tmp

```

```

315 #
316 #
317 #
318 #
319 #
320 #
321 class Sym(Col):
322     "Summarize symbolic columns."
323     def __init__(i,**kw):
324         super().__init__(**kw)
325         i.has, i.mode, i.most = {}, None, 0
326
327     def add(i,x,inc):
328         if x != "?":
329             i.n += inc
330             tmp = i.has[x] = inc + i.has.get(x,0)
331             if tmp > i.most: i.most, i.mode = tmp, x
332         return x
333
334     def dist(i,x,y): return 0 if x==y else 1
335
336     def div(i):
337         p = lambda x: x/i.n
338         return sum( -p(x)*math.log(p(x),2) for x in i.has.values() )
339
340     def mid(i): return i.mode
341
342     def merge(i,j):
343         k = Sym(at=i.at, txt=i.txt)
344         for k,n in i.has.items(): k.add(x,n)
345         for k,n in j.has.items(): k.add(x,n)
346         return k
347
348     def spans(i,j, all):
349         tmp = {}
350         for x,n in i.has.items():
351             s = tmp[x] = (tmp[x] if x in tmp else Span(i,x,x))
352             s.add(x,0,n)
353         for x,n in j.has.items():
354             s = tmp[x] = (tmp[x] if x in tmp else Span(i,x,x))
355             s.add(x,1,n)
356         tmp = [second(x) for x in sorted(tmp.items(), key=first)]
357         if len(tmp) > 1 : all + tmp
358
359 #
360 #
361 #
362 #
363 #
364 #
365 class Sample(o):
366     "Load, then manage, a set of examples."
367     def __init__(i,init=[]):
368         i.rows, i.cols, i.x, i.y = [], [], [], []
369         if str==type(init): [i + row for row in file(init)]
370         if list==type(init): [i + row for row in init]
371
372     def __add__(i,a):
373         def col(at,txt):
374             what = Num if txt[0].isupper() else Sym
375             now = what(at=at, txt=txt)
376             where = i.y if "+" in txt or "-" in txt or "!" in txt else i.x
377             if txt[-1] != ":" : where += [now]
378             return now
379
380             if i.cols: i.rows += [[col.add(a[col.at]) for col in i.cols]]
381             else: i.cols = [col(at,txt) for at,txt in enumerate(a)]
382
383     def mid(i,cols=None): return [col.mid() for col in (cols or i.all)]
384     def div(i,cols=None): return [col.div() for col in (cols or i.all)]
385
386     def clone(i,init=[]):
387         out = Sample()
388         out + [col.txt for col in i.cols]
389         [out + x for x in init]
390         return out
391
392     def dist(i,x,y):
393         d = sum( col.dist(x[col.at], y[col.at])**the.p for col in i.x )
394         return (d/len(i.x)) ** (1/the.p)
395
396     def far(i, x, rows=None):
397         tmp= sorted([(i.dist(x,y),y) for y in (rows or i.rows)],key=first)
398         return tmp[ int(len(tmp)*the.far) ]
399
400     def proj(i,row,x,y,c):
401         a = i.dist(row,x)
402         b = i.dist(row,y)
403         return ((a**2 + c**2 - b**2) / (2*c) , row)
404
405     def half(i, top=None):
406         top = top or i
407         some = random.choices(i.rows, k=the.Some)
408         w = some[0]
409         _,x = top.far(w, some)
410         c,y = top.far(x, some)
411         left, right = i.clone(), i.clone()
412         for n, (_,r) in enumerate(
413             sorted([top.proj(r,x,y,c) for r in i.rows],key=first)):
414             (left if n <= len(i.rows)//2 else right)._add__(r)
415         return left,right
416
417     def split(i,top=None):
418         here = Tree(i)
419         top = top or i
420         if len(i.rows) >= 2*len(top.rows)**the.enough:
421             left0, right0 = i.half(top)
422             spans = []
423             [icol.spans(rcol,spans) for lcol,rcol in zip(left0.x, right0.x)]
424             if len(spans) > 0:
425                 here.when = Span.sort(spans)[0]
426                 left, right = i.clone(), i.clone()
427                 [(left if span.selects(row) else right).add(row) for row in i.rows]
428                 if len(left.rows) < len(i.rows): here.left = left.split(top)
429                 if len(right.rows) < len(i.rows): here.right = right.split(top)
430             return here
431
432 #
433 #
434 #
435 #
436 #
437 class Tree(i):
438     def __init__(i,here):
439         i.here, i.when, i.yes, i.no = here, None, None, None
440
441     def show(i,pre=""):
442         "Print tree with indents."
443         print(f"[pre]{i.here.ys.n}")
444         if i.yes:
445             print(f"[pre] {i.when.show(True)}") ; i.yes.show(pre + "|. ")
446         if i.no:
447             print(f"[pre] {i.when.show(False)}") ; i.no.show( pre + "|. ")

```

```

448 #
449 #
450 #
451 #
452 #
453 #
454 #
455
456 class Demos:
457     "Possible start-up actions."
458     def opt():
459         print(the)
460
461     def num():
462         n=Num()
463         for x in range(10000): n.add(x)
464         print(sorted(n._all),n)
465
466     def sym():
467         s=Sym()
468         for x in range(10000): s.add( int(r()*20))
469         print(s)
470
471     def rows():
472         for row in file(the.data): print(row)
473
474     def sample(): s=Sample(the.data); print(len(s.rows))
475
476     def done(): s=Sample(the.data); s.dist(s.rows[1], s.rows[2])
477
478     def dist():
479         s=Sample(the.data)
480         for row in s.rows: print(s.dist(s.rows[0], row))
481
482     def far():
483         s=Sample(the.data)
484         for row in s.rows: print(row,s.far(row))
485
486     def clone():
487         s=Sample(the.data); s1 = s.clone(s.rows)
488         print(s.x[0])
489         print(s1.x[0])
490
491     def half():
492         s=Sample(the.data); s1,s2 = s.half()
493         print(s1.mid(s1.y))
494         print(s2.mid(s2.y))
495
496 if __name__ == "__main__":
497     demo(the.todo,Demos)

```