

```

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:
34 let's infer minimal explanations.
35
36 OPTIONS:
37
38 -Max      max numbers to keep      : 512
39 -Some     find 'far' in this many eggs : 512
40 -cautious On any crash, stop+show stack : False
41 -data     data file                  : ./data/aut093.csv
42 -enough   min leaf size              : 5
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 ## See Also
51
52 [issues](issues) â-^@% [repo](github)
53
54 ## Algorithm
55
56 Stochastic clustering to generate tiny models. Uses random projections
57 then unsupervised iterative dichotomization using ranges that
58 most distinguish sibling clusters.
59
60 ## License
61
62 This is free and unencumbered software released into the public
63 domain.
64
65 Anyone is free to copy, modify, publish, use, compile, sell, or
66 distribute this software, either in source code form or as a compiled
67 binary, for any purpose, commercial or non-commercial, and by any
68 means.
69
70 In jurisdictions that recognize copyright laws, the author or authors
71 of this software dedicate any and all copyright interest in the
72 software to the public domain. We make this dedication for the
73 benefit of the public at large and to the detriment of our heirs
74 and successors. We intend this dedication to be an overt act of
75 relinquishment in perpetuity of all present and future rights to
76 this software under copyright law.
77
78 THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND,
79 EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF
80 MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT.
81 IN NO EVENT SHALL THE AUTHORS BE LIABLE FOR ANY CLAIM, DAMAGES OR
82 OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE,
83 ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE
84 OR OTHER DEALINGS IN THE SOFTWARE.
85
86 For more information, please refer to <http://unlicense.org/>
87
88 """
89 import traceback, random, math, sys, re
90 from random import random as r
91 from typing import Any

```

```

92 #
93 #
94 #
95 #
96 #
97 #
98 #
99 #
100 def any(a:list) -> Any:
101     "Return a random item."
102     return a[anywhere(a)]
103
104 def anywhere(a:list) -> int:
105     "Return a random index of list 'a'."
106     return random.randint(0, len(a)-1)
107
108 big = sys.maxsize
109
110 def atom(x):
111     "Return a number or trimmed string."
112     x=x.strip()
113     if x=="True": return True
114     elif x=="False": return False
115     else:
116         try: return int(x)
117         except:
118             try: return float(x)
119             except: return x.strip()
120
121 def demo(do,all):
122     "Maybe run a demo, if we want it, resetting random seed first."
123     todo = dir(all)
124     if do and do != "all":
125         todo = [x for x in dir(all) if x.startswith(do)]
126     for one in todo:
127         fun = all.__dict__.get(one,"")
128         if type(fun)==type(demo):
129             random.seed(the.seed)
130             doc = re.sub(r"^\n\s+", "\n", fun.__doc__ or "")
131             try:
132                 fun()
133                 print("PASS:", doc)
134             except Exception as e:
135                 all.fails += 0
136                 if the.cautious: traceback.print_exc(); exit(1)
137             else:
138                 print("FAIL:", doc, e)
139             exit(all.fails)
140
141 def file(f):
142     "Iterator. Returns one row at a time, as cells."
143     with open(f) as fp:
144         for line in fp:
145             line = re.sub(r'([\n\r\v\ ]#*)', '', line)
146             if line:
147                 yield [atom(cell.strip()) for cell in line.split(",")]
148
149 def first(a:list) -> Any:
150     "Return first item."
151     return a[0]
152
153 def merge(b4:list) -> list:
154     "While we can find similar adjacent things, merge them."
155     j,n,now = -1,len(b4),[]
156     while j < n-1:
157         j += 1
158         a = b4[j]
159         if j < n-2:
160             if merged := a.merge(b4[j+1]):
161                 a = merged
162                 j += 1 # we will continue, after missing one
163                 now += [a]
164             # if 'now' is same size as 'b4', look for any other merges.
165             return b4 if len(now)==len(b4) else merge(now)
166
167 class o(object):
168     "Class that can pretty print its slots, with fast inits."
169     def __init__(i, **d): i.__dict__.update(**d)
170     def __repr__(i):
171         pre = i.__class__.__name__ if isinstance(i,o) else ""
172         return pre+str(
173             {k: v for k, v in sorted(i.__dict__.items()) if str(k)[0] != "_"})
174
175 def options(doc:str) -> o:
176     """Convert 'doc' to options dictionary using command line args.
177     Args cause two 'shortands': (1) boolean flags have no arguments (and mentioning
178     those on the command line means 'flip the default value'; (2) args need only
179     mention the first few of a key (e.g. -s is enough to select for -seed)."""
180     d={}
181     for line in doc.splitlines():
182         if line and line.startswith(" -"):
183             key, *, x = line.strip()[1:].split(" ") # get 1st,last word on each line
184             for j,flag in enumerate(sys.argv):
185                 if flag and flag[0]=="-" and key.startswith(flag[1:]):
186                     x= "True" if x=="False" else ("False" if x=="True" else sys.argv[j+1])
187             d[key] = atom(x)
188     if d["help"]: exit(print(re.sub(r'\n#.*', "", doc, flags=re.S)))
189     return o(**d)
190
191 def r() -> float:
192     "Return random number 0..1"
193     return random.random()
194
195 def second(a:list) -> Any:
196     "Return second item."
197     return a[1]
198
199 the = options(__doc__)

```

```

199 #
200 #
201 #
202 #
203 #
204 #
205 #
206 #
207 #
208 #
209 #
210 #
211 #
212 class Span(o):
213     """Given two 'Sample's and some 'x' range 'lo..hi'.
214     a 'Span' holds often that range appears in each 'Sample'."""
215     def __init__(i,col, lo, hi, ys=None):
216         i.col, i.lo, i.hi, i.ys = col, lo, hi,  ys or Sym()
217
218     def add(i, x:float, y:Any, inc=1) -> None:
219         "y" is a label identifying one 'Sample' or another."
220         i.lo = min(x,i.lo)
221         i.hi = max(x,i.hi)
222         i.ys.add(y,inc)
223
224     def merge(i, j): # -> Span|None
225         "If the merged span is simpler, return that merge."
226         a, b, c = i.ys, j.ys, i.ys.merge(j.ys)
227         if c.div()*.99 <= (a.n*a.div() + b.n*b.div())/(a.n + b.n):
228             return Span(i.col, min(i.lo,j.lo),max(i.hi,j.hi), ys=c)
229
230     def selects(i,row:list) -> bool:
231         "True if the range accepts the row."
232         x = row[col.at]; return x=="?" or i.lo<=x and x<i.hi
233
234     def show(i, positive=True) -> None:
235         "Show the range."
236         txt = i.col.txt
237         if positive:
238             if i.lo == i.hi: return f"[txt] == {i.lo}"
239             elif i.lo == -big: return f"[txt] < {i.hi}"
240             elif i.hi == big: return f"[txt] >= {i.lo}"
241             else: return f"[i.lo] <= [txt] < {i.hi}"
242         else:
243             if i.lo == i.hi: return f"[txt] != {i.lo}"
244             elif i.lo == -big: return f"[txt] >= {i.hi}"
245             elif i.hi == big: return f"[txt] < {i.lo}"
246             else: return f"[txt] < {i.lo} or [txt] >= {i.hi}"
247
248     def support(i) -> float:
249         "Returns 0..1."
250         return i.ys.n / i.col.n
251
252     @staticmethod
253     def sort(spans : list) -> list:
254         "Good spans have large support and low diversity."
255         divs, supports = Num(), Num()
256         sn = lambda s: supports.norm( s.support())
257         dn = lambda s: divs.norm( s.ys.div())
258         f = lambda s: ((1 - sn(s))*2 + dn(s))*2**.5
259         for s in spans:
260             divs.add( s.ys.div())
261             supports.add(s.support())
262         return sorted(spans, key=f)
263
264 #
265 #
266 #
267 #
268 #
269 class Col(o):
270     "Summarize columns."
271     def __init__(i,at=0,txt=""):
272         i.n,i.at,i.txt,i.w=0,at,txt, (-1 if "<" in txt else 1)
273
274     def dist(i,x:Any, y:Any) -> float:
275         return 1 if x=="?" and y=="?" else i.dist1(x,y)
276
277 #
278 #
279 #
280 #
281 class Num(Col):
282     "Summarize numeric columns."
283     def __init__(i,**kw):
284         super().__init__(**kw)
285         i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, the.Max, False
286
287     def add(i,x: float ,inc=1):
288         "Reservoir sampler. If '_all' is full, sometimes replace an item at random."
289         if x != " ":
290             i.n += inc
291             i.lo = min(x,i.lo)
292             i.hi = max(x,i.hi)
293             if len(i._all) < i.max : i.ok=False; i._all += [x]
294             elif r() < i.max/i.n: i.ok=False; i._all[anywhere(i._all)] = x
295         return x
296
297     def all(i):
298         "Return '_all', sorted."
299         if not i.ok: i.ok=True; i._all.sort()
300         return i._all
301
302     def dist1(i,x,y):
303         if x=="?": y=i.norm(y); x=(1 if y<.5 else 0)
304         elif y=="?": x=i.norm(x); y=(1 if x<.5 else 0)
305         else: x,y = i.norm(x), i.norm(y)
306         return abs(x-y)
307
308     def div(i):
309         """Report the diversity of this distribution (using standard deviation).
310         &pm;2, 2.56, 3 &sigma; is 66.90,95% of the mass. 2&sigma; &sigma;. So one
311         standard deviation is (90-10)th divide by 2.4 times &sigma;. """
312         return (i.per(.9) - i.per(.1)) / 2.56
313
314     def merge(i,j):
315         "Return two 'Num's."
316         k = Num(at=i.at, txt=i.txt)
317         for x in i._all: k.add(x)
318         for x in j._all: k.add(x)
319         return k
320
321     def mid(i):
322         "Return central tendency of this distribution (using median)."
323         return i.per(.5)
324
325     def norm(i,x):
326         "Normalize 'x' to the range 0..1."
327         return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)
328
329     def per(i,p:float=.5) -> float:
330         "Return the p-th ranked item."
331         a = i.all(); return a[ int(p*len(a)) ]
332
333     def spans(i,j, all):
334         """Divide the whole space 'lo' to 'hi' into, say, 'xsmall'=16 bin,

```

```

335     then count the number of times we the bin on other side.
336     Then merge similar adjacent bins."""
337     lo = min(i.lo, j.lo)
338     hi = max(i.hi, j.hi)
339     gap = (hi-lo) / (6/the.xsmall)
340     at = lambda z: lo + int((z-lo)/gap)*gap
341     tmp = {}
342     for x in map(at, i._all):
343         s = tmp[x] = tmp[x] if x in tmp else Span(i,x,x+gap)
344         s.add(x,0)
345     for x in map(at, j._all):
346         s = tmp[x] = tmp[x] if x in tmp else Span(i,x,x+gap)
347         s.add(x,1)
348     tmp = merge([x for x in sorted(tmp.items(),key=first)])
349     if len(tmp) > 1 : all + tmp
350
351 #
352 #
353 #
354 #
355 class Sym(Col):
356     "Summarize symbolic columns."
357     def __init__(i,**kw):
358         super().__init__(**kw)
359         i.has, i.mode, i.most = {}, None, 0
360
361     def add(i, x:str, inc:int=1) -> str:
362         "Update symbol counts in 'has', updating 'mode' as we go."
363         if x != " ":
364             i.n += inc
365             tmp = i.has[x] = inc + i.has.get(x,0)
366             if tmp > i.most: i.most, i.mode = tmp, x
367         return x
368
369     def dist(i,x:str, y:str) ->float:
370         "Distance between two symbols."
371         return 0 if x==y else 1
372
373     def div(i):
374         "Return diversity of this distribution (using entropy)."
375         p = lambda x: x/i.n
376         return sum( -p(x)*math.log(p(x),2) for x in i.has.values() )
377
378     def merge(i,j):
379         "Merge two 'Sym's."
380         k = Sym(at=i.at, txt=i.txt)
381         for k,n in i.has.items(): k.add(x,n)
382         for k,n in j.has.items(): k.add(x,n)
383         return k
384
385     def mid(i):
386         "Return central tendency of this distribution (using mode)."
387         return i.mode
388
389     def spans(i,j, all):
390         """For each symbol in 'i' and 'j', count the
391         number of times we see it on either side."""
392         tmp = {}
393         for x,n in i.has.items():
394             s = tmp[x] = (tmp[x] if x in tmp else Span(i,x,x))
395             s.add(x,0,n)
396         for x,n in j.has.items():
397             s = tmp[x] = (tmp[x] if x in tmp else Span(i,x,x))
398             s.add(x,1,n)
399         tmp = [second(x) for x in sorted(tmp.items(), key=first)]
400         if len(tmp) > 1 : all + tmp
401
402 #
403 #
404 #
405 #
406 #
407 #
408 class Sample(o):
409     "Load, then manage, a set of examples."
410     def __init__(i, inits=[]):
411         i.rows, i.cols, i.x, i.y = [], [], [], []
412         if str == type(inits): [i + row for row in file(inits)]
413         if list == type(inits): [i + row for row in inits]
414
415     def __add__(i,a):
416         def col(at,txt):
417             what = Num if txt[0].isupper() else Sym
418             now = what(at=at, txt=txt)
419             where = i.y if "+" in txt or "-" in txt or "!" in txt else i.x
420             if txt[-1] != ".": where += [now]
421             return now
422         if i.cols: i.rows += [[col.at(a[col.at]) for col in i.cols]]
423         else: i.cols = [col(at,txt) for at,txt in enumerate(a)]
424
425     def clone(i, inits=[]):
426         out = Sample()
427         out + [col.txt for col in i.cols]
428         [out + x for x in inits]
429         return out
430
431     def dist(i,x,y):
432         d = sum( col.dist(x[col.at], y[col.at])*the.p for col in i.x )
433         return (d/len(i.x)) ** (1/the.p)
434
435     def div(i,cols=None): return [col.div() for col in (cols or i.all)]
436
437     def far(i, x, rows=None):
438         tmp = sorted([(i.dist(x,y),y) for y in (rows or i.rows)],key=first)
439         return tmp[ int(len(tmp)*the.far) ]
440
441     def half(i, top=None):
442         top = top or i
443         some = random.choices(i.rows, k=the.Some)
444         w = some[0]
445         _,x = top.far(w, some)
446         c,y = top.far(x, some)
447         left, right = i.clone(), i.clone()
448         for n,_,r in enumerate(
449             sorted([top.proj(r,x,y,c) for r in i.rows],key=first)):
450             (left if n <= len(i.rows)//2 else right).__add__(r)
451         return left,right
452
453     def mid(i,cols=None):
454         return [col.mid() for col in (cols or i.all)]
455
456     def proj(i,row,x,y,c):
457         "Find the distance of a 'row' on a line between 'x' and 'y'."
458         a = i.dist(row,x)
459         b = i.dist(row,y)
460         return ((a**2 + c**2 - b**2) / (2*c) , row)
461
462     def split(i,top=None):
463         """Split the data using random projections. Find the span that most
464         separates the data. Divide data on that span."""
465         here = Tree(i)
466         top = top or i
467         if len(i.rows) >= 2*len(top.rows)**the.enough:
468             left, right = i.half(top)
469

```

```

471 spans = []
472 [lcol.spans(rcol,spans) for lcol,rcol in zip(left.x, right.x)]
473 if len(spans) > 0:
474     here.span = Span.sort(spans)[0]
475     yes, no = i.clone(), i.clone()
476     [(yes if span.selects(row) else no).add(row) for row in i.rows]
477     if len(yes.rows) < len(i.rows): here.yes = yes.split(top)
478     if len(no.rows) < len(i.rows): here.no = no.split(top)
479     return here
480
481 #
482 #
483 #
484 #
485 #
486 class Tree(o):
487     """Binary tree that splits into 'yes','no' branches containing samples
488     that do/do not match a 'span'."""
489     def __init__(i,here):
490         i.here, i.span, i.yes, i.no = here, None, None, None
491
492     def show(i,pre=""):
493         """Print tree with indents."""
494         print(f"{pre}{i.span.ys.n}")
495         if i.yes:
496             print(f"{pre} {i.span.show(True)}") ; i.yes.show(pre + "|. ")
497         if i.no:
498             print(f"{pre} {i.span.show(False)}") ; i.no.show( pre + "|. ")

```

```

#
#
#
#
#

```

```

499 #
500 #
501 #
502 #
503 #
504 #
505 #
506 #
507 class Demos:
508     """Possible start-up actions."
509     fails=0
510     def _opt():
511         """show the config"
512         [print(f"{k}>10}={v}") for k,v in the.__dict__.items()]
513
514     def seed():
515         """seed"
516         assert .494 <= r() <= .495
517
518     def num():
519         """check 'Num'."
520         n = Num()
521         for _ in range(100): n.add(r())
522         assert .30 <= n.div() <= .31, "in range"
523
524     def sym():
525         """check 'Sym'."
526         s = Sym()
527         for x in "aaaabbc": s.add(x)
528         assert 1.37 <= s.div() <= 1.38, "entropy"
529         assert 'a' == s.mid(), "mode"
530
531     def rows():
532         """count rows in a file."
533         assert 399 == len([row for row in file(the.data)])
534
535     def sample():
536         """sampling."
537         s=Sample(the.data)
538         assert 398 == len(s.rows), "length of rows"
539         assert 249 == s.x[-1].has[1], "symbol counts"
540
541     def dist():
542         """distance between rows"
543         s=Sample(the.data)
544         assert .84 <= s.dist(s.rows[1], s.rows[-1]) <= .842
545
546     def far():
547         """distant items"
548         s=Sample(the.data)
549         for _ in range(32):
550             a,_ = s.far(any(s.rows))
551             assert a>.5, "large?"
552
553     def clone():
554         """cloning"
555         s=Sample(the.data)
556         s1 = s.clone(s.rows)
557         d1,d2 = s.x[0].__dict__, s1.x[0].__dict__
558         for k,v in d1.items():
559             assert d2[k] == v, "clone test"
560
561     def half():
562         """divide data in two"
563         s=Sample(the.data); s1,s2 = s.half()
564         print(s1.mid(s1.y))
565         print(s2.mid(s2.y))
566
567 if __name__ == "__main__": demo(the.todo,Demos)

```

```

#
#
#
#
#
#
#
#

```