

```

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 #
32 #
33 #
34 #
35 #
36 #
37 #
38 #
39 #
40 #
41 #
42 #
43 #
44 #
45 #
46 #
47 #
48 #
49 #
50 #
51 #
52 #
53 #
54 #
55 #
56 #
57 #
58 #
59 #
60 #
61 #
62 #
63 #
64 #
65 #
66 #
67 #
68 #
69 #
70 #
71 #
72 #
73 #
74 #
75 #
76 #
77 #
78 #
79 #
80 #
81 #
82 #
83 #
84 #
85 #
86 #
87 #
88 #
89 #
90 #
91 #
92 #
93 #
94 #
95 #
96 #
97 #
98 #
99 #
100 #
101 #
102 #
103 #
104 #
105 #
106 #
107 #
108 #
109 #
110 #

```

I L E X n s
 m i n a m

OPTIONS:
 -Max max numbers to keep : 512
 -Some find 'far' in this many eggs : 512
 -cautious On any crash, stop+show stack : False
 -data data file : data/auto93.csv
 -enough min leaf size : .5
 -help show help : False
 -far how far to look in 'Some' : .9
 -p distance coefficient : 2
 -seed random number seed : 10019
 -todo start up task : nothing
 -xsmall Cohen's small effect : .35

See Also

[issues](https://github.com/timm/sublime/issues)
 :: [repo](https://github.com/timm/sublime)
 :: [view source](https://github.com/timm/sublime/blob/main/docs/pdf)

<img
 src=https://github.com/timm/sublime/actions/workflows/main.yml/badge.svg>

](https://doi.org/10.5281/zenodo.5912461)

Algorithm

Stochastic clustering to generate tiny models. Uses random projections
 to divide the space. Then, optionally, explain the clusters by
 unsupervised iterative dichotomization using ranges that most
 distinguish sibling clusters.

Example1: just bi-cluster on two distant points

```

56 ...
57
58 /sublime.py -c -s $RANDOM -t cluster
59
60 : 398
61 : 199
62 : 99
63 : 49 Lbs- Acc+ Mpg+
64 : 24 : [2255, 15.5, 30]
65 : 25 : [2575, 16.4, 30]
66 : 50
67 : 25 : [2110, 16.4, 30] <== best
68 : 25 : [2205, 16, 30]
69 : 100
70 : 50
71 : 25 : [2234, 15.5, 30]
72 : 25 : [2278, 16.5, 30]
73 : 50
74 : 25 : [2220, 15.5, 30]
75 : 25 : [2320, 15.8, 30]
76 : 199
77 : 99
78 : 49
79 : 24 : [2451, 16.5, 20]
80 : 25 : [3021, 15.5, 20]
81 : 50
82 : 25 : [3425, 17.6, 20]
83 : 25 : [3155, 16.7, 20]
84 : 100
85 : 50
86 : 25 : [4141, 13.5, 10]
87 : 25 : [4054, 13.2, 20]
88 : 50
89 : 25 : [4425, 11, 10]
90 : 25 : [4129, 13, 10]
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110

```

Example2: as above but split on range that most divides data

```

93 ...
94
95 /sublime.py -c -s $RANDOM -t xplain
96
97 : 398 : [2807, 15.5, 20]
98 : 167 : [3725, 14.5, 20]
99 : 34 : [3609, 13, 20]
100 : 133 : [3735, 14.9, 20]
101 : 56 : [3336, 17, 20]
102 : 22 : [3410, 17.1, 20]
103 : 34 : [3233, 17, 20]
104 : 77 : [4129, 13.2, 20]
105 : 37 : [4274, 13, 10]
106 : 40 : [3962, 13.5, 20]
107 : 35 : [4054, 13.2, 20]
108 : 231 : [2290, 16, 30] <== best
109
110

```

```

111 ## License
112
113 **BSD 2-clause license:**
114 Redistribution and use in source and binary forms, with or without
115 modification, are permitted provided that the following conditions are met:
116 1. Redistributions of source code must retain the above copyright notice, this
117 list of conditions and the following disclaimer.
118 2. Redistributions in binary form must reproduce the above copyright notice,
119 this list of conditions and the following disclaimer in the documentation
120 and/or other materials provided with the distribution.
121
122 THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND
123 ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
124 IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR
125 PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR
126 CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL,
127 EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO,
128 PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR
129 PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF
130 LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING
131 NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS
132 SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
133
134
135 import traceback, random, copy, math, sys, re
136 import random as rnd
137 from typing import Any
138 r = rnd.random

```

```

139 #
140 #
141 #
142 #
143 #
144 #
145 #
146 #
147 def any(a:list) -> Any:
148     "Return a random item."
149     return a[anywhere(a)]
150
151 def anywhere(a:list) -> int:
152     "Return a random index of list 'a'."
153     return rnd.randint(0, len(a)-1)
154
155 big = sys.maxsize
156
157 def atom(x):
158     "Return a number or trimmed string."
159     x=x.strip()
160     if x=="True": return True
161     elif x=="False": return False
162     else:
163         try: return int(x)
164         except:
165             try: return float(x)
166             except: return x.strip()
167
168 def demo(do,all):
169     "Maybe run a demo, if we want it, resetting random seed first."
170     todo = dir(all)
171     if do and do != "all":
172         todo = [x for x in dir(all) if x.startswith(do)]
173     for one in todo:
174         fun = all.__dict__.get(one,"")
175         if type(fun)==type(demo):
176             rnd.seed(the.seed)
177             doc = re.sub(r'\n\s+', "\n", fun.__doc__ or "")
178             try:
179                 fun()
180                 print("PASS:", doc)
181             except Exception as e:
182                 all.fails += 0
183                 if the.cautious: traceback.print_exc(); exit(1)
184                 else: print("FAIL:", doc, e)
185     exit(all.fails)
186
187 def file(f):
188     "Iterator. Returns one row at a time, as cells."
189     with open(f) as fp:
190         for line in fp:
191             line = re.sub(r'([\n\r\v\`"]#*)', '', line)
192             if line:
193                 yield [cell.strip() for cell in line.split(",")]
194
195 def first(a:list) -> Any:
196     "Return first item."
197     return a[0]
198
199 def merge(b4:list) -> list:
200     "While we can find similar adjacent things, merge them."
201     j,n,now = -1,len(b4),[]
202     while j < n-1:
203         j += 1
204         a = b4[j]
205         if j < n-2:
206             if merged := a.merge(b4[j+1]):
207                 a = merged
208                 j += 1 # we will continue, after missing one
209             now += [a]
210     # if 'now' is same size as 'b4', look for any other merges.
211     return b4 if len(now)==len(b4) else merge(now)
212
213 class o(object):
214     "Class that can pretty print its slots, with fast inits."
215     def __init__(i, **d): i.__dict__.update(**d)
216     def __repr__(i):
217         pre = i.__class__.__name__ if isinstance(i,o) else ""
218         return pre+'('+' '.join([f'{k} {v}' for k, v in
219                                 sorted(i.__dict__.items()) if str(k)[0] != "_"])+')'+'?'
220
221 def options(doc:str) -> o:
222     """Convert 'doc' to options dictionary using command line args.
223     Args cause two 'shortands': (1) boolean flags have no arguments (and mentioning
224     those on the command line means 'flip the default value'; (2) args need only
225     mention the first few of a key (e.g. -s is enough to select for -seed)."""
226     d={}
227     for line in doc.splitlines():
228         if line and line.startswith(" -"):
229             key, *, x = line.strip()[1:].split("#") # get 1st,last word on each line
230             for j,flag in enumerate(sys.argv):
231                 if flag and flag[0]=="-" and key.startswith(flag[1:]):
232                     x="True" if x=="False" else "False" if x=="True" else sys.argv[j+1])
233             d[key] = atom(x)
234     if d["help"]: exit(print(re.sub(r'\n#.*', "",doc,flags=re.S)))
235     return o(**d)
236
237 def per(a, p=.5):
238     "Return the p-th item in 'a'."
239     return a[ int(p*len(a)) ]
240
241 def r() -> float:
242     "Return random number 0..1"
243     return rnd.random()
244
245 def rn(x:float, n=3) -> float:
246     "Round a number to three decimals."
247     return round(x,n)
248
249 def rN(a:list, n=3) -> list:
250     "Round a list of numbers to three decimals."
251     return [rn(x,n=n) for x in a]
252
253 def second(a:list) -> Any:
254     "Return second item."
255     return a[1]

```

```

256 #
257 #
258 #
259 #
260 #
261 #
262 #
263 #
264 #
265 #
266 #
267 #
268 #
269 #
270 #
271 #
272 #
273 #
274 #
275 #
276 #
277 #
278 #
279 #
280 #
281 #
282 #
283 #
284 #
285 #
286 #
287 #
288 #
289 #
290 #
291 #
292 #
293 #
294 #
295 #
296 #
297 #
298 #
299 #
300 #
301 #
302 #
303 #
304 #
305 #
306 #
307 #
308 #
309 #
310 #
311 #
312 #
313 #
314 #
315 #
316 #
317 #
318 #
319 #
320 #
321 #
322 #
323 #
324 #
325 #
326 #
327 #
328 #
329 #
330 #
331 #
332 #
333 #
334 #
335 #
336 #
337 #
338 #
339 #
340 #
341 #
342 #
343 #
344 #
345 #
346 #
347 #
348 #
349 #
350 #
351 #
352 #
353 #
354 #
355 #
356 #
357 #
358 #
359 #
360 #
361 #
362 #
363 #
364 #
365 #
366 #
367 #
368 #
369 #
370 #
371 #
372 #
373 #
374 #
375 #
376 #
377 #
378 #
379 #
380 #
381 #
382 #
383 #
384 #
385 #
386 #
387 #
388 #
389 #
390 #
391 #

```

```

392 class Num(Col):
393     "Summarize numeric columns."
394     def __init__(i, size,**kw):
395         super().__init__(**kw)
396         i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, size, False
397
398
399 def add(i,x: float ,inc=1):
400     "Reservoir sampler. If '_all' is full, sometimes replace an item at random."
401     if x != "?":
402         i.n += inc
403         i.lo = min(x,i.lo)
404         i.hi = max(x,i.hi)
405         if len(i._all) < i.max : i.ok=False; i._all += [x]
406         elif x() < i.max/i.n: i.ok=False; i._all[anywhere(i._all)] = x
407     return x
408
409 def all(i):
410     "Return '_all'.sorted."
411     if not i.ok: i.ok=True; i._all.sort()
412     return i._all
413
414 def dist1(i,x,y):
415     if x=="?": y=i.norm(y); x=(1 if y<.5 else 0)
416     elif y=="?": x=i.norm(x); y=(1 if x<.5 else 0)
417     else : x,y = i.norm(x), i.norm(y)
418     return abs(x-y)
419
420 def div(i):
421     """Report the diversity of this distribution (using standard deviation).
422     &pm;2, 2.56, 3 &sigma; is 66,90,95% of the mass. 2&sigma;. So one
423     standard deviation is (90-10)th divide by 2.4 times &sigma;."""
424     return (i.per(.9) - i.per(.1)) / 2.56
425
426 def merge(i,j):
427     "Return two 'Num's."
428     k = Num(i.max, at=i.at, txt=i.txt)
429     for x in i._all: k.add(x)
430     for x in j._all: k.add(x)
431     return k
432
433 def mid(i):
434     "Return central tendency of this distribution (using median)."
435     return i.per(.5)
436
437 def norm(i,x):
438     "Normalize 'x' to the range 0..1."
439     return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)
440
441 def per(i,p:float=.5) -> float:
442     "Return the p-th ranked item."
443     return per(i.all(), p)
444
445 def prep(i,x):
446     "Return 'x' as a float."
447     return x if x=="?" else float(x)
448
449 def spans(i,j, bins, out):
450     """Divide the whole space 'lo' to 'hi' into, say, 'xsmall'=16 bin,
451     then count the number of times we the bin on other side.
452     Then merge similar adjacent bins."""
453     lo = min(i.lo, j.lo)
454     hi = max(i.hi, j.hi)
455     gap = (hi-lo) / bins
456     xys = [(x,"this",1) for x in i._all] + [
457         (x,"that",1) for x in j._all]
458     one = Span(i.lo,lo)
459     all = [one]
460     for x,y,n in sorted(xys, key=first):
461         if one.hi - one.lo > gap:
462             one = Span(i, one.hi,x)
463             all += [one]
464         one.add(x,y,n)
465     all = merge(all)
466     all[0].lo = -big
467     all[-1].hi = big
468     if len(all) > 1: out += all
469 #
470 #
471 #
472 #
473 #
474 class Example(o):
475     def __init__(i,cells):
476         "One example stores a list of cells."
477         i.cells=cells
478     def __getitem__(i,k):
479         "Accessor."
480         return i.cells[k]
481
482 def dist(i,j, sample):
483     "Separation of two examples."
484     cols, p = sample.the.p
485     d = sum(col.dist(i[col.at], j[col.at])**p for col in cols)
486     return (d/len(cols)) ** (1/p)
487
488 def better(i,j, sample):
489     "Compare different goals."
490     n = len(cols)
491     for col in cols:
492         a,b = col.norm( i[col.at] ), col.norm( j[col.at] )
493         s1 == math.e**(col.w*(a-b)/n)
494         s2 == math.e**(col.w*(b-a)/n)
495     return s1/n < s2/n
496 #
497 #
498 #
499 #
500 #
501 class Explain(o):
502     """Split the data using random projections. Find the span that most
503     separates the data. Divide data on that span."""
504     def __init__(i,sample, top=None):
505         i.here, i.span, i.yes, i.no = sample, None, None, None
506         top = top or sample
507         enough = len(top.rows)**top.the.enough
508         if len(sample.rows) >= 2*enough:
509             left, right, *_ = sample.half(top)
510             spans = []
511             bins = 6/top.the.xsmall
512             [lcol.spans(rcol, bins, spans) for lcol,rcol in zip(left.x,right.x)]
513             if len(spans) > 0:
514                 i.span = Span.sort(spans)[0]
515                 yes, no = sample.clone(), sample.clone()
516                 [(yes if i.span.selects(row) else no).add(row) for row in sample.rows]
517                 if enough <= len(yes.rows) < len(sample.rows): i.yes= Explain(yes,top)
518                 if enough <= len(no.rows) < len(sample.rows): i.no = Explain(no, top)
519
520 def show(i,pre=""):
521     "Pretty print"
522     if not pre:
523         tmp= i.here.mid(i.here.y)
524         print(f"{'':40} : {len(i.here.rows):5} : {tmp}")
525     for (status,kid) in [(True,i.yes), (False, i.no)]:
526         if kid:
527             s=f"[pre]{lspan.show(status)}"

```

```

528         tmp= kid.here.mid(kid.here.y)
529         print(f"{'s':40} : {len(kid.here.rows):5} : {tmp}")
530         kid.show(pre + " |. ")
531 #
532 #
533 #
534 #
535 #
536 class Cluster(o):
537     "Tree with 'left','right' samples, broken at median between far points."
538     def __init__(i, sample, top=None):
539         i.left, i.right, i.x, i.y, i.c, i.mid = None,None,None,None,None,None
540         i.here = sample
541         top = top or sample
542         enough = len(top.rows)**top.the.enough
543         if len(sample.rows) >= 2*enough:
544             left, right, i.x, i.y, i.c, i.mid = sample.half(top)
545             if len(left.rows) < len(sample.rows):
546                 i.left = Cluster(left, top)
547                 i.right = Cluster(right,top)
548
549 def show(i,pre=""):
550     "pretty print"
551     s= f"{'pre':40} : {len(i.here.rows):5}"
552     print(f"{'s':s}" if i.left else f"{'s':s} : {i.here.mid(i.here.y)}")
553     for kid in [i.left,i.right]:
554         if kid: kid.show(pre + " |. ")
555 #
556 #
557 #
558 #
559 #
560 class Sample(o):
561     "Load, then manage, a set of examples."
562
563     def __init__(i, the, inits=[]):
564         """Samples hold 'rows', summarized in 'col'umns. The non-skipped columns
565         are stored in 'x,y' lists for independent and dependent columns. Also
566         stored is the 'klass' column and 'the' configuration options."""
567         i.the = the
568         i.rows, i.cols, i.x, i.y, i.klass = [], [], [], [],None
569         if str ==type(inits): [i.add(row, True) for row in file(inits)]
570         if list==type(inits): [i.add(row) for row in inits]
571
572 def add(i, a, raw=False):
573     """If we have no 'cols', this 'a' is the first row with the column names.
574     Otherwise 'a' is another row of data."""
575     def prep(a,c) : return c.prep(a[c.at]) if raw else a[c.at]
576     def is_num(x) : return x[0].isupper()
577     def is_skip(x) : return x[-1]=="-"
578     def is_klass(x) : return "!" in x
579     def is_goal(x) : return "+" in x or "-" in x or is_klass(x)
580     def col(at,txt) :
581         now = Num(i.the.Max,at=at,txt=txt) if is_num(txt) else Sym(at=at,txt=txt)
582         where= i.y if is_goal(txt) else i.x
583         if not is_skip(txt):
584             where += [now]
585             if is_klass(txt): i.klass = now
586         return now
587     #-----
588     if i.cols: i.rows += [Example([col.add(pre(a,col)) for col in i.cols])]
589     else: i.cols = [col(at,txt) for at,txt in enumerate(a)]
590
591 def clone(i,inits=[]):
592     "Generate a new 'Sample' with the same structure as this 'Sample'."
593     out = Sample(i.the)
594     out.add([col.txt for col in i.cols])
595     [out.add(x) for x in inits]
596     return out
597
598 def far(i,x,rows):
599     "Return something 'far' percent away from 'x' in 'rows'."
600     return per(sorted([(x.dist(y,i),y) for y in rows],key=first), i.the.far)
601
602 def half(i, top=None):
603     "Using two faraway points 'x,y' break data at median distance."
604     some= i.rows if len(i.rows)<i.the.Some else rnd.choices(i.rows,k=i.the.Some)
605     top= top or i
606     w = any(some)
607     _x= top.far(w, some)
608     c,y= top.far(x, some)
609     tmp= [row for _,row in sorted([(top.project(row,x,y,c), row)
610                                     for row in i.rows],key=first))]
611     mid= len(tmp)//2
612     return i.clone(tmp[:mid]), i.clone(tmp[mid:]), x, y, c, tmp[mid]
613
614 def mid(i,cols=None):
615     "Return a list of the mids of some columns."
616     return [col.mid() for col in (cols or i.all)]
617
618 def project(i,row,x,y,c):
619     "Find the distance of a 'row' on a line between 'x' and 'y'."
620     a = row.dist(x,i)
621     b = row.dist(y,i)
622     return (a**2 + c**2 - b**2) / (2*c)

```

```

623 #
624 #
625 #
626 #
627 #
628 #
629 #
630
631 class Demos:
632     "Possible start-up actions."
633     fails=0
634     "Number of errors; returned to operating system as our exit code"
635     def opt():
636         "show the config."
637         print(the)
638
639     def seed():
640         "seed"
641         assert .494 <= r() <= .495
642
643     def num():
644         "check 'Num'."
645         n = Num(512)
646         for _ in range(100): n.add(r())
647         assert .30 <= n.div() <= .31, "in range"
648
649     def sym():
650         "check 'Sym'."
651         s = Sym()
652         for x in "aaaabbc": s.add(x)
653         assert 1.37 <= s.div() <= 1.38, "entropy"
654         assert 'a' == s.mid(), "mode"
655
656     def rows():
657         "count rows in a file."
658         assert 399 == len([row for row in file(the.data)])
659
660     def sample():
661         "sampling"
662         s = Sample(the, the.data)
663         print(the.data, len(s.rows))
664         print(s.x[3], s.rows[-1])
665         assert 398 == len(s.rows), "length of rows"
666         assert 249 == s.x[-1].has['l'], "symbol counts"
667
668     def dist():
669         "distance between rows"
670         s = Sample(the, the.data)
671         assert .84 <= s.rows[1].dist(s.rows[-1],s) <= .842
672
673     def clone():
674         "cloning"
675         s = Sample(the, the.data)
676         s1 = s.clone(s.rows)
677         d1,d2 = s.x[0].__dict__, s1.x[0].__dict__
678         for k,v in d1.items():
679             print(d2[k],v)
680             assert d2[k] == v, "clone test"
681
682     def half():
683         "divide data in two"
684         s = Sample(the, the.data)
685         s1,s2,*_ = s.half()
686         print(s1.mid(s1.y))
687         print(s2.mid(s2.y))
688
689     def cluster():
690         "divide data in two"
691         s = Sample(the, the.data)
692         Cluster(s).show(); print("")
693
694     def xplain():
695         "divide data in two"
696         s = Sample(the, the.data);
697         Explain(s).show(); print("")
698
699 #-----
700 the = options(__doc__)
701 if __name__ == "__main__": demo(the.todo,Demos)
702
703 """
704 Example class
705 """

```