

```

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: return float(x)
165         except: return x.strip()
166
167 def demo(do,all):
168     "Maybe run a demo, if we want it, resetting random seed first."
169     todo = dir(all)
170     if do and do != "all":
171         todo = [x for x in dir(all) if x.startswith(do)]
172     for one in todo:
173         fun = all.__dict__.get(one,"")
174         if type(fun)==type(demo):
175             rnd.seed(the.seed)
176             doc = re.sub(r'\n\s+', "\n", fun.__doc__ or "")
177             try:
178                 fun()
179                 print("PASS:", doc)
180             except Exception as e:
181                 all.fails += 0
182                 if the.cautious: traceback.print_exc(); exit(1)
183                 else: print("FAIL:", doc, e)
184             exit(all.fails)
185
186 def file(f):
187     "Iterator. Returns one row at a time, as cells."
188     with open(f) as fp:
189         for line in fp:
190             line = re.sub(r'([\n\r\v\`"]#*)', '', line)
191             if line:
192                 yield [cell.strip() for cell in line.split(",")]
193
194 def first(a:list) -> Any:
195     "Return first item."
196     return a[0]
197
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, *_ = 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 r() -> float:
238     "Return random number 0..1"
239     return rnd.random()
240
241 def rn(x:float, n=3) -> float:
242     "Round a number to three decimals."
243     return round(x,n)
244
245 def rN(a:list, n=3) -> list:
246     "Round a list of numbers to three decimals."
247     return [rn(x,n=n) for x in a]
248
249 def second(a:list) -> Any:
250     "Return second item."
251     return a[1]

```

```

252 #
253 #
254 #
255 #
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 class Num(Col):
389     "Summarize numeric columns."
390     def __init__(i, size, **kw):
391         super().__init__(**kw)
392         i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, size, False
393
394
395 def add(i, x: float, inc=1):
396     "Reservoir sampler. If 'all' is full, sometimes replace an item at random."
397     if x != "?":
398         i.n += inc
399         i.lo = min(x, i.lo)
400         i.hi = max(x, i.hi)
401         if len(i._all) < i.max: i.ok=False; i._all += [x]
402         elif r() < i.max/i.n: i.ok=False; i._all[anywhere(i._all)] = x
403     return x
404
405 def all(i):
406     "Return 'all', sorted."
407     if not i.ok: i.ok=True; i._all.sort()
408     return i._all
409
410 def dist1(i, x, y):
411     if x=="?": y=i.norm(y); x=(1 if y<.5 else 0)
412     elif y=="?": x=i.norm(x); y=(1 if x<.5 else 0)
413     else: x, y = i.norm(x), i.norm(y)
414     return abs(x-y)
415
416 def div(i):
417     """Report the diversity of this distribution (using standard deviation).
418     &pm;2.256, 3 &sigma; is 66,90,95%, of the mass. 2&sigma;. So one
419     standard deviation is (90-10)th divide by 2.4 times &sigma;."""
420     return (i.per(.9) - i.per(.1)) / 2.56
421
422 def merge(i, j):
423     "Return two 'Num's."
424     k = Num(i.max, at=i.at, txt=i.txt)
425     for x in i._all: k.add(x)
426     for x in j._all: k.add(x)
427     return k
428
429 def mid(i):
430     "Return central tendency of this distribution (using median)."
431     return i.per(.5)
432
433 def norm(i, x):
434     "Normalize 'x' to the range 0..1."
435     return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)
436
437 def per(i, p:float=.5) -> float:
438     "Return the p-th ranked item."
439     a = i.all(); return a[ int(p*len(a)) ]
440
441 def prep(i, x):
442     "Return 'x' as a float."
443     return x if x=="?" else float(x)
444
445 def spans(i, j, bins, out):
446     """Divide the whole space 'lo' to 'hi' into, say, 'xsmall'=16 bin,
447     then count the number of times we the bin on other side.
448     Then merge similar adjacent bins."""
449     lo = min(i.lo, j.lo)
450     hi = max(i.hi, j.hi)
451     gap = (hi-lo) / bins
452     xys = [(x, "this", 1) for x in i._all] + [
453         (x, "that", 1) for x in j._all]
454     one = Span(i.lo, lo)
455     all = [one]
456     for x, y, n in sorted(xys, key=first):
457         if one.hi - one.lo > gap:
458             one = Span(i, one.hi, x)
459             all += [one]
460         one.add(x, y, n)
461     all = merge(all)
462     all[0].lo = -big
463     all[-1].hi = big
464     if len(all) > 1: out += all
465
466 #
467 # example
468 #
469
470 class Example(o):
471     def __init__(i, cells):
472         "One example stores a list of cells."
473         i.cells=cells
474     def __getitem__(i, k):
475         "Accessor."
476         return i.cells[k]
477
478 def dist(i, j, sample):
479     "Separation of two examples."
480     cols, p = sample, sample.the.p
481     d = sum(col.dist(i[col.at], j[col.at])**p for col in cols)
482     return (d/len(cols)) ** (1/p)
483
484 def better(i, j, sample):
485     "Compare different goals."
486     n = len(cols)
487     for col in cols:
488         a, b = col.norm( i[col.at] ), col.norm( j[col.at] )
489         s1 -= math.e**(col.w*(a-b)/n)
490         s2 -= math.e**(col.w*(b-a)/n)
491     return s1/n < s2/n
492
493 #
494 # explain
495 #
496
497 class Explain(o):
498     "Tree with 'yes'/'no' branches for samples that do/do not match a 'span'."
499     def __init__(i, here):
500         i.here, i.span, i.yes, i.no = here, None, None, None
501
502 def show(i, pre=""):
503     "Pretty print"
504     if not pre:
505         tmp = i.here.mid(i.here.y)
506         print(f"{pre:40}: {len(i.here.rows):5}: {tmp}")
507     for (status, kid) in [(True, i.yes), (False, i.no)]:
508         if kid:
509             s=f"{pre}{i.span.show(status)}"
510             tmp = kid.here.mid(kid.here.y)
511             print(f"{s:40}: {len(kid.here.rows):5}: {tmp}")
512             kid.show(pre + "|. ")
513
514 #
515 # cluster
516 #
517
518 class Cluster(o):
519     "Tree with 'left','right' samples, broken at median between far points."
520     def __init__(i, sample, top=None):
521         i.here, i.left, i.right, i.x, i.y, i.c, i.mid = sample, None, None, None, None, None, None
522         top = top or sample
523
524 if len(sample.rows) >= 2*(len(top.rows)**top.the.enough):
525     left, right, i.x, i.y, i.c, i.mid = sample.half(top)
526     if len(left.rows) < len(sample.rows):
527         i.left = Cluster(left, top)
528         i.right = Cluster(right, top)
529
530 def show(i, pre=""):
531     s=f"{pre:40}: {len(i.here.rows):5}"
532     print(f"{s} if i.left else f'{s} : {i.here.mid(i.here.y)}")
533     for kid in [i.left, i.right]:
534         if kid: kid.show(pre + "|. ")
535
536 #
537 # sample
538 #
539
540 class Sample(o):
541     "Load, then manage, a set of examples."
542
543     def __init__(i, the, inits=[]):
544         """Samples hold 'rows', summarized in 'col'umns. The non-skipped columns
545         are stored in 'x,y' lists for independent and dependent columns. Also
546         stored is the 'klass' column and 'the' configuration options."""
547         i.the = the
548         i.rows, i.cols, i.x, i.y, i.klass = [], [], [], [], None
549         if str == type(inits): [i.add(row, True) for row in file(inits)]
550         if list == type(inits): [i.add(row) for row in inits]
551
552 def add(i, a, raw=False):
553     """If we have no 'cols', this 'a' is the first row with the column names.
554     Otherwise 'a' is another row of data."""
555     def pre(a, c): return c.prep(a[c.at]) if raw else a[c.at]
556     def is_num(x): return x[0].isupper()
557     def is_skip(x): return x[-1]=="-"
558     def is_klass(x): return "!" in x
559     def is_goal(x): return "+" in x or "-" in x or is_klass(x)
560     def col(at, txt):
561         now = Num(i.the.Max, at=at, txt=txt)
562         where = i.y if is_goal(txt) else i.x
563         if not is_skip(txt):
564             where += [now]
565         if is_klass(txt): i.klass = now
566         return now
567     if i.cols: i.rows += [Example([col.add(pre(a, col)) for col in i.cols])]
568     else: i.cols = [col(at, txt) for at, txt in enumerate(a)]
569
570 def clone(i, inits=[]):
571     "Generate a new 'Sample' with the same structure as this 'Sample'."
572     out = Sample(i.the)
573     out.add([col.txt for col in i.cols])
574     [out.add(x) for x in inits]
575     return out
576
577 def cluster(i, top=None):
578     """Split the data using random projections. Find the span that most
579     separates the data. Divide data on that span."""
580     here = Cluster(i)
581     top = top or i
582     if len(i.rows) >= 2*(len(top.rows)**i.the.enough):
583         left, right, x, y, c, mid = i.half(top)
584         if len(left.rows) < len(i.rows):
585             here = Cluster(i, x, y, c, mid)
586             here.left = left.cluster(top)
587             here.right = right.cluster(top)
588     return here
589
590 def far(i, x, rows):
591     tmp = sorted([(x.dist(y, i), y) for y in (rows or i.rows)], key=first)
592     return tmp[ int(len(tmp)*i.the.far) ]
593
594 def half(i, top=None):
595     "Using two faraway points 'x,y' break data at median distance."
596     some=i.rows if len(i.rows)<i.the.Some else rnd.choices(i.rows, k=i.the.Some)
597     top = top or i
598     w = any(some)
599     _, x = top.far(w, some)
600     c, y = top.far(x, some)
601     tmp = [r for _, r in sorted([(top.proj(r, x, y, c), r)
602                                for r in i.rows], key=first))]
603     mid = len(tmp)//2
604     return i.clone(tmp[:mid]), i.clone(tmp[mid:]), x, y, c, tmp[mid]
605
606 def mid(i, cols=None):
607     "Return a list of the mids of some columns."
608     return [col.mid() for col in (cols or i.all)]
609
610 def proj(i, row, x, y, c):
611     "Find the distance of a 'row' on a line between 'x' and 'y'."
612     a = row.dist(x, i)
613     b = row.dist(y, i)
614     return (a**2 + c**2 - b**2) / (2*c)
615
616 def xplain(i, top=None):
617     """Split the data using random projections. Find the span that most
618     separates the data. Divide data on that span."""
619     here = Explain(i)
620     top = top or i
621     tiny = len(top.rows)**top.the.enough
622     if len(i.rows) >= 2*tiny:
623         left, right, _ = i.half(top)
624         spans = []
625         [icol.spans(rcol, 6/top.the.xsmall, spans) for lcol, rcol
626          in zip(left.x, right.x)]
627         if len(spans) > 0:
628             here.span = Span.sort(spans)[0]
629             yes, no = i.clone(), i.clone()
630             [(yes if here.span.selects(row) else no).add(row) for row in i.rows]
631             if tiny <= len(yes.rows) < len(i.rows): here.yes = yes.xplain(top=top)
632             if tiny <= len(no.rows) < len(i.rows): here.no = no.xplain(top=top)
633     return here

```

```

634 #
635 #
636 #
637 #
638 #
639 #
640 #
641
642 class Demos:
643     "Possible start-up actions."
644     fails=0
645     "Number of errors; returned to operating system as our exit code"
646     def opt():
647         "show the config."
648         print(the)
649
650     def seed():
651         "seed"
652         assert .494 <= r() <= .495
653
654     def num():
655         "check 'Num'."
656         n = Num(512)
657         for _ in range(100): n.add(r())
658         assert .30 <= n.div() <= .31, "in range"
659
660     def sym():
661         "check 'Sym'."
662         s = Sym()
663         for x in "aaaabbc": s.add(x)
664         assert 1.37 <= s.div() <= 1.38, "entropy"
665         assert 'a' == s.mid(), "mode"
666
667     def rows():
668         "count rows in a file."
669         assert 399 == len([row for row in file(the.data)])
670
671     def sample():
672         "sampling"
673         s = Sample(the, the.data)
674         print(the.data, len(s.rows))
675         print(s.x[3], s.rows[-1])
676         assert 398 == len(s.rows), "length of rows"
677         assert 249 == s.x[-1].has['l'], "symbol counts"
678
679     def dist():
680         "distance between rows"
681         s = Sample(the, the.data)
682         assert .84 <= s.rows[1].dist(s.rows[-1],s) <= .842
683
684     def clone():
685         "cloning"
686         s = Sample(the, the.data)
687         s1 = s.clone(s.rows)
688         d1,d2 = s.x[0].__dict__, s1.x[0].__dict__
689         for k,v in d1.items():
690             print(d2[k],v)
691             assert d2[k] == v, "clone test"
692
693     def half():
694         "divide data in two"
695         s = Sample(the, the.data)
696         s1,s2,*_ = s.half()
697         print(s1.mid(s1.y))
698         print(s2.mid(s2.y))
699
700     def cluster():
701         "divide data in two"
702         s = Sample(the, the.data)
703         Cluster(s).show(); print("")
704
705     def xplain():
706         "divide data in two"
707         s = Sample(the, the.data);
708         s.xplain().show(); print("")
709
710 #-----
711 the = options(__doc__)
712 if __name__ == "__main__": demo(the.todo,Demos)
713
714 """
715 Example class
716 """

```