```python
#!/usr/bin/env python3.10
# vim: ts=2 sw=2 sts=2 et :
#
#     dew    ~~~~~~~~~~~~~~~~~~~~~~~~\\W~~~~~~~~~~~~\|/~~
#               ~~~w/w~"~~,\ ' `:/,~(~ '"~~~~~~~"~o~\~/~w|/~
#                 ) ___/#\::'/ (O "==._____  O, (O  /'
#              O        |:::/{ }  |                  (o
#                       |:::\(\  \
#                       |:::\  \  `.
#            .          |::.   {\
#                  . (|::.      .                       .
#                     /:. _/ , |
#                      ) :_(:;     \
#                     `.:.  /:'  }         .
#                    \::.  :\/:'  /               +
#              ".:._::\____    /:'  /         .       .
#          .       \:    X `-_| _,\/'  _-'
#        `.  (   \: \,-,' + /`,_'" _--_._---"::._/
#           ,='.  \ `  ` ,' ,:'  '/`--",--"---._,/'7
#          /:+- - + +- : :- + + -:'  /(o-) \)
#       .  \/:/`-'  , \ ' ` ` ` ) : , /_  —=-'   .
#          .,  ,-==-. ,\ +#./'   \:.  / /
#          ,     './ \:. '.); \ _/ /
#    `         \  _| '"=:_::.'.);  \ __/ /
#   ((      .     (_:#::.:::. `-._  /:, /-.,_ `._,
#     ,              /;-._,-.____  ,------.__
#                       _      .                   .
#   .             .        +            .       .
# """
./sublime.py [OPTIONS]
(c)2022 Tim Menzies <timm@ieee.org>, BSD license
S.U.B.L.I.M.E. =
Sublime's unsupervised bifurcation: let's infer minimal explanations.

OPTIONS:

  −Max     max numbers to keep        : 512
  −Some    find 'far' in this many egs  : 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)

[![DOI](https://zenodo.org/badge/DOI/10.5281/zenodo.5912461.svg)](https://doi.org/10.5281/zenodo.5912461)
![](https://img.shields.io/badge/purpose−se−−ai−blueviolet)
![](https://img.shields.io/badge/language−python3−orange)
![](https://img.shields.io/badge/platform−osx,linux−pink)
<a href=https://github.com/timm/sublime/actions/workflows/main.yml><img
src=https://github.com/timm/sublime/actions/workflows/main.yml/badge.svg></a>

## 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

```
/sublime.py −c −s $RANDOM −t cluster
                              :  398
|..                           :  199
|.. |..                       :   99
|.. |.. |..                   :   49   Lbs−  Acc+  Mpg+
|.. |.. |.. |..               :   24  : [2255, 15.5, 30]
|.. |.. |.. |..               :   25  : [2575, 16.4, 30]
|.. |.. |..                   :   50
|.. |.. |.. |..               :   25  : [2110, 16.4, 30] <== best
|.. |.. |.. |..               :   25  : [2205, 16, 30]
|.. |..                       :  100
|.. |.. |..                   :   50
|.. |.. |.. |..               :   25  : [2234, 15.5, 30]
|.. |.. |.. |..               :   25  : [2278, 16.5, 30]
|.. |.. |..                   :   50
|.. |.. |.. |..               :   25  : [2220, 15.5, 30]
|.. |.. |.. |..               :   25  : [2320, 15.8, 30]
|..                           :  199
|.. |..                       :   99
|.. |.. |..                   :   49
|.. |.. |.. |..               :   24  : [2451, 16.5, 20]
|.. |.. |.. |..               :   25  : [3021, 15.5, 20]
|.. |.. |..                   :   50
|.. |.. |.. |..               :   25  : [3425, 17.6, 20]
|.. |.. |.. |..               :   25  : [3155, 16.7, 20]
|.. |..                       :  100
|.. |.. |..                   :   50
|.. |.. |.. |..               :   25  : [4141, 13.5, 10]
|.. |.. |.. |..               :   25  : [4054, 13.2, 20]
|.. |.. |..                   :   50
|.. |.. |.. |..               :   25  : [4425, 11, 10]
|.. |.. |.. |..               :   25  : [4129, 13, 10]
```

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

```
./sublime.py −c −s $RANDOM −t xplain
                                Lbs−  Acc+  Mgg+
                              :  398 : [2807, 15.5, 20]
198 <= Lbs < 454              :  167 : [3725, 14.5, 20]
|.. Modl < 72                 :   34 : [3609, 13,  20]
|.. Modl >= 72                :  133 : [3735, 14.9, 20]
|.. |.. Cylr < 8              :   56 : [3336, 17,  20]
|.. |.. |.. 77 <= Modl < 82   :   22 : [3410, 17.1, 20]
|.. |.. |.. Modl < 77 or Modl >= 82  :   34 : [3233, 17,  20]
|.. |.. Cylr >= 8             :   77 : [4129, 13.2, 20]
|.. |.. |.. Modl < 75         :   37 : [4274, 13,  10]
|.. |.. |.. Modl >= 75        :   40 : [3962, 13.5, 20]
|.. |.. |.. |.. Lbs >= 302    :   35 : [4054, 13.2, 20]
Lbs < 198 or Lbs >= 454       :  231 : [2290, 16,  30] <== best
```

## License

Redistribution and use in source and binary forms, with or without
modification, are permitted provided that the following conditions are met:
1. Redistributions of source code must retain the above copyright notice, this
   list of conditions and the following disclaimer.
2. Redistributions in binary form must reproduce the above copyright notice,
   this list of conditions and the following disclaimer in the documentation
   and/or other materials provided with the distribution.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND
ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR
PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR
CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL,
EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO,
PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR
PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF
LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING
NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS
SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
"""

import traceback, random, math, sys, re
from random import random as r
from typing import Any
```

```
151  #      ___              ___
152  #  /\  \            /\  \
153  #  \/\  \          /  \  \ ___
154  #     \   \        /  /\  \ ___
155  #      \   \      /  /  \  \ \_\  L \
156  #      /\   \    /  /    \  \ \___
157  #      \/___/    \/_/     \/_/  \/_/

158
159  def any(a: list) -> Any:
160      "Return a random item."
161      return a[anywhere(a)]
162
163  def anywhere(a: list) -> int:
164      "Return a random index of list 'a'."
165      return random.randint(0, len(a)-1)
166
167
168  big = sys.maxsize
169
170  def atom(x):
171      "Return a number or trimmed string."
172      x = x.strip()
173      if x == "True" :
174          return True
175      elif x == "False":
176          return False
177      else:
178          try:
179              return int(x)
180          except:
181              try:
182                  return float(x)
183              except:
184                  return x.strip()
185
186  def demo(do, all):
187      "Maybe run a demo, if we want it, resetting random seed first."
188      todo = dir(all)
189      if do and do != "all":
190          todo = [x for x in dir(all) if x.startswith(do)]
191      for one in todo:
192          fun = all.__dict__.get(one, "")
193          if type(fun) == type(demo):
194              random.seed(the.seed)
195              doc = re.sub(r'\n\s+', "\n", fun.__doc__ or "")
196              try:
197                  fun()
198                  print("PASS:", doc)
199              except Exception as e:
200                  all.fails += 0
201                  if the.cautious :
202                      traceback.print_exc()
203                      exit(1)
204                  else :
205                      print("FAIL:", doc, e)
206      exit(all.fails)
207
208  def file(f):
209      "Iterator. Returns one row at a time, as cells."
210      with open(f) as fp:
211          for line in fp:
212              line = re.sub(r'([\n\t\r"\' ]|#.*)', '', line)
213              if line:
214                  yield [atom(cell.strip()) for cell in line.split(",")]
215
216  def first(a: list) -> Any:
217      "Return first item."
218      return a[0]
219
220  def merge(b4: list) -> list:
221      "While we can find similar adjacent things, merge them."
222      j, n, now = -1, len(b4), []
223      while j < n-1:
224          j += 1
225          a = b4[j]
226          if j < n-2:
227              if merged := a.merge(b4[j+1]):
228                  a = merged
229                  j += 1  # we will continue, after missing one
230          now += [a]
231      # if 'now' is same size as 'b4', look for any other merges.
232      return b4 if len(now) == len(b4) else merge(now)
233
234  class o(object):
235      "Class that can pretty print its slots, with fast inits."
236      def __init__(i, **d): i.__dict__.update(**d)
237
238      def __repr__(i):
239          pre = i.__class__.__name__ if isinstance(i, o) else ""
240          return pre+str(
241              {k: v for k, v in sorted(i.__dict__.items()) if str(k)[0] != "_"})
242
243  def options(doc: str) -> o:
244      """Convert 'doc' to options dictionary using command line args.
245      Args canuse two 'shorthands': (1) boolean flags have no arguments (and mentioning
246      those on the command line means 'flip the default value'; (2) args need only
247      mention the first few of a key (e.g. -s is enough to select for -seed)."""
248      d = {}
249      for line in doc.splitlines():
250          if line and line.startswith("  -"):
251              # get 1st,last word on each line
252              key, *_, x = line.strip()[1:].split(" ")
253              for j, flag in enumerate(sys.argv):
254                  if flag and flag[0] == "-" and key.startswith(flag[1:]):
255                      x = "True" if x == "False" else(
256                          "False" if x == "True" else sys.argv[j+1])
257              d[key] = atom(x)
258      if d["help"]:
259          exit(print(re.sub(r'\n#.*', "", doc, flags=re.S)))
260      return o(**d)
261
262  def r() -> float:
263      "Return random number 0..1"
264      return random.random()
265
266  def rn(x: float, n=3) -> float:
267      "Round a number to three decimals."
268      return round(x, n)
269
270  def rN(a: list, n=3) -> list:
271      "Round a list of numbers to three decimals."
272      return [rn(x, n=n) for x in a]
273
274  def second(a: list) -> Any:
275      "Return second item."
276      return a[1]
```

```
277  #      ___
278  #  /\  \
279  #  \/\  \
280  #     \   \      ___      ___      ___      ___      ___      ___
281  #     /\   \    /\  \    / L \    \___      /\   \    /\   \    /\   \
282  #  \ \___\    /  \  \    \/_/. \    \/\  \    /  \  \    /  \  \    /  \  \
283  #      \/___/    \/_/\/_/ \/_/\/_/  \/_/  \/_/  \/_/  \/_/

284
285  #       ___      ___
286  #  /  __<    / '_ \    / '_ \    / _   _ \
287  #  \___/    \__, /    \___/    /_/ /_/
288  #          /_/

289
290  class Span(o):
291      """Given two 'Sample's and some 'x' range 'lo..hi',
292      a 'Span' holds often that range appears in each 'Sample'."""
293      def __init__(i, col, lo, hi, ys=None):
294          i.col, i.lo, i.hi, i.ys = col, lo, hi, ys or Sym()
295
296      def add(i, x: float, y: Any, inc=1) -> None:
297          "'y' is a label identifying, one 'Sample' or another."
298          i.lo = min(x, i.lo)
299          i.hi = max(x, i.hi)
300          i.ys.add(y, inc)
301
302      def merge(i, j):  # -> Span|None
303          "If the merged span is simpler, return that merge."
304          a, b, c = i.ys, j.ys, i.ys.merge(j.ys)
305          if (i.ys.n == 0 or j.ys.n == 0 or
306                  c.div()*.99 <= (a.n*a.div() + b.n*b.div())/(a.n + b.n)):
307              return Span(i.col, min(i.lo, j.lo), max(i.hi, j.hi), ys=c)
308
309      def selects(i, row: list) -> bool:
310          "True if the range accepts the row."
311          x = row[i.col.at]
312          return x == "?" or i.lo <= x and x < i.hi
313
314      def show(i, positive=True) -> None:
315          "Show the range."
316          txt = i.col.txt
317          if positive:
318              if i.lo == i.hi:
319                  return f"{txt} == {i.lo}"
320              elif i.lo == -big:
321                  return f"{txt} < {i.hi}"
322              elif i.hi == big:
323                  return f"{txt} >= {i.lo}"
324              else :
325                  return f"{i.lo} <= {txt} < {i.hi}"
326          else:
327              if i.lo == i.hi:
328                  return f"{txt} != {i.lo}"
329              elif i.lo == -big:
330                  return f"{txt} >= {i.hi}"
331              elif i.hi == big:
332                  return f"{txt} < {i.lo}"
333              else :
334                  return f"{txt} < {i.lo} or {txt} >= {i.hi}"
335
336      def support(i) -> float:
337          "Returns 0..1."
338          return i.ys.n / i.col.n
339
340      @staticmethod
341      def sort(spans : list) -> list:
342          "Good spans have large support and low diversity."
343          divs, supports = Num(), Num()
344          def sn(s): return supports.norm( s.support())
345          def dn(s): return divs.norm(    s.ys.div())
346          def f(s): return ((1 - sn(s))**2 + dn(s)**2)**.5/2**.5
347          for s in spans:
348              divs.add(    s.ys.div())
349              supports.add(s.support())
350          return sorted(spans, key=f)
351
352  #      ___
353  #  /  __  \
354  #  \  |_/_/    / _ \    / _ \
355  #      \___|    \_/    /_/

356
357  class Col(o):
358      "Summarize columns."
359      def __init__(i, at=0, txt=""):
360          i.n, i.at, i.txt, i.w = 0, at, txt, (-1 if "-" in txt else 1)
361
362      def dist(i, x: Any, y: Any) -> float:
363          return 1 if x == "?" and y == "?" else i.dist1(x, y)
364
365  #      ___      ___      ___
366  #  /  __<    / |_| \    / '_ \
367  #  \___/    \__, /    /_/_/_/
368  #          /_/

369
370  class Sym(Col):
371      "Summarize symbolic columns."
372      def __init__(i, **kw):
373          super().__init__(**kw)
374          i.has, i.mode, i.most = {}, None, 0
375
376      def add(i, x: str, inc: int = 1) -> str:
377          "Update symbol counts in 'has', updating 'mode' as we go."
378          if x != "?":
379              i.n += inc
380              tmp = i.has[x] = inc + i.has.get(x, 0)
381              if tmp > i.most:
382                  i.most, i.mode = tmp, x
383          return x
384
385      def dist(i, x: str, y: str) -> float:
386          "Distance between two symbols."
387          return 0 if x == y else 1
388
389      def div(i):
390          "Return diversity of this distribution (using entropy)."
391          def p(x): return x / (1E-31 + i.n)
392          return sum( -p(x)*math.log(p(x), 2) for x in i.has.values() )
393
394      def merge(i, j):
395          "Merge two 'Sym's."
396          k = Sym(at=i.at, txt=i.txt)
397          for x, n in i.has.items():
398              k.add(x, n)
399          for x, n in j.has.items():
400              k.add(x, n)
401          return k
402
403      def mid(i):
404          "Return central tendancy of this distribution (using mode)."
405          return i.mode
406
407      def spans(i, j, out):
408          """For each symbol in 'i' and 'j', count the
409      number of times we see it on either side."""
410          xys = [(x, "this", n) for x, n in i.has.items()] + [
411              (x, "that", n) for x, n in j.has.items()]
412          one, last = None, None
```

```python
      all = []
      for x, y, n in sorted(xys, key=first):
        if x != last:
          last = x
          one = Span(i, x, x)
          all += [one]
        one.add(x, y, n)
      if len(all) > 1 :
        out += all

#    _,  __,  ___
#   | \| | | |' \|
#   |_||\_,_||_|_|
#

class Num(Col):
    "Summarize numeric columns."
    def __init__(i, **kw):
      super().__init__(**kw)
      i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, the.Max, False

    def add(i, x: float , inc=1):
      "Reservoir sampler. If '_all' is full, sometimes replace an item at random."
      if x != "?":
        i.n += inc
        i.lo = min(x, i.lo)
        i.hi = max(x, i.hi)
        if len(i._all) < i.max    :
          i.ok = False
          i._all += [x]
        elif r() < i.max/i.n:
          i.ok = False
          i._all[anywhere(i._all)] = x
      return x

    def all(i):
      "Return '_all', sorted."
      if not i.ok:
        i.ok = True
        i._all.sort()
      return i._all

    def dist1(i, x, y):
      if x == "?":
        y = i.norm(y)
        x = (1 if y < .5 else 0)
      elif y == "?":
        x = i.norm(x)
        y = (1 if x < .5 else 0)
      else :
        x, y = i.norm(x), i.norm(y)
      return abs(x-y)

    def div(i):
      """Report the diversity of this distribution (using standard deviation).
      &pm;2, 2,56, 3 &sigma; is 66,90,95%, of the mass.  2&8&sigma;. So one
      standard deviation is (90-10)th divide by  2.4 times &sigma;."""
      return  (i.per(.9) - i.per(.1)) / 2.56

    def merge(i, j):
      "Return two 'Num's."
      k = Num(at=i.at, txt=i.txt)
      for x in i._all:
        k.add(x)
      for x in j._all:
        k.add(x)
      return k

    def mid(i):
      "Return central tendency of this distribution (using median)."
      return i.per(.5)

    def norm(i, x):
      "Normalize 'x' to the range 0..1."
      return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)

    def per(i, p: float = .5) -> float:
      "Return the p-th ranked item."
      a = i.all()
      return a[ int(p*len(a)) ]

    def spans(i, j, out):
      """Divide the whole space 'lo' to 'hi' into, say, 'xsmall'=16 bin,
      then count the number of times we the bin on other side.
      Then merge similar adjacent bins."""
      lo = min(i.lo, j.lo)
      hi = max(i.hi, j.hi)
      gap = (hi-lo)  / (6/the.xsmall)
      xys = [(x, "this", 1) for x in i._all] + [
             (x, "that", 1) for x in j._all]
      one = Span(i, lo, lo)
      all = [one]
      for x, y, n in sorted(xys, key=first):
        if one.hi - one.lo > gap:
          one = Span(i, one.hi, x)
          all += [one]
        one.add(x, y, n)
      all = merge(all)
      all[ 0].lo = -big
      all[-1].hi = big
      if len(all) > 1:
        out += all

#    ___       _             _
#   |_  \/ _  |' _) | _'  _  _ __
#   <_-)/\/\ |'._>  | <_,_| | | |'
#          |_|

class Explain(o):
    "Tree with 'yes','no' branches for samples that do/do not match a 'span'."
    def __init__(i, here):
      i.here, i.span, i.yes, i.no = here, None, None, None

    def show(i, pre=""):
      if not pre:
        tmp = i.here.mid(i.here.y)
        print(f"{'':40} : {len(i.here.rows):5} : {tmp}")
      if i.yes:
        s = f"{pre}{i.span.show(True)}"
        tmp = i.yes.here.mid(i.yes.here.y)
        print(f"{s:40} : {len(i.yes.here.rows):5} : {tmp}")
        i.yes.show(pre + "|.. ")
      if i.no:
        s = f"{pre}{i.span.show(False)}"
        tmp = i.no.here.mid(i.no.here.y)
        print(f"{s:40} : {len(i.no.here.rows):5} : {tmp}")
        i.no.show(pre + "|.. ")

#    __  _        _    _
#   /  ||  | _,_ <-_ _|' _   _ _
#   \_ ||_ |_,_| _> |_<-' | |`-'
#

class Cluster(o):
    "Tree with 'left','right' samples, broken at median between far points."
    def __init__(i, here, x=None, y=None, c=None, mid=None):
      i.here, i.x, i.y, i.c, i.mid, i.left, i.right = here, x, y, c, mid, None, None
```

```python
    def show(i, pre=""):
      s = f"{pre:40} : {len(i.here.rows):5}"
      print(f"{s}" if i.left else f"{s} : {i.here.mid(i.here.y)}")
      for kid in [i.left, i.right]:
        if kid:
          kid.show(pre + "|.. ")

#          _            _
#    __   _       _    _|' _   _
#   /  \ |'_) |\/|'|_)| <-' | |`-'
#   \__/ |'_| |  | |_ |_ _| | |`-'
#                  |_|

class Sample(o):
    "Load, then manage, a set of examples."
    def __init__(i, inits=[]):
      i.rows, i.cols, i.x, i.y, i.klass = [], [], [], [], None
      if str == type(inits):
        [i.add(row) for row in file(inits)]
      if list == type(inits):
        [i.add(row) for row in inits]

    def add(i, a):
      def col(at, txt):
        what = Num if txt[0].isupper() else Sym
        now = what(at=at, txt=txt)
        where = i.y if "+" in txt or "-" in txt or "!" in txt else i.x
        if txt[-1] != ":":
          where += [now]
          if "!" in txt:
            i.klass = now
        return now
      # -----------
      if i.cols:
        i.rows += [[col.add(a[col.at]) for col in i.cols]]
      else:
        i.cols = [col(at, txt) for at, txt in enumerate(a)]

    def clone(i, inits=[]):
      out = Sample()
      out.add([col.txt for col in i.cols])
      [out.add(x) for x in inits]
      return out

    def cluster(i, top=None):
      """Split the data using random projections. Find the span that most
      separates the data. Divide data on that span."""
      here = Cluster(i)
      top = top or i
      if len(i.rows) >= 2*(len(top.rows)**the.enough):
        left, right, x, y, c, mid = i.half(top)
        if len(left.rows) < len(i.rows):
          here = Cluster(i, x, y, c, mid)
          here.left = left.cluster(top)
          here.right = right.cluster(top)
      return here

    def dist(i, x, y):
      d = sum( col.dist(x[col.at], y[col.at])**the.p for col in i.x )
      return (d/len(i.x)) ** (1/the.p)

    def div(i, cols=None):
      return [col.div() for col in (cols or i.all)]

    def far(i, x, rows=None):
      tmp = sorted([(i.dist(x, y), y) for y in (rows or i.rows)], key=first)
      return tmp[ int(len(tmp)*the.far) ]

    def half(i, top=None):
      "Using two faraway points 'x,y' break data at median distance."
      some = i.rows if len(
             i.rows) < the.Some else random.choices(i.rows, k=the.Some)
      top = top or i
      w = any(some)
      _, x = top.far(w, some)
      c, y = top.far(x, some)
      tmp = [r for _, r in sorted([(top.proj(r, x, y, c), r)
                            for r in i.rows], key=first)]
      mid = len(tmp)//2
      return i.clone(tmp[:mid]), i.clone(tmp[mid:]), x, y, c, tmp[mid]

    def mid(i, cols=None):
      return [col.mid() for col in (cols or i.all)]

    def proj(i, row, x, y, c):
      "Find the distance of a 'row' on a line between 'x' and 'y'."
      a = i.dist(row, x)
      b = i.dist(row, y)
      return (a**2 + c**2 - b**2) / (2*c)

    def xplain(i, top=None):
      """Split the data using random projections. Find the span that most
      separates the data. Divide data on that span."""
      here = Explain(i)
      top = top or i
      tiny = len(top.rows)**the.enough
      if len(i.rows) >= 2*tiny:
        left, right, *_ = i.half(top)
        spans = []
        [lcol.spans(rcol, spans) for lcol, rcol in zip(left.x, right.x)]
        if len(spans) > 0:
          here.span = Span.sort(spans)[0]
          yes, no = i.clone(), i.clone()
          [(yes if here.span.selects(row) else no).add(row) for row in i.rows]
          if tiny <= len(yes.rows) < len(i.rows):
            here.yes = yes.xplain(top=top)
          if tiny <= len(no.rows ) < len(i.rows):
            here.no = no.xplain(top=top)
      return here
```

```
657  #      _
658  #     /\ \
659  #     \_\ \         ___      ___      ___       ___
660  #     /'_` \      /'___\   /'___\ '\    /'___\  /'__`\
661  #    /\ \L\ \    /\ \__/  /\ \__/   \  /\ \L\ \ \/\ \  '
662  #    \ \___,_\   \ \____\ \ \____\   \ \ \____/\ \___,_\
663  #     \/__,_ /    \/____/  \/____/    \ \/___/  \/__,_ /

664
665  class Demos:
666    "Possible start-up actions."
667    fails = 0
668
669    def opt():
670      "show the config."
671      [print(f"{k:>10} = {v}") for k, v in the.__dict__.items()]
672
673    def seed():
674      "seed"
675      assert .494 <= r() <= .495
676
677    def num():
678      "check 'Num'."
679      n = Num()
680      for _ in range(100):
681        n.add(r())
682      assert .30 <= n.div() <= .31, "in range"
683
684    def sym():
685      "check 'Sym'."
686      s = Sym()
687      for x in "aaaabbc":
688        s.add(x)
689      assert 1.37 <= s.div() <= 1.38, "entropy"
690      assert 'a' == s.mid(), "mode"
691
692    def rows():
693      "count rows in a file."
694      assert 399 == len([row for row in file(the.data)])
695
696    def sample():
697      "sampling."
698      s = Sample(the.data)
699      assert 398 == len(s.rows), "length of rows"
700      assert 249 == s.x[-1].has[1], "symbol counts"
701
702    def dist():
703      "distance between rows"
704      s = Sample(the.data)
705      assert .84 <= s.dist(s.rows[1], s.rows[-1]) <= .842
706
707    def far():
708      "distant items"
709      s = Sample(the.data)
710      for _ in range(32):
711        a, _ = s.far(any(s.rows))
712        assert a > .5, "large?"
713
714    def clone():
715      "cloning"
716      s = Sample(the.data)
717      s1 = s.clone(s.rows)
718      d1, d2 = s.x[0].__dict__, s1.x[0].__dict__
719      for k, v in d1.items():
720        assert d2[k] == v, "clone test"
721
722    def half():
723      "divide data in two"
724      s = Sample(the.data)
725      s1, s2, *_ = s.half()
726      print(s1.mid(s1.y))
727      print(s2.mid(s2.y))
728
729    def cluster():
730      "divide data in two"
731      s = Sample(the.data)
732      s.cluster().show()
733      print("")
734
735    def xplain():
736      "divide data in two"
737      s = Sample(the.data)
738      s.xplain().show()
739      print("")
740
741
742  # ----------------------------------------------------
743  the = options(__doc__)
744  if __name__ == "__main__":
745    demo(the.todo, Demos)
746
747  """
748  all config local to Sample
749  Example class
750  """
```