

```

104 ## License
105
106 This is free and unencumbered software released into the public
107 domain.
108
109 Anyone is free to copy, modify, publish, use, compile, sell, or
110 distribute this software, either in source code form or as a compiled
111 binary, for any purpose, commercial or non-commercial, and by any
112 means.
113
114 In jurisdictions that recognize copyright laws, the author or authors
115 of this software dedicate any and all copyright interest in the
116 software to the public domain. We make this dedication for the
117 benefit of the public at large and to the detriment of our heirs
118 and successors. We intend this dedication to be an overt act of
119 relinquishment in perpetuity of all present and future rights to
120 this software under copyright law.
121
122 THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY
123 EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE
124 MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND
125 IN NO EVENT SHALL THE AUTHORS BE LIABLE FOR ANY OTHER
126 LIABILITY, WHETHER IN AN ACTION OF CONTRACT OR OTHERWISE,
127 ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR
128 FOR OTHER DEALINGS IN THE SOFTWARE.
129
130 For more information, please refer to <http://unlicense.org/>
131
132 ***

```

```
152
153 import traceback, random, math, sys, re
154 from random import random as r
155 from typing import Any
```

```

156 #
157 #
158 #
159 #
160 #
161 #
162 #
163
164 def any(a:list) -> Any:
165     "Return a random item."
166     return a[anywhere(a)]
167
168 def anywhere(a:list) -> int:
169     "Return a random index of list 'a'."
170     return random.randint(0, len(a)-1)
171
172 big = sys.maxsize
173
174 def atom(x):
175     "Return a number or trimmed string."
176     x=x.strip()
177     if x=="True": return True
178     elif x=="False": return False
179     else:
180         try: return int(x)
181         except:
182             try: return float(x)
183             except: return x.strip()
184
185 def demo(do,all):
186     "Maybe run a demo, if we want it, resetting random seed first."
187     todo = dir(all)
188     if do and do != "all":
189         todo = [x for x in dir(all) if x.startswith(do)]
190     for one in todo:
191         fun = all.__dict__.get(one,"")
192         if type(fun)==type(demo):
193             random.seed(the.seed)
194             doc = re.sub(r'\n\s+', "\n", fun.__doc__ or "")
195             try:
196                 fun()
197                 print("PASS:", doc)
198             except Exception as e:
199                 all.fails += 0
200                 if the.cautious: traceback.print_exc(); exit(1)
201                 else: print("FAIL:", doc, e)
202             exit(all.fails)
203
204 def file(f):
205     "Iterator. Returns one row at a time, as cells."
206     with open(f) as fp:
207         for line in fp:
208             line = re.sub(r'(\n|\r|\v)|#.*', '', line)
209             if line:
210                 yield [atom(cell.strip()) for cell in line.split(",")]
211
212 def first(a:list) -> Any:
213     "Return first item."
214     return a[0]
215
216 def merge(b4:list) -> list:
217     "While we can find similar adjacent things, merge them."
218     j,n,now = -1,len(b4),[]
219     while j < n-1:
220         j += 1
221         a = b4[j]
222         if j < n-2:
223             if merged := a.merge(b4[j+1]):
224                 a = merged
225                 j += 1 # we will continue, after missing one
226             now += [a]
227         # if 'now' is same size as 'b4', look for any other merges.
228     return b4 if len(now)==len(b4) else merge(now)
229
230 class o(object):
231     "Class that can pretty print its slots, with fast inits."
232     def __init__(i, **d): i.__dict__.update(**d)
233     def __repr__(i):
234         pre = i.__class__.__name__ if isinstance(i,o) else ""
235         return pre+str(
236             {k: v for k, v in sorted(i.__dict__.items()) if str(k)[0] != "_"})
237
238 def options(doc:str) -> o:
239     """Convert 'doc' to options dictionary using command line args.
240     Args cause two 'shorthands': (1) boolean flags have no arguments (and mentioning
241     those on the command line means 'flip the default value'; (2) args need only
242     mention the first few of a key (e.g. -s is enough to select for -seed)."""
243     d={}
244     for line in doc.splitlines():
245         if line and line.startswith(" -"):
246             key, _, x = line.strip()[1:].split("#") # get 1st,last word on each line
247             for j,flag in enumerate(sys.argv):
248                 if flag and flag[0]=="-" and key.startswith(flag[1:]):
249                     x= "True" if x=="False" else "False" if x=="True" else sys.argv[j+1])
250             d[key] = atom(x)
251     if d["help"]: exit(print(re.sub(r'\n#.*', "", doc, flags=re.S)))
252     return o(**d)
253
254 def r() -> float:
255     "Return random number 0..1"
256     return random.random()
257
258 def rn(x:float, n=3) -> float:
259     "Round a number to three decimals."
260     return round(x,n)
261
262 def rN(a:list, n=3) -> list:
263     "Round a list of numbers to three decimals."
264     return [rn(x,n=n) for x in a]
265
266 def second(a:list) -> Any:
267     "Return second item."
268     return a[1]

```

```

269 #
270 #
271 #
272 #
273 #
274 #
275 #
276
277 #
278 #
279 #
280 #
281 #
282
283 class Span(o):
284     """Given two 'Sample's and some 'x' range 'lo..hi'.
285     a 'Span' holds often that range appears in each 'Sample'."""
286     def __init__(i,col, lo, hi, ys=None,):
287         i.col, i.lo, i.hi, i.ys = col, lo, hi, ys or Sym()
288
289 def add(i, x:float, y:Any, inc=1) -> None:
290     "y' is a label identifying one 'Sample' or another."
291     i.lo = min(x, i.lo)
292     i.hi = max(x, i.hi)
293     i.ys.add(y,inc)
294
295 def merge(i, j): # -> Span|None
296     "If the merged span is simpler, return that merge."
297     a, b, c = i.ys, j.ys, i.ys.merge(j.ys)
298     if (i.ys.n==0 or j.ys.n==0 or
299         c.div()*.99 <= (a.n*a.div() + b.n*b.div())/(a.n + b.n)):
300         return Span(i.col, min(i.lo,j.lo),max(i.hi,j.hi), ys=c)
301
302 def selects(i,row:list) -> bool:
303     "True if the range accepts the row."
304     x = row[i.col.at]; return x=="?" or i.lo<=x and x<i.hi
305
306 def show(i, positive=True) -> None:
307     "Show the range."
308     txt = i.col.txt
309     if positive:
310         if i.lo == i.hi: return f"[txt] == {i.lo}"
311         elif i.lo == -big: return f"[txt] < {i.hi}"
312         elif i.hi == big: return f"[txt] >= {i.lo}"
313         else:
314             return f"[i.lo] <= [txt] < {i.hi}"
315     else:
316         if i.lo == i.hi: return f"[txt] != {i.lo}"
317         elif i.lo == -big: return f"[txt] >= {i.hi}"
318         elif i.hi == big: return f"[txt] < {i.lo}"
319         else:
320             return f"[txt] < {i.lo} or [txt] >= {i.hi}"
321
322 def support(i) -> float:
323     "Returns 0..1."
324     return i.ys.n / i.col.n
325
326 @staticmethod
327 def sort(spans : list) -> list:
328     "Good spans have large support and low diversity."
329     divs, supports = Num(), Num()
330     sn = lambda s: supports.norm( s.support())
331     dn = lambda s: divs.norm( s.ys.div())
332     f = lambda s: ((1 - sn(s))*2 + dn(s)**2)**.5/2***.5
333     for s in spans:
334         divs.add( s.ys.div())
335         supports.add(s.support())
336     return sorted(spans, key=f)
337
338 #
339 #
340 #
341
342 class Col(o):
343     "Summarize columns."
344     def __init__(i,at=0,txt=""):
345         i.n,i.at,i.txt,i.w=0,at,txt,(-1 if "-" in txt else 1)
346
347 def dist(i,x:Any, y:Any) -> float:
348     return 1 if x=="?" and y=="?" else i.dist1(x,y)
349
350 #
351 #
352 #
353
354 class Sym(Col):
355     "Summarize symbolic columns."
356     def __init__(i,**kw):
357         super().__init__(**kw)
358         i.has, i.mode, i.mode = {}, None, 0
359
360 def add(i, x:str, inc:int=1) -> str:
361     "Update symbol counts in 'has', updating 'mode' as we go."
362     if x != " ":
363         i.n += inc
364         tmp = i.has[x] = inc + i.has.get(x,0)
365         if tmp > i.mode: i.mode, i.mode = tmp, x
366     return x
367
368 def dist(i,x:str, y:str) ->float:
369     "Distance between two symbols."
370     return 0 if x==y else 1
371
372 def div(i):
373     "Return diversity of this distribution (using entropy)."
374     p = lambda x: x / (1E-31 + i.n)
375     return sum( -p(x)*math.log(p(x),2) for x in i.has.values() )
376
377 def merge(i,j):
378     "Merge two 'Sym's."
379     k = Sym(at=i.at, txt=i.txt)
380     for x,n in i.has.items(): k.add(x,n)
381     for x,n in j.has.items(): k.add(x,n)
382     return k
383
384 def mid(i):
385     "Return central tendency of this distribution (using mode)."
386     return i.mode
387
388 def spans(i,j, out):
389     """For each symbol in 'i' and 'j', count the
390     number of times we see it on either side."""
391     xys = [(x,"this",n) for x,n in i.has.items()] + [
392         (x,"that",n) for x,n in j.has.items()]
393     one, last = None,None
394     all = []
395     for x,y,n in sorted(xys, key=first):
396         if x != last:
397             last = x
398             one = Span(i, x,x)
399             all += [one]
400             one.add(x,y,n)
401         if len(all) > 1 : out += all

```

```

401 #
402 #
403 #
404 #
405 class Num(Col):
406     "Summarize numeric columns."
407     def __init__(i,**kw):
408         super().__init__(**kw)
409         i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, the.Max, False
410
411     def add(i,x: float ,inc=1):
412         "Reservoir sampler. If '_all' is full, sometimes replace an item at random."
413         if x != "":
414             i.n += inc
415             i.lo = min(x,i.lo)
416             i.hi = max(x,i.hi)
417             if len(i._all) < i.max : i.ok=False; i._all += [x]
418             elif r() < i.max/i.n: i.ok=False; i._all[anywhere(i._all)] = x
419         return x
420
421     def all(i):
422         "Return '_all', sorted."
423         if not i.ok: i.ok=True; i._all.sort()
424         return i._all
425
426     def dist1(i,x,y):
427         if x=="?": y=i.norm(y); x=(1 if y<.5 else 0)
428         elif y=="?": x=i.norm(x); y=(1 if x<.5 else 0)
429         else : x,y = i.norm(x), i.norm(y)
430         return abs(x-y)
431
432     def div(i):
433         """"Report the diversity of this distribution (using standard deviation).
434         &pm;2, 2.56, 3 &sigma; is 66,90,95%, of the mass. 2&sigma;. So one
435         standard deviation is (90-10)th divide by 2.4 times &sigma;."""
436         return (i.per(.9) - i.per(.1)) / 2.56
437
438     def merge(i,j):
439         "Return two 'Num's."
440         k = Num(at=i.at, txt=i.txt)
441         for x in i._all: k.add(x)
442         for x in j._all: k.add(x)
443         return k
444
445     def mid(i):
446         "Return central tendency of this distribution (using median)."
447         return i.per(.5)
448
449     def norm(i,x):
450         "Normalize 'x' to the range 0..1."
451         return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)
452
453     def per(i,p:float=.5) -> float:
454         "Return the p-th ranked item."
455         a = i.all(); return a[ int(p*len(a)) ]
456
457     def spans(i,j, out):
458         """"Divide the whole space 'lo' to 'hi' into, say, 'xsmall'=16 bin,
459         then count the number of times we hit the bin on other side.
460         Then merge similar adjacent bins.""
461         lo = min(i.lo, j.lo)
462         hi = max(i.hi, j.hi)
463         gap = (hi-lo) / (6/the.xsmall)
464         xys = [(x,"this",1) for x in i._all] + [
465             (x,"that",1) for x in j._all]
466         one = Span(i.lo,lo)
467         all = [one]
468         for x,y,n in sorted(xys, key=first):
469             if one.hi - one.lo > gap:
470                 one = Span(i, one.hi,x)
471                 all += [one]
472                 one.add(x,y,n)
473         all = merge(all)
474         all[0].lo = -big
475         all[-1].hi = big
476         if len(all) > 1: out += all
477 #
478 #
479 #
480 #
481 #
482 class Explain(o):
483     "Tree with 'yes','no' branches for samples that do/do not match a 'span'."
484     def __init__(i,here):
485         i.here, i.span, i.yes, i.no = here, None, None, None
486
487     def show(i,pre=""):
488         if not pre:
489             tmp = i.here.mid(i.here.y)
490             print(f"[{pre:40}]: {len(i.here.rows):5} : {tmp}")
491         if i.yes:
492             s=f"[pre]{i.span.show(True)}"
493             tmp = i.yes.here.mid(i.yes.here.y)
494             print(f"[{s:40}]: {len(i.yes.here.rows):5} : {tmp}")
495             i.yes.show(pre + "[. ")
496         if i.no:
497             s=f"[pre]{i.span.show(False)}"
498             tmp = i.no.here.mid(i.no.here.y)
499             print(f"[{s:40}]: {len(i.no.here.rows):5} : {tmp}")
500             i.no.show(pre + "[. ")
501 #
502 #
503 #
504 #
505 #
506 #
507 #
508 class Cluster(o):
509     "Tree with 'left','right' samples, broken at median between far points."
510     def __init__(i,here,x=None,y=None,c=None,mid=None):
511         i.here,i.x,i.y,i.c,i.mid,i.left,i.right = here,x,y,c,mid,None,None
512
513     def show(i,pre=""):
514         s = f"[pre:40]: {len(i.here.rows):5}"
515         print(f"[{s}] if i.left else f"[{s}]: {i.here.mid(i.here.y)}")
516         for kid in [i.left,i.right]:
517             if kid: kid.show(pre + "[. ")

```

```

518 #
519 #
520 #
521 #
522 #
523 #
524 class Sample(o):
525     "Load, then manage, a set of examples."
526     def __init__(i,init=[]):
527         i.rows, i.cols, i.x, i.y, i.klass = [], [], [], [], None
528         if str == type(inits): [i.add(row) for row in file(inits)]
529         if list == type(inits): [i.add(row) for row in inits]
530
531     def add(i,a):
532         def col(at,txt):
533             what = Num if txt[0].isupper() else Sym
534             now = what(at=at, txt=txt)
535             where = i.y if "x" in txt or "-" in txt or "!" in txt else i.x
536             if txt[-1] != " ":
537                 where += [now]
538             if "!" in txt: i.klass = now
539             return now
540         #-----
541         if i.cols: i.rows += [[col.add(a[col.at]) for col in i.cols]]
542         else: i.cols = [col(at,txt) for at,txt in enumerate(a)]
543
544     def clone(i,init=[]):
545         out = Sample()
546         out.add([col.txt for col in i.cols])
547         [out.add(x) for x in inits]
548         return out
549
550     def cluster(i,top=None):
551         """"Split the data using random projections. Find the span that most
552         separates the data. Divide data on that span.""
553         here = Cluster(i)
554         top = top or i
555         if len(i.rows) >= 2*(len(top.rows)**the.enough):
556             left,right,x,y,c,mid = i.half(top)
557             if len(left.rows) < len(i.rows):
558                 here = Cluster(i,x,y,c,mid)
559             here.left = left.cluster(top)
560             here.right = right.cluster(top)
561         return here
562
563     def dist(i,x,y):
564         d = sum( col.dist(x[col.at], y[col.at])**the.p for col in i.x )
565         return (d/len(i.x)) ** (1/the.p)
566
567     def div(i,cols=None):
568         return [col.div() for col in (cols or i.all)]
569
570     def far(i, x, rows=None):
571         tmp = sorted([(i.dist(x,y),y) for y in (rows or i.rows)],key=first)
572         return tmp[ int(len(tmp)*the.far) ]
573
574     def half(i, top=None):
575         "Using two faraway points 'x,y' break data at median distance."
576         some= i.rows if len(i.rows)<the.Some else random.choices(i.rows, k=the.Some)
577         top = top or i
578         w = any(some)
579         _,x= top.far(w, some)
580         c,y= top.far(x, some)
581         tmp = [r for _,r in sorted([(top.proj(r,x,y,c),r)
582                                     for r in i.rows],key=first))]
583         mid= len(tmp)//2
584         return i.clone(tmp[:mid]), i.clone(tmp[mid:]), x, y, c, tmp[mid]
585
586     def mid(i,cols=None):
587         return [col.mid() for col in (cols or i.all)]
588
589     def proj(i,row,x,y,c):
590         "Find the distance of a 'row' on a line between 'x' and 'y'."
591         a = i.dist(row,x)
592         b = i.dist(row,y)
593         return (a**2 + c**2 - b**2) / (2*c)
594
595     def xplain(i,top=None):
596         """"Split the data using random projections. Find the span that most
597         separates the data. Divide data on that span.""
598         here = Explain(i)
599         top = top or i
600         tiny = len(top.rows)**the.enough
601         if len(i.rows) >= 2*tiny:
602             left, right, *_ = i.half(top)
603             spans = []
604             [icol.spans(rcol,spans) for lcol,rcol in zip(left.x, right.x)]
605             if len(spans) > 0:
606                 here.span = Span.sort(spans)[0]
607                 yes, no = i.clone(), i.clone()
608                 [yes if here.span.selects(row) else no].add(row) for row in i.rows
609                 if tiny <= len(yes.rows) < len(i.rows): here.yes = yes.xplain(top=top)
610                 if tiny <= len(no.rows ) < len(i.rows): here.no = no.xplain(top=top)
611             return here
612
613

```

```

614 #
615 #
616 #
617 #
618 #
619 #
620 #
621
622 class Demos:
623     "Possible start-up actions."
624     fails=0
625     def opt():
626         "show the config"
627         [print(f"{k}>10}={v}") for k,v in the.__dict__.items()]
628
629     def seed():
630         "seed"
631         assert .494 <= r() <= .495
632
633     def num():
634         "check 'Num'."
635         n = Num()
636         for _ in range(100): n.add(r())
637         assert .30 <= n.div() <= .31, "in range"
638
639     def sym():
640         "check 'Sym'."
641         s = Sym()
642         for x in "aaaabbc": s.add(x)
643         assert 1.37 <= s.div() <= 1.38, "entropy"
644         assert 'a' == s.mid(), "mode"
645
646     def rows():
647         "count rows in a file."
648         assert 399 == len([row for row in file(the.data)])
649
650     def sample():
651         "sampling"
652         s = Sample(the.data)
653         assert 398 == len(s.rows), "length of rows"
654         assert 249 == s.x[-1].has[1], "symbol counts"
655
656     def dist():
657         "distance between rows"
658         s = Sample(the.data)
659         assert .84 <= s.dist(s.rows[1], s.rows[-1]) <= .842
660
661     def far():
662         "distant items"
663         s = Sample(the.data)
664         for _ in range(32):
665             a,_ = s.far(any(s.rows))
666             assert a>.5, "large?"
667
668     def clone():
669         "cloning"
670         s = Sample(the.data)
671         s1 = s.clone(s.rows)
672         d1,d2 = s.x[0].__dict__, s1.x[0].__dict__
673         for k,v in d1.items():
674             assert d2[k] == v, "clone test"
675
676     def half():
677         "divide data in two"
678         s = Sample(the.data); s1,s2,*_ = s.half()
679         print(s1.mid(s1.y))
680         print(s2.mid(s2.y))
681
682     def cluster():
683         "divide data in two"
684         s = Sample(the.data)
685         s.cluster().show(); print("")
686
687     def xplain():
688         "divide data in two"
689         s = Sample(the.data)
690         s.xplain().show(); print("")
691
692 #-----
693 the=options(__doc__)
694 if __name__ == "__main__": demo(the.todo,Demos)
695
696 """
697 all config local to Sample
698 Example class
699 """

```