


```

138 #
139 #
140 #
141 #
142 #
143 #
144 #
145
146 def any(a:list) -> Any:
147     "Return a random item."
148     return a[anywhere(a)]
149
150 def anywhere(a:list) -> int:
151     "Return a random index of list 'a'."
152     return random.randint(0, len(a)-1)
153
154 big = sys.maxsize
155
156 def atom(x):
157     "Return a number or trimmed string."
158     x=x.strip()
159     if x=="True": return True
160     elif x=="False": return False
161     else:
162         try: return int(x)
163         except: return float(x)
164         except: return x.strip()
165
166 def demo(do,all):
167     "Maybe run a demo, if we want it, resetting random seed first."
168     todo = dir(all)
169     if do and do != "all":
170         todo = [x for x in dir(all) if x.startswith(do)]
171     for one in todo:
172         fun = all.__dict__[one, ""]
173         if type(fun)==type(demo):
174             random.seed(the.seed)
175             doc = re.sub(r'\n\s+', "\n", fun.__doc__ or "")
176             try:
177                 fun()
178                 print("PASS:", doc)
179             except Exception as e:
180                 all.fails += 0
181                 if the.cautious: traceback.print_exc(); exit(1)
182                 else: print("FAIL:", doc, e)
183             exit(all.fails)
184
185 def file(f):
186     "Iterator. Returns one row at a time, as cells."
187     with open(f) as fp:
188         for line in fp:
189             line = re.sub(r'(\n|\r|\t)|#.*', '', line)
190             if line:
191                 yield [cell.strip() for cell in line.split(",")]
192
193 def first(a:list) -> Any:
194     "Return first item."
195     return a[0]
196
197 def merge(b4:list) -> list:
198     "While we can find similar adjacent things, merge them."
199     j,n,now = -1, len(b4), []
200     while j < n-1:
201         j += 1
202         a = b4[j]
203         if j < n-2:
204             if merged := a.merge(b4[j+1]):
205                 a = merged
206                 j += 1 # we will continue, after missing one
207             now += [a]
208         # if 'now' is same size as 'b4', look for any other merges.
209     return b4 if len(now)==len(b4) else merge(now)
210
211 class o(object):
212     "Class that can pretty print its slots, with fast inits."
213     def __init__(i, **d): i.__dict__.update(**d)
214     def __repr__(i):
215         pre = i.__class__.__name__ if isinstance(i,o) else ""
216         return pre+'('+'(''.join([f":{k} {v}" for k, v in
217             sorted(i.__dict__.items()) if str(k)[0] != "_"]))+')?'
218
219 def options(doc:str) -> o:
220     """Convert 'doc' to options dictionary using command line args.
221     Args cause two 'shorthands': (1) boolean flags have no arguments (and mentioning
222     those on the command line means 'flip the default value'; (2) args need only
223     mention the first few of a key (e.g. -s is enough to select for -seed)."""
224     d={}
225     for line in doc.splitlines():
226         if line and line.startswith(" -"):
227             key, *_ = line.strip()[1:].split(" ") # get 1st,last word on each line
228             for j,flag in enumerate(sys.argv):
229                 if flag and flag[0]=="-" and key.startswith(flag[1:]):
230                     x="True" if x=="False" else "False" if x=="True" else sys.argv[j+1])
231                     d[key] = atom(x)
232             if d["help"]: exit(print(re.sub(r'\n#.*', "", doc, flags=re.S)))
233     return o(**d)
234
235 def r() -> float:
236     "Return random number 0..1"
237     return random.random()
238
239 def rn(x:float, n=3) -> float:
240     "Round a number to three decimals."
241     return round(x,n)
242
243 def rN(a:list, n=3) -> list:
244     "Round a list of numbers to three decimals."
245     return [rn(x,n=n) for x in a]
246
247 def second(a:list) -> Any:
248     "Return second item."
249     return a[1]
250

```

```

251 #
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 #
389 #
390 #
391 class Num(Col):
392     "Summarize numeric columns."
393     def __init__(i, size,**kw):
394         super().__init__(**kw)
395         i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, size, False
396
397     def add(i, x: float, inc=1):
398         "Reservoir sampler. If '_all' is full, sometimes replace an item at random."
399         if x != "":
400             i.n += inc
401             i.lo = min(x, i.lo)
402             i.hi = max(x, i.hi)
403             if len(i._all) < i.max : i.ok=False; i._all += [x]
404             elif r() < i.max/i.n: i.ok=False; i._all[anywhere(i._all)] = x
405         return x
406
407     def all(i):
408         "Return '_all'. sorted."
409         if not i.ok: i.ok=True; i._all.sort()
410         return i._all
411
412     def dist1(i, x, y):
413         if x=="?": y=i.norm(y); x=(1 if y<.5 else 0)
414         elif y=="?": x=i.norm(x); y=(1 if x<.5 else 0)
415         else : x, y = i.norm(x), i.norm(y)
416         return abs(x-y)
417
418     def div(i):
419         """Report the diversity of this distribution (using standard deviation).
420         &pm;2, 2.56, 3 &sigma; is 66,90,95%, of the mass. 2&sigma;. So one
421         standard deviation is (90-10)th divide by 2.4 times &sigma;."""
422         return (i.per(.9) - i.per(.1)) / 2.56
423
424     def merge(i, j):
425         "Return two 'Num's."
426         k = Num(i.max, at=i.at, txt=i.txt)
427         for x in i._all: k.add(x)
428         for x in j._all: k.add(x)
429         return k
430
431     def mid(i):
432         "Return central tendency of this distribution (using median)."
433         return i.per(.5)
434
435     def norm(i, x):
436         "Normalize 'x' to the range 0..1."
437         return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)
438
439     def per(i, p:float=.5) -> float:
440         "Return the p-th ranked item."
441         a = i.all(); return a[ int(p*len(a)) ]
442
443     def prep(i, x):
444         "Return 'x' as a float."
445         return x if x=="?" else float(x)
446
447     def spans(i, j, bins, out):
448         """Divide the whole space 'lo' to 'hi' into, say, 'xsmall'=16 bin.
449         then count the number of times we bin on other side.
450         Then merge similar adjacent bins."""
451         lo = min(i.lo, j.lo)
452         hi = max(i.hi, j.hi)
453         gap = (hi-lo) / bins
454         xys = [(x,"this",1) for x in i._all] + [
455             (x,"that",1) for x in j._all]
456         one = Span(i.lo, lo)
457         all = [one]
458         for x, y, n in sorted(xys, key=first):
459             if one.hi - one.lo > gap:
460                 one = Span(i, one.hi, x)
461                 all += [one]
462             one.add(x, y, n)
463         all = merge(all)
464         all[0].lo = -big
465         all[-1].hi = big
466         if len(all) > 1: out += all
467
468 #
469 #
470 #
471 #
472 #
473 #
474 class Example(o):
475     def __init__(i, cells) : i.cells=cells
476     def __getitem__(i, k) : return i.cells[k]
477
478     def dist(i, j, sample):
479         cols, p = sample.the.p
480         d = sum(col.dist(i[col.at], j[col.at])**p for col in cols)
481         return (d/len(cols)) ** (1/p)
482
483     def better(i, j, sample):
484         n = len(cols)
485         for col in cols:
486             a, b = col.norm( i[col.at] ), col.norm( j[col.at] )
487             s1 = math.e**(col.w*(a-b)/n)
488             s2 = math.e**(col.w*(b-a)/n)
489             return s1/n < s2/n
490
491 #
492 #
493 #
494 #
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         if not pre:
504             tmp = i.here.mid(i.here.y)
505             print(f"{pre[:40]}: {len(i.here.rows):5} : {tmp}")
506         if i.yes:
507             s=f"{pre}{i.span.show(True)}"
508             tmp = i.yes.here.mid(i.yes.here.y)
509             print(f"{s[:40]}: {len(i.yes.here.rows):5} : {tmp}")
510             i.yes.show(pre + "|. ")
511         if i.no:
512             s=f"{pre}{i.span.show(False)}"
513             tmp = i.no.here.mid(i.no.here.y)
514             print(f"{s[:40]}: {len(i.no.here.rows):5} : {tmp}")
515             i.no.show(pre + "|. ")
516
517 #
518 #
519 #
520 #
521 #
522 class Cluster(o):

```

```

523 "Tree with 'left','right' samples, broken at median between far points."
524 def __init__(i, here, x=None, y=None, c=None, mid=None):
525     i.here, i.x, i.y, i.c, i.mid, i.left, i.right = here, x, y, c, mid, None, None
526
527     def show(i, pre=""):
528         s=f"{pre[:40]}: {len(i.here.rows):5}"
529         print(f"{s}" if i.left else f"{s} : {i.here.mid(i.here.y)}")
530         for kid in [i.left, i.right]:
531             if kid: kid.show(pre + "|. ")
532
533 #
534 #
535 #
536 #
537 #
538 class Sample(o):
539     "Load, then manage, a set of examples."
540
541     def __init__(i, the, inits=[]):
542         i.the = the
543         i.rows, i.cols, i.x, i.y, i.klass = [], [], [], [], None
544         if atr==type(inits): [i.add(row, True) for row in file(inits)]
545         if list==type(inits): [i.add(row) for row in inits]
546
547     def add(i, a, raw=False):
548         def is_num(x) : return x[0].isupper()
549         def is_skip(x) : return x[-1]=="."
550         def is_klass(x): return "!" in x
551         def is_goal(x) : return "+" in x or "-" in x or is_klass(x)
552         col(at,txt):
553             now = Num(i.the.Max, at=at, txt=txt) if is_num(txt) else Sym(at=at, txt=txt)
554             where= i.y if is_goal(txt) else i.x
555             if not is_skip(txt):
556                 where += [now]
557             if is_klass(txt): i.klass = now
558             return now
559
560         pre = lambda a, c: c.prep(a[c.at]) if raw else a[c.at]
561         if i.cols: i.rows += [Example([col.add(pre(a, col)) for col in i.cols])]
562         else: i.cols = [col(at,txt) for at,txt in enumerate(a)]
563
564     def clone(i, inits=[]):
565         out = Sample(i.the)
566         out.add([col.txt for col in i.cols])
567         [out.add(x) for x in inits]
568         return out
569
570     def cluster(i, top=None):
571         """Split the data using random projections. Find the span that most
572         separates the data. Divide data on that span."""
573         here = Cluster(i)
574         top = top or i
575         if len(i.rows) >= 2*(len(top.rows)**i.the.enough):
576             left, right, x, y, c, mid = i.half(top)
577             if len(left.rows) < len(i.rows):
578                 here = Cluster(i, x, y, c, mid)
579             here.left = left.cluster(top)
580             here.right = right.cluster(top)
581             return here
582
583     def far(i, x, rows):
584         tmp= sorted([(x.dist(y,i),y) for y in (rows or i.rows)], key=first)
585         return tmp[ int(len(tmp)*i.the.far) ]
586
587     def half(i, top=None):
588         "Using two faraway points 'x,y' break data at median distance."
589         some= i.rows if len(i.rows)<i.the.Some else random.choices(i.rows, k=the.Som
590 e)
591         top= top or i
592         w = any(some)
593         _, x= top.far(w, some)
594         c, y= top.far(x, some)
595         tmp= [r for _, r in sorted([(top.proj(r, x, y, c), r)
596             for r in i.rows], key=first))]
597         mid= len(tmp)//2
598         return i.clone(tmp[:mid]), i.clone(tmp[mid:]), x, y, c, tmp[mid]
599
600     def mid(i, cols=None):
601         return [col.mid() for col in (cols or i.all)]
602
603     def proj(i, row, x, y, c):
604         "Find the distance of a 'row' on a line between 'x' and 'y'."
605         a = row.dist(x, i)
606         b = row.dist(y, i)
607         return (a**2 + c**2 - b**2) / (2*c)
608
609     def xplain(i, top=None):
610         """Split the data using random projections. Find the span that most
611         separates the data. Divide data on that span."""
612         here = Explain(i)
613         top = top or i
614         tiny = len(top.rows)**top.the.enough
615         if len(i.rows) >= 2*tiny:
616             left, right, *_ = i.half(top)
617             spans = []
618             [lcol.spans(rcol, 6/top.the.xsmall, spans) for lcol, rcol
619                 in zip(left.x, right.x)]
620             if len(spans) > 0:
621                 here.span = Span.sort(spans)[0]
622                 yes, no = i.clone(), i.clone()
623                 [yes if here.span.selects(row) else no].add(row) for row in i.rows]
624                 if tiny <= len(yes.rows) < len(i.rows): here.yes = yes.xplain(top=top)
625                 if tiny <= len(no.rows) < len(i.rows): here.no = no.xplain(top=top)
626             return here

```

```

626 #
627 #
628 #
629 #
630 #
631 #
632 #
633
634 class Demos:
635     "Possible start-up actions."
636     fails=0
637     def opt():
638         "show the config."
639         print(the)
640
641     def seed():
642         "seed"
643         assert .494 <= r() <= .495
644
645     def num():
646         "check 'Num'."
647         n = Num(512)
648         for _ in range(100): n.add(r())
649         assert .30 <= n.div() <= .31, "in range"
650
651     def sym():
652         "check 'Sym'."
653         s = Sym()
654         for x in "aaaabbc": s.add(x)
655         assert 1.37 <= s.div() <= 1.38, "entropy"
656         assert 'a' == s.mid(), "mode"
657
658     def rows():
659         "count rows in a file."
660         assert 399 == len([row for row in file(the.data)])
661
662     def sample():
663         "sampling."
664         s = Sample(the, the.data)
665         print(the.data, len(s.rows))
666         print(s.x[3], s.rows[-1])
667         assert 398 == len(s.rows), "length of rows"
668         assert 249 == s.x[-1].has['l'], "symbol counts"
669
670     def dist():
671         "distance between rows"
672         s = Sample(the, the.data)
673         assert .84 <= s.rows[1].dist(s.rows[-1],s) <= .842
674
675     def clone():
676         "cloning"
677         s = Sample(the, the.data)
678         s1 = s.clone(s.rows)
679         d1,d2 = s.x[0].__dict__, s1.x[0].__dict__
680         for k,v in d1.items():
681             print(d2[k],v)
682             assert d2[k] == v, "clone test"
683
684     def half():
685         "divide data in two"
686         s = Sample(the, the.data)
687         s1,s2,*_ = s.half()
688         print(s1.mid(s1.y))
689         print(s2.mid(s2.y))
690
691     def cluster():
692         "divide data in two"
693         s = Sample(the, the.data)
694         s.cluster().show(); print("")
695
696     def xplain():
697         "divide data in two"
698         s = Sample(the, the.data);
699         s.xplain().show(); print("")
700
701 #-----
702 the=options(__doc__)
703 if __name__ == "__main__": demo(the.todo,Demos)
704
705 """
706 Example class
707 """

```