

93 #
94 #
95 #
96 #
97 #
98 #
99 #

```

1  def any(a:list) -> Any:
2      "Return a random item."
3      return a[anywhere(a)]
4
5
6  def anywhere(a:list) -> int:
7      "Return a random index of list 'a'."
8      return random.randint(0, len(a)-1)
9
10
11 big = sys.maxsize
12
13
14 def atom(x):
15     "Return a number or trimmed string."
16     x=x.strip()
17     if x=="True": return True
18     elif x=="False": return False
19     else:
20         try: return int(x)
21         except:
22             try: return float(x)
23             except: return x.strip()
24
25
26 def demo(do,all):
27     "Maybe run a demo, if we want it, resetting random seed first."
28     todo = dir(all)
29     if do and do != "all":
30         todo = [x for x in dir(all) if x.startswith(do)]
31     for one in todo:
32         fun = all.__dict__.get(one,"")
33         if type(fun)==type(demo):
34             random.seed(the.seed)
35             doc = re.sub(r'\n\s+', "\n", fun.__doc__ or "")
36             try:
37                 fun()
38                 print("PASS:", doc)
39             except Exception as e:
40                 all.fails += 0
41                 if the.cautious: traceback.print_exc(); exit(1)
42                 else: print("FAIL:", doc, e)
43     exit(all.fails)
44
45
46 def file(f):
47     "Iterator. Returns one row at a time, as cells."
48     with open(f) as fp:
49         for line in fp:
50             line = re.sub(r'(\n|\\r|\\t|\\#|\\.|\\'|\\")', '', line)
51             if line:
52                 yield [atom(cell.strip()) for cell in line.split(",")]
53
54
55 def first(a:list) -> Any:
56     "Return first item."
57     return a[0]
58
59
60 def merge(b4:list) -> list:
61     "While we can find similar adjacent things, merge them."
62     j,n,now = -1,len(b4),[]
63     while j < n-1:
64         j += 1
65         a = b4[j]
66         if j < n-2:
67             if merged := a.merge(b4[j+1]):
68                 a = merged
69                 j += 1 # we will continue, after missing one
70         now += [a]
71     # if 'now' is same size as 'b4', look for any other merges.
72     return b4 if len(now)==len(b4) else merge(now)
73
74
75 class o(object):
76     "Class that can pretty print its slots, with fast inits."
77     def __init__(i,**d): i.__dict__.update(**d)
78     def __repr__(i):
79         pre = i.__class__.__name__ if isinstance(i,o) else ""
80         return pre+str(
81             {k: v for k, v in sorted(i.__dict__.items()) if str(k)[0] != "_"})
82
83
84 def options(doc:str) -> o:
85     """Convert 'doc' to options dictionary using command line args.
86     Args cause two 'shorthands': (1) boolean flags have no arguments (and mentioning
87     those on the command line means 'flip the default value'; (2) args need only
88     mention the first few of a key (e.g. -s is enough to select for -seed)."""
89     d={}
90     for line in doc.splitlines():
91         if line and line.startswith(" -"):
92             key, *, x = line.strip()[1:].split("#") # get 1st,last word on each line
93             for j,flag in enumerate(sys.argv):
94                 if flag and flag[0]=="-" and key.startswith(flag[1:]):
95                     x = "True" if x=="False" else ("False" if x=="True" else sys.argv[j+1])
96             d[key] = atom(x)
97     if d["help"]: exit(print(re.sub(r'\n#','',""),doc,flags=re.S))
98     return o(**d)
99
100
101 def r() -> float:
102     "Return random number 0..1"
103     return random.random()
104
105
106 def rn(x:float, n=3) -> float:
107     "Round a number to three decimals."
108     return round(x,n)
109
110
111 def rN(a:list, n=3) -> list:
112     "Round a list of numbers to three decimals."
113     return [rn(x,n=n) for x in a]
114
115
116 def second(a:list) -> Any:
117     "Return second item."
118     return a[1]
119
120
121 #the = options(__doc__)

```

```

208 #
209 #
210 #
211 #
212 #
213 #
214 #
215 #
216 #
217 #
218 #
219 #
220 #
221 class Span(o):
222     """Given two 'Sample's and some 'x' range 'lo,hi'.
223     a 'Span' holds often that range appears in each 'Sample'."""
224     def __init__(i,col, lo, hi, ys=None,):
225         i.col, i.lo, i.hi, i.ys = col, lo, hi,  ys or Sym()
226
227     def add(i, x:float, y:Any, inc=1) -> None:
228         "y" is a label identifying one 'Sample' or another."
229         i.lo = min(x, i.lo)
230         i.hi = max(x, i.hi)
231         i.ys.add(y,inc)
232
233     def merge(i, j): # -> Span|None
234         "If the merged span is simpler, return that merge."
235         a, b, c = i.ys, j.ys, i.ys.merge(j.ys)
236         if (i.ys.n==0 or j.ys.n==0 or
237             c.div()*0.99 <= (a.n*a.div() + b.n*b.div())/(a.n + b.n)):
238             return Span(i.col, min(i.lo,j.lo),max(i.hi,j.hi), ys=c)
239
240     def selects(i,row:list) -> bool:
241         "True if the range accepts the row."
242         x = row[i.col.at]; return x=="?" or i.lo<=x and x<i.hi
243
244     def show(i, positive=True) -> None:
245         "Show the range."
246         txt = i.col.txt
247         if positive:
248             if i.lo == i.hi: return f"[{txt}] == {i.lo}"
249             elif i.lo == -big: return f"[{txt}] < {i.hi}"
250             elif i.hi == big: return f"[{txt}] >= {i.lo}"
251             else : return f"[{i.lo}] <= {txt] < {i.hi}"
252         else:
253             if i.lo == i.hi: return f"[{txt}] != {i.lo}"
254             elif i.lo == -big: return f"[{txt}] >= {i.hi}"
255             elif i.hi == big: return f"[{txt}] < {i.lo}"
256             else : return f"[{txt}] < {i.lo} or {txt] >= {i.hi}"
257
258     def support(i) -> float:
259         "Returns 0..1."
260         return i.ys.n / i.col.n
261
262     @staticmethod
263     def sort(spans : list) -> list:
264         "Good spans have large support and low diversity."
265         divs, supports = Num(), Num()
266         sn = lambda s: supports.norm( s.support())
267         dn = lambda s: divs.norm( s.ys.div())
268         f = lambda s: ((1 - sn(s))*2 + dn(s)**2)**.5/2**.5
269         for s in spans:
270             divs.add( s.ys.div())
271             supports.add(s.support())
272         return sorted(spans, key=f)
273
274 #
275 #
276 #
277 #
278 #
279 class Col(o):
280     "Summarize columns."
281     def __init__(i,at=0,txt=""):
282         i.n,i.at,i.txt,i.w=0,at,txt,(-1 if "-" in txt else 1)
283
284     def dist(i,x:Any, y:Any) -> float:
285         return 1 if x=="?" and y=="?" else i.dist1(x,y)
286
287 #
288 #
289 #
290 #
291 #
292 class Sym(Col):
293     "Summarize symbolic columns."
294     def __init__(i,**kw):
295         super().__init__(**kw)
296         i.has, i.mode, i.most = {}, None, 0
297
298     def add(i, x:str, inc:int=1) -> str:
299         "Update symbol counts in 'has', updating 'mode' as we go."
300         if x != "?":
301             i.n += inc
302             tmp = i.has[x] = inc + i.has.get(x,0)
303             if tmp > i.most: i.most, i.mode = tmp, x
304             return x
305
306     def dist(i,x:str, y:str) ->float:
307         "Distance between two symbols."
308         return 0 if x==y else 1
309
310     def div(i):
311         "Return diversity of this distribution (using entropy)."
312         p = lambda x: x / (1E-31 + i.n)
313         return sum( -p(x)*math.log(p(x),2) for x in i.has.values() )
314
315     def merge(i, j):
316         "Merge two 'Sym's."
317         k = Sym(at=i.at, txt=i.txt)
318         for x,n in i.has.items(): k.add(x,n)
319         for x,n in j.has.items(): k.add(x,n)
320         return k
321
322     def mid(i):
323         "Return central tendency of this distribution (using mode)."
324         return i.mode
325
326     def spans(i, j, out):
327         """For each symbol in 'i' and 'j', count the
328         number of times we see it on either side."""
329         xys = [(x,"this",n) for x,n in i.has.items()] + [
330             (x,"that",n) for x,n in j.has.items()]
331         one, last = None,None
332         all = []
333         for x,y,n in sorted(xys, key=first):
334             if x != last:
335                 last = x
336                 one = Span(i, x,x)
337                 all += [one]
338                 one.add(x,y,n)
339             if len(all) > 1 : out += all

```

```

340 #
341 #
342 #
343 #
344 class Num(Col):
345     "Summarize numeric columns."
346     def __init__(i,**kw):
347         super().__init__(**kw)
348         i._all, i.lo, i.hi, i.max, i.ok = [], 1E32, -1E32, the.Max, False
349
350     def add(i,x: float ,inc=1):
351         "Reservoir sampler. If '_all' is full, sometimes replace an item at random."
352         if x != "":
353             i.n += inc
354             i.lo = min(x,i.lo)
355             i.hi = max(x,i.hi)
356             if len(i._all) < i.max : i.ok=False; i._all += [x]
357             elif r() < i.max/i.n: i.ok=False; i._all[anywhere(i._all)] = x
358             return x
359
360     def all(i):
361         "Return '_all', sorted."
362         if not i.ok: i.ok=True; i._all.sort()
363         return i._all
364
365     def dist1(i,x,y):
366         if x=="?": y=i.norm(y); x=(1 if y<.5 else 0)
367         elif y=="?": x=i.norm(x); y=(1 if x<.5 else 0)
368         else : x,y = i.norm(x), i.norm(y)
369         return abs(x-y)
370
371     def div(i):
372         """Report the diversity of this distribution (using standard deviation).
373         &pm;2, 2.56, 3 &sigma; is 66,90,95%, of the mass. 2&sigma;. So one
374         standard deviation is (90-10)th divide by 2.4 times &sigma;."""
375         return (i.per(.9) - i.per(.1)) / 2.56
376
377     def merge(i,j):
378         "Return two 'Num's."
379         k = Num(at=i.at, txt=i.txt)
380         for x in i._all: k.add(x)
381         for x in j._all: k.add(x)
382         return k
383
384     def mid(i):
385         "Return central tendency of this distribution (using median)."
386         return i.per(.5)
387
388     def norm(i,x):
389         "Normalize 'x' to the range 0..1."
390         return 0 if i.hi-i.lo < 1E-9 else (x-i.lo)/(i.hi-i.lo)
391
392     def per(i,p:float=.5) -> float:
393         "Return the p-th ranked item."
394         a = i.all(); return a[ int(p*len(a)) ]
395
396     def spans(i, j, out):
397         """Divide the whole space 'lo' to 'hi' into, say, 'xsmall'=16 bin,
398         then count the number of times we the bin on other side.
399         Then merge similar adjacent bins."""
400         lo = min(i.lo, j.lo)
401         hi = max(i.hi, j.hi)
402         gap = (hi-lo) / (6/the.xsmall)
403         xys = [(x,"this",1) for x in i._all] + [
404             (x,"that",1) for x in j._all]
405         one = Span(i,lo,lo)
406         all = [one]
407         for x,y,n in sorted(xys, key=first):
408             if one.hi - one.lo > gap:
409                 one = Span(i, one.hi,x)
410                 all += [one]
411                 one.add(x,y,n)
412             all = merge(all)
413             all[0].lo = -big
414             all[-1].hi = big
415             if len(all) > 1: out += all
416
417 #
418 #
419 #
420 #
421 #
422 class Explain(o):
423     "Tree with 'yes','no' branches for samples that do/do not match a 'span'."
424     def __init__(i,here):
425         i.here, i.span, i.yes, i.no = here, None, None, None
426
427     def show(i,pre=""):
428         if not pre:
429             tmp = i.here.mid(i.here.y)
430             print(f"[{pre}:40]: {len(i.here.rows):5} : {tmp}")
431         if i.yes:
432             s=f"[pre]{i.span.show(True)}"
433             tmp = i.yes.here.mid(i.yes.here.y)
434             print(f"[{s:40}: {len(i.yes.here.rows):5} : {tmp}")
435             i.yes.show(pre + "[. ")
436         if i.no:
437             s=f"[pre]{i.span.show(False)}"
438             tmp = i.no.here.mid(i.no.here.y)
439             print(f"[{s:40}: {len(i.no.here.rows):5} : {tmp}")
440             i.no.show(pre + "[. ")
441
442 #
443 #
444 #
445 #
446 #
447 class Cluster(o):
448     "Tree with 'left','right' samples, broken at median between far points."
449     def __init__(i,here,x=None,y=None,c=None,mid=None):
450         i.here,i.x,i.y,i.c,i.mid,i.left,i.right = here,x,y,c,mid,None,None
451
452     def show(i,pre=""):
453         s = f"[pre:40]: {len(i.here.rows):5}"
454         print(f"[{s}] " if i.left else f"[{s}] : {i.here.mid(i.here.y)}")
455         for kid in [i.left,i.right]:
456             if kid: kid.show(pre + "[. ")

```

```

457 #
458 #
459 #
460 #
461 #
462
463 class Sample(o):
464     "Load, then manage, a set of examples."
465     def __init__(i, inits=[]):
466         i.rows, i.cols, i.x, i.y, i.klass = [], [], [], [], None
467         if str == type(inits): [i.add(row) for row in file(inits)]
468         if list == type(inits): [i.add(row) for row in inits]
469
470     def add(i, a):
471         def col(at, txt):
472             what = Num if txt[0].isupper() else Sym
473             now = what(at=at, txt=txt)
474             where = i.y if "+" in txt or "-" in txt or "!" in txt else i.x
475             if txt[-1] != " ":
476                 where += [now]
477             if "!" in txt: i.klass = now
478             return now
479         #-----
480         if i.cols: i.rows += [[col.add(a[col.at]) for col in i.cols]]
481         else: i.cols = [col(at,txt) for at,txt in enumerate(a)]
482
483     def clone(i, inits=[]):
484         out = Sample()
485         out.add([col.txt for col in i.cols])
486         [out.add(x) for x in inits]
487         return out
488
489     def cluster(i, top=None):
490         """Split the data using random projections. Find the span that most
491         separates the data. Divide data on that span."""
492         here = Cluster(i)
493         top = top or i
494         if len(i.rows) >= 2*(len(top.rows)**the.enough):
495             left, right, x, y, c, mid = i.half(top)
496             if len(left.rows) < len(i.rows):
497                 here = Cluster(i, x, y, c, mid)
498             here.left = left.cluster(top)
499             here.right = right.cluster(top)
500             return here
501
502     def dist(i, x, y):
503         d = sum( col.dist(x[col.at], y[col.at])**the.p for col in i.x )
504         return (d/len(i.x)) ** (1/the.p)
505
506     def div(i, cols=None):
507         return [col.div() for col in (cols or i.all)]
508
509     def far(i, x, rows=None):
510         tmp = sorted([(i.dist(x, y), y) for y in (rows or i.rows)], key=first)
511         return tmp[ int(len(tmp)*the.far) ]
512
513     def half(i, top=None):
514         "Using two faraway points 'x,y' break data at median distance."
515         some = i.rows if len(i.rows) < the.Some else random.choices(i.rows, k=the.Some)
516         top = top or i
517         w = any(some)
518         _, x = top.far(w, some)
519         c, y = top.far(x, some)
520         tmp = [r for _, r in sorted([(top.proj(r, x, y, c), r)
521                                     for r in i.rows], key=first))]
522         mid = len(tmp)//2
523         return i.clone(tmp[:mid]), i.clone(tmp[mid:]), x, y, c, tmp[mid]
524
525     def mid(i, cols=None):
526         return [col.mid() for col in (cols or i.all)]
527
528     def proj(i, row, x, y, c):
529         "Find the distance of a 'row' on a line between 'x' and 'y'."
530         a = i.dist(row, x)
531         b = i.dist(row, y)
532         return (a**2 + c**2 - b**2) / (2*c)
533
534     def xplain(i, top=None):
535         """Split the data using random projections. Find the span that most
536         separates the data. Divide data on that span."""
537         here = Explain(i)
538         top = top or i
539         tiny = len(top.rows)**the.enough
540         if len(i.rows) >= 2*tiny:
541             left, right, *_ = i.half(top)
542             spans = []
543             [lcol.spans(rcol, spans) for lcol, rcol in zip(left.x, right.x)]
544             if len(spans) > 0:
545                 here.span = Span.sort(spans)[0]
546                 yes, no = i.clone(), i.clone()
547                 [yes if here.span.selects(row) else no].add(row) for row in i.rows]
548                 if tiny <= len(yes.rows) < len(i.rows): here.yes = yes.xplain(top=top)
549                 if tiny <= len(no.rows) < len(i.rows): here.no = no.xplain(top=top)
550             return here
551
552 #
553 #
554 #
555 #
556 #
557 #
558 #
559 #
560
561 class Demos:
562     "Possible start-up actions."
563     fails=0
564     def opt():
565         "show the config"
566         [print(f"{k}>10]={v}") for k,v in the.__dict__.items()]
567
568     def seed():
569         "seed"
570         assert .494 <= r() <= .495
571
572     def num():
573         "check 'Num'."
574         n = Num()
575         for _ in range(100): n.add(r())
576         assert .30 <= n.div() <= .31, "in range"
577
578     def sym():
579         "check 'Sym'."
580         s = Sym()
581         for x in "aaaabbc": s.add(x)
582         assert 1.37 <= s.div() <= 1.38, "entropy"
583         assert 'a' == s.mid(), "mode"
584
585     def rows():
586         "count rows in a file."
587         assert 399 == len([row for row in file(the.data)])
588
589     def sample():
590         "sampling"
591         s = Sample(the.data)
592         assert 398 == len(s.rows), "length of rows"
593         assert 249 == s.x[-1].has[1], "symbol counts"
594
595     def dist():
596         "distance between rows"
597         s = Sample(the.data)
598         assert .84 <= s.dist(s.rows[1], s.rows[-1]) <= .842
599
600     def far():
601         "distant items"
602         s = Sample(the.data)
603         for _ in range(32):
604             a, _ = s.far(any(s.rows))
605             assert a > .5, "large?"
606
607     def clone():
608         "cloning"
609         s = Sample(the.data)
610         s1 = s.clone(s.rows)
611         d1, d2 = s.x[0].__dict__, s1.x[0].__dict__
612         for k,v in d1.items():
613             assert d2[k] == v, "clone test"
614
615     def half():
616         "divide data in two"
617         s = Sample(the.data); s1,s2,*_ = s.half()
618         print(s1.mid(s1.y))
619         print(s2.mid(s2.y))
620
621     def cluster():
622         "divide data in two"
623         s = Sample(the.data)
624         s.cluster().show(); print("")
625
626     def xplain():
627         "divide data in two"
628         s = Sample(the.data)
629         s.xplain().show(); print("")
630
631 #-----
632 the=options(__doc__)
633 if __name__ == "__main__": demo(the.todo, Demos)
634
635 """
636 all config local to Sample
637 Example class
638 """

```