```python
  1:  #!/usr/bin/env python3.9
  2:  # vim: ts=2 sw=2 sts=2 et :
  3:  # autopep8 --exclude 'E20,E401,E226,E301,E302,E41'
  4:  """
  5:    ,-_|\    Contrast set learning
  6:   /   \   (c) Tim Menzies, 2021, unlicense.org
  7:   \_,-._*   Cluster, then report deltas between similiar
  8:      v     clusters (which will be few).
  9:  """
 10:  import re, sys, math, copy, argparse, random, itertools
 11:
 12:  def config(): return dict(
 13:     BINS=(float, .5, 'bins are of size n**BINS'),
 14:     COLS=(str, 'x', 'columns to use for inference'),
 15:     DATA=(str, '../data/auto2.csv', 'where to read data'),
 16:     FAR=(float, .9, 'where to look for far things'),
 17:     GOAL=(str, 'best', 'learning goals: best|rest|other'),
 18:     IOTA=(float, .3, 'small = sd**iota'),
 19:     K=(int, 2, 'bayes low class frequency hack'),
 20:     M=(int, 1, 'bayes low range frequency hack'),
 21:     P=(int, 2, 'distance calculation exponent'),
 22:     SAMPLES=(int, 20, '#samples to find far things?'),
 23:     SEED=(int, 10013, 'seed for random numbers'),
 24:     VERBOSE=(bool, False, 'set verbose'),
 25:     TOP=(int, 20, 'focus on this many'),
 26:     WILD=(int, False, 'run example, no protection'),
 27:     XAMPLE=(str, "", "egs: '-x ls' lists all, '-x all' runs all"))
 28:
 29:  class o(object):
 30:    def __init__(i, **k): i.__dict__.update(**k)
 31:    def __setitem__(i, k, v): i.__dict__[k] = v
 32:    def __getitem__(i, k): return i.__dict__[k]
 33:    def __repr__(i): return i.__class__.__name__ + str(
 34:      {k: v for k, v in i.__dict__.items() if k[0] != "_"})
 35:
 36:  # Columns _____
 37:  class Col(o):
 38:    "Store columns in 'Col', 'Skip', 'Sym', 'Num'."
 39:    def __init__(i, at=0, txt="", inits=[]):
 40:      i.n, i.at, i.txt = 0, at, txt
 41:      i.w = -1 if "-" in txt else 1
 42:      [i.add(x) for x in inits]
 43:
 44:    def add(i, x, n=1):
 45:      if x != "?":
 46:        i.n += n
 47:        x = i.add1(x, n)
 48:      return x
 49:  # _____
 50:  class Skip(Col):
 51:    def add1(i, x, n=1): return x
 52:  # _____
 53:  class Sym(Col):
 54:    def __init__(i, **kw):
 55:      i.has, i.mode, i.most = {}, None, 0
 56:      super().__init__(**kw)
 57:
 58:    def add1(i, x, n=1):
 59:      new = inc(i.has, x, n)
 60:      if new > i.most:
 61:        i.most, i.mode = new, x
 62:      return x
 63:

 64:    def bins(i, j, _):
 65:      for k in (i.has | j.has):
 66:        yield True, i.has.get(k,0), i.at, (k, k)
 67:        yield False, j.has.get(k,0), j.at, (k, k)
 68:
 69:    def dist(i, x, y): return 0 if x == y else 1
 70:
 71:    def ent(i):
 72:      return sum(-v/i.n * math.log(v/i.n) for v in i.has.values())
 73:
 74:    def merge(i, j):
 75:      k = Sym(at=i.at, txt=i.txt)
 76:      [k.add(x, n) for has in (i.has, j.has) for x, n in has.items()]
 77:      return k
 78:
 79:    def merged(i, j):
 80:      k = i.merge(j)
 81:      e1, n1, e2, n2, e, n = i.ent(), i.n, j.ent(), j.n, k.ent(), k.n
 82:      tmp = n1/n*e1 + n2/n*e2
 83:      if e1 + e2 < 0.01 or e * .95 < tmp:
 84:        return k
 85:
 86:    def mid(i): return i.mode
 87:  # _____
 88:  class Num(Col):
 89:    def __init__(i, **kw):
 90:      i._all, i.ok = [], False
 91:      super().__init__(**kw)
 92:
 93:    def add1(i, x, n):
 94:      x, i.ok = float(x), False
 95:      for _ in range(n):
 96:        i._all += [x]
 97:      return x
 98:
 99:    def all(i):
100:      if not i.ok:
101:        i.ok = True
102:        i._all = sorted(i._all)
103:      return i._all
104:
105:    def bins(i, j, the):
106:      xy = [(z, True) for z in i._all]+[(z, False) for z in j._all]
107:      iota = the.IOTA * (i.n*i.sd() + j.n*j.sd()) / (i.n + j.n)
108:      for ((lo, hi), sym) in bins(xy, iota=iota, size=len(xy)**the.BINS):
109:        yield True, sym.has.get(True,0), i.at, (lo, hi)
110:        yield False, sym.has.get(False,0), j.at, (lo, hi)
111:
112:    def dist(i, x, y):
113:      if x == "?":
114:        y = i.norm(y)
115:        x = 1 if y < 0.5 else 0
116:      elif y == "?":
117:        x = i.norm(x)
118:        y = 1 if x < 0.5 else 0
119:      else :
120:        x, y = i.norm(x), y.norm(y)
121:      return abs(x-y)
122:
123:    def mid(i): return per(i.all(), p=.5)
124:
125:    def norm(i, x):
126:      if x == "?":

127:      return x
128:      a = i.all()
129:      return max(0, min(1, (x-first(a))/(last(a)-first(a)+1E-32)))
130:
131:    def sd(i)    : return (per(i.all(), .9) - per(i.all(), .1))/2.56
132:    def span(i)  : return (first(i.all()), last(i.all()))
133:    def wide(i, n=0): return last(i.all()) - first(i.all()) >= n
134:
135:  # Row and Rows _____
136:  class Row(o):
137:    def __init__(i, lst, rows=None): i.rows, i.cells = rows, lst
138:
139:    def __lt__(i, j):
140:      goals = i.rows.cols.y
141:      s1, s2, n = 0, 0, len(goals)
142:      for col in goals:
143:        a = col.norm(i.cells[col.at])
144:        b = col.norm(j.cells[col.at])
145:        s1 -= math.e**(col.w * (a - b) / n)
146:        s2 -= math.e**(col.w * (b - a) / n)
147:      return s1 / n < s2 / n
148:
149:    def dist(i, j, the):
150:      d = n = 1E-32
151:      for col in i.rows.cols[the.COLS]:
152:        n += 1
153:        x, y = i.cells[at], j.cells[at]
154:        d += 1 if x == "?" and y == "?" else col.dist(x, y) ** the.P
155:      return (d/n) ** (1/the.P)
156:
157:    def far(i, rows, the):
158:      tmp = [(dist(i, j), j) for _ in range(the.SAMPLE)]
159:      return per(sorted(tmp, key=first), the.FAR)
160:
161:    def ys(i): return [i.cells[col.at] for col in i.rows.cols.y]
162:  # _____
163:  class Rows(o):
164:    def __init__(i, inits=[]):
165:      i.rows = []
166:      i.cols = o(all=[], names=[], x=[], y=[], klass=None)
167:      [i.add(x) for x in inits]
168:
169:    def add(i, a): i.data(a) if i.cols.names else i.header(a)
170:
171:    def best(i, the):
172:      i.rows.sort()
173:      ds = [the.IOTA*y.sd() for y in i.cols.y]
174:      best, rest = i.clone(), i.clone()
175:      for n, row in enumerate(i.rows):
176:        bestp = False
177:        for n1, n2, d in zip(i.rows[0].ys(), row.ys(), ds):
178:          bestp |= abs(n1-n2) <= d
179:        (best if bestp else rest).add(row)
180:      return best, rest
181:
182:    def clone(i, inits=[]): return Rows([i.cols.names] + inits)
183:
184:    def data(i, a):
185:      a = a.cells if type(a) == Row else a
186:      i.rows += [Row([col.add(a[col.at]) for col in i.cols.all],
187:           rows=i)]
188:
189:    def header(i, a):
```

```
190:      i.cols.names = a
191:      for at, x in enumerate(a):
192:        new = Skip if i.skipp(x) else (Num if i.nump(x) else Sym)
193:        new = new(at=at, txt=x)
194:        i.cols.all += [new]
195:        if not i.skipp(x):
196:          i.cols["y" if i.yp(x) else "x"] += [new]
197:          if i.klassp(x):
198:            i.cols.klass = new
199:
200:    def klassp(i, x): return "!" in x
201:    def nump(i, x): return x[0].isupper()
202:    def skipp(i, x): return "?" in x
203:    def yp(i, x): return "-" in x or "+" in x or i.klassp(x)
204:    def ys(i): return [col.mid() for col in i.cols.y]
205:    def ysd(i): return [col.sd() for col in i.cols.y]
206:  # _____
207:  def stratify(src):
208:    all, klass = None, {}
209:    for n, row in enumerate(src):
210:      if all:
211:        kl = row[all.cols.klass.at]
212:        here = klass[kl] = klass.get(kl, None) or all.clone()
213:        here.add(row)
214:        all.add(row)
215:      else:
216:        all = Rows([row])
217:    return o(all=all, klass=klass)
218:
219:  # Discretizations _____
220:  # Use 'bins' to divide numeric data into ranges.
221:  def bins(xy, iota=0, size=30):
222:    def merge(b4):
223:      j, tmp, n = 0, [], len(b4)
224:      while j < n:
225:        ((lo, _), ay) = a = b4[j]
226:        if j < n - 1:
227:          ((_, hi), by) = b4[j + 1]
228:          if cy := ay.merged(by):
229:            a = ((lo, hi), cy)
230:            j += 1
231:        tmp += [a]
232:        j += 1
233:      return merge(tmp) if len(tmp) < len(b4) else b4
234:
235:    def divide(xy):
236:      bin = o(x=Num(), y=Sym())
237:      bins = [bin]
238:      for i, (x, y) in enumerate(xy):
239:        if bin.x.n >= size and x != b4:
240:          if i < len(xy)-size and bin.x.wide(iota):
241:            bin = o(x=Num(), y=Sym())
242:            bins += [bin]
243:        bin.x.add(x)
244:        bin.y.add(y)
245:        b4 = x
246:      return bins
247:
248:    xy = sorted(xy, key=first)
249:    return merge([(bin.x.span(), bin.y) for bin in divide(xy)])
250:
251:  # Learn class deltas _____
252:  def contrasts(here, there, the):
```

```
253:    goal = {'best' : lambda b, r: b**2/b+r,
254:            'rest' : lambda b, r: r**2/(b+r),
255:            'other': lambda b, r: 1/(b+r)    }[the.GOAL]
256:
257:    def like(d, kl):
258:      out = prior = (hs[kl] + the.K) / (n + the.K*2)
259:      for at, span in d.items():
260:        f = has.get((kl, (at, span)), 0)
261:        out *= (f + the.M*prior) / (hs[kl] + M)
262:      return out
263:
264:    def val(d): return goal(like(d, True), like(d, False)), d
265:    def top(a): return sorted(a, reversed=True, key=first)[:the.TOP]
266:
267:    has = {(kl, (at, (lo, hi))): f
268:           for col1, col2 in zip(here.cols.x, there.cols.x)
269:           for at, (lo, hi) in col1.bins(col2, the)}
270:    n = len(here.rows, there.rows)
271:    hs = {True: len(here.rows), False: len(there.rows)}
272:    solos = [val(dict(at=x)) for at, x in set([z for _, z in has])]
273:    ranges = {}
274:    for _, d in top(solos):
275:      for k in d:
276:        ranges[k] = ranges.get(k, set()).add(d[k])
277:    for rule in top([val(d) for d in dict_product(ranges)]):
278:      print(rule)
279:
280:  # Misc utils _____
281:  # string stuff
282:  def color(end="\n", **kw):
283:    s, a, z = "", "\u001b[", ";1m"
284:    c = dict(black=30, red=31, green=32, yellow=33,
285:             purple=34, pink=35, blue=36, white=37)
286:    for col, txt in kw.items():
287:      s = s+a + str(c[col]) + z+txt+"\033[0m"
288:    print(s, end=end)
289:
290:  def mline(m): m += [["-"*len(str(x)) for x in m[-1]]]
291:
292:  def printm(matrix):
293:    s = [[str(e) for e in row] for row in matrix]
294:    lens = [max(map(len, col)) for col in zip(*s)]
295:    fmt = ' | '.join('{{:>{}}}'.format(x) for x in lens)
296:    for row in [fmt.format(*row) for row in s]:
297:      print(row)
298:
299:  # maths stuff
300:  def r3(a): return [round(x, 3) for x in a]
301:
302:  # dictionary stuff
303:  def has(d, k): return d.get(k, 0)
304:  def inc(d, k, n=1): tmp = d[k] = n + d.get(k, 0); return tmp
305:
306:  def dict_product(d):
307:    keys = d.keys()
308:    for p in itertools.product(*d.values()):
309:      yield dict(zip(keys, p))
310:
311:  # list stuff
312:  def first(a): return a[0]
313:  def last(a): return a[-1]     # $\label{comment}$
314:  def per(a, p=.5): return a[int(p*len(a))]
315:
```

```
316:  # file stuff
317:  def csv(f=None, sep=","):
318:    def prep(s): return re.sub(r'([\n\t\r ]|#.*)', '', s)
319:    if f:
320:      with open(f) as fp:
321:        for s in fp:
322:          if s := prep(s):
323:            yield s.split(sep)
324:    else:
325:      for s in sys.stdin:
326:        if s := prep(s):
327:          yield s.split(sep)
328:
329:  # command-line stuff
330:  def cli(use, txt, config):
331:    fmt = argparse.RawTextHelpFormatter
332:    used, p = {}, argparse.ArgumentParser(prog=use, description=txt,
333:                       formatter_class=fmt)
334:    for k, (_, b4, h) in config.items():
335:      k0 = k[0]
336:      used[k0] = c = k0 if k0 in used else k0.lower()
337:      if b4 == False:
338:        p.add_argument("-"+c, dest=k, default=False,
339:                help=h,
340:                action="store_true")
341:      else:
342:        p.add_argument("-"+c, dest=k, default=b4,
343:                help=h + " [" + str(b4) + "]",
344:                type=type(b4), metavar=k)
345:    return o( **p.parse_args().__dict__ )
346:
347:  # Unit tests _____
348:  class Eg:
349:    def ls(the):
350:      "list  all examples."
351:      print("\nexamples:")
352:      for k, f in vars(Eg).items():
353:        if k[0] != "_":
354:          print(f"  {k:<13} {f.___doc___}")
355:
356:    def _fail(the):
357:      "testing failure"
358:      assert False, "failing"
359:
360:    def data(the, file="../data/vote.csv"):
361:      "simple load of data into  a table"
362:      r = Rows(csv(file))
363:      assert 435 == len(r.rows)
364:      assert 195 == r.cols.all[1].has['y']
365:
366:    def nclasses(the, file="../data/diabetes.csv", kl="positive"):
367:      "read data with nclasses"
368:      rs = stratify(csv(file))
369:      assert 2 == len(rs.klass)
370:      assert 268 == len(rs.klass[kl].rows)
371:      assert 768 == len(rs.all.rows)
372:      assert 3.90625 == rs.klass[kl].cols.all[0].sd()
373:
374:    def bins(the, file="../data/diabetes.csv",
375:         k1="positive", k2="negative"):
376:      "discretize some data"
377:      rs = stratify(csv(file))
378:      bins1(rs.klass[k1], rs.klass[k2], the)
```

```
379:
380:    def bestrest(the, file="../data/auto93.csv"):
381:      "discretize some multi-goal data"
382:      r = Rows(csv(file))
383:      goods, bads = r.best(the)
384:      bins1(goods, bads, the)
385:
386:  def bins1(goods, bads, the):
387:    for good, bad in zip(goods.cols.x, bads.cols.x):
388:      bins = sorted(good.bins(bad, the))
389:      if len(bins) > 1:
390:        print(f"\n{good.txt}")
391:        for bin in bins:
392:          print("\t", bin)
393:
394:  # Main program _____
395:  def main(the):
396:    def run(fun, fails, the):
397:      s = f" {fun.__name__:<12}"
398:      if the.WILD:
399:        print("wild")
400:        fun(copy.deepcopy(the))
401:        sys.exit()
402:      try:
403:        fun(copy.deepcopy(the))
404:        random.seed(the.SEED)
405:        color(green=(chr(10003) + s), white=fun.__doc__)
406:      except Exception as err:
407:        fails = fails + 1
408:        color(red=(chr(10007) + s), white=str(err))
409:      return fails
410:    # _____
411:    fails = 0
412:    if the.XAMPLE == "all":
413:      for k, f in vars(Eg).items():
414:        if k[0] != "_" and k != "ls":
415:          fails = run(f, fails, the)
416:    else:
417:      if the.XAMPLE and the.XAMPLE in vars(Eg):
418:        f = vars(Eg)[the.XAMPLE]
419:        if the.XAMPLE == "ls":
420:          f(the)
421:        else :
422:          fails = run(f, fails, the)
423:    sys.exit(fails)
424:
425:
426:  if __name__ == "__main__":
427:    main( cli("./keys", __doc__, config()) )
```