

```

1: #!/usr/bin/env python3.9
2: # vim: ts=2 sw=2 sts=2 et :
3: # autotepe8 --exclude 'E20,E401,E226,E301,E302,E41'
4: """
5:     \ Contrast set learning
6:     \ (c) Tim Menzies, 2021, unlicense.org
7:     \ Cluster, then reports, just the
8:     \ v deltas between similar clusters.
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:     TOP=( int, 10, 'focus on this many'),
25:     UNSAFE=( bool, False, 'run example, no protection'),
26:     VERBOSE=(bool, False, 'set verbose'),
27:     XAMPLE=( str, "", "'-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:     def merge(i, j):
51:         k = i.__class__(at=i.at, txt=i.txt)
52:         [k.add(x, n) for has in (i.has, j.has) for x, n in has.items()]
53:         return k
54:
55:     def merged(i, j):
56:         k = i.merge(j)
57:         vi, vj, vk = i.var(), j.var(), k.var()
58:         if vk < (i.n*vi + j.n*vj) / (i.n + j.n):
59:             #if vi + vj < 0.01 or vk *.95 < (i.n*vi + j.n*vj) / (i.n + j.n):
60:             return k
61:
62: #
63: class Skip(Col):
64:     def add1(i, x, n=1): return x
65: #
66: class Sym(Col):
67:     def __init__(i, **kw):
68:         i.has, i.mode, i.most = {}, None, 0
69:         super().__init__(**kw)
70:
71:     def add1(i, x, n=1):
72:         new = inc(i.has, x, n)
73:         if new > i.most:
74:             i.most, i.mode = new, x
75:         return x
76:
77:     def bins(i, j, _):
78:         for k in (i.has | j.has):
79:             yield True, i.has.get(k, 0), i.at, (k, k)
80:             yield False, j.has.get(k, 0), j.at, (k, k)
81:
82:     def dist(i, x, y): return 0 if x == y else 1
83:     def mid(i): return i.mode
84:
85:     def var(i):
86:         return sum(-v/i.n * math.log(v/i.n) for v in i.has.values())
87:
88: #
89: class Num(Col):
90:     def __init__(i, **kw):
91:         i._all, i.ok = [], False
92:         super().__init__(**kw)
93:
94:     def add1(i, x, n):

```

```

95:         x, i.ok = float(x), False
96:     for _ in range(n):
97:         i._all += [x]
98:     return x
99:
100: def all(i):
101:     if not i.ok:
102:         i.ok = True
103:         i._all = sorted(i._all)
104:     return i._all
105:
106: def bins(i, j, the):
107:     xy = [(z, True) for z in i._all] + [(z, False) for z in j._all]
108:     iota = the.IOTA * (i.n*i.sd() + j.n*j.sd()) / (i.n + j.n)
109:     for ((lo, hi), sym) in bins(xy, iota=iota, size=len(xy)**the.BINS):
110:         yield True, sym.has.get(True, 0), i.at, (lo, hi)
111:     yield False, sym.has.get(False, 0), j.at, (lo, hi)
112:
113: def dist(i, x, y):
114:     if x == "?":
115:         y = i.norm(y)
116:         x = 1 if y < 0.5 else 0
117:     elif y == "?":
118:         x = i.norm(x)
119:         y = 1 if x < 0.5 else 0
120:     else:
121:         x, y = i.norm(x), i.norm(y)
122:     return abs(x-y)
123:
124: def mid(i): return per(i.all(), p=.5)
125:
126: def norm(i, x):
127:     if x == "?":
128:         return x
129:     a = i.all()
130:     return max(0, min(1, (x-first(a))/(last(a)-first(a)+1E-32)))
131:
132: def sd(i) : return (per(i.all(), .9) - per(i.all(), .1))/2.56
133: def span(i) : return (first(i.all()), last(i.all()))
134: def var(i) : return i.sd()
135: def wide(i, n=0): return last(i.all()) - first(i.all()) >= n
136:
137: # Row and Rows
138: class Row(o):
139:     def __init__(i, lst, rows=None): i.rows, i.cells = rows, lst
140:
141:     def __lt__(i, j):
142:         goals = i.rows.cols.y
143:         s1, s2, n = 0, 0, len(goals)
144:         for col in goals:
145:             a = col.norm(i.cells[col.at])
146:             b = col.norm(j.cells[col.at])
147:             s1 -= math.e**(col.w * (a - b) / n)
148:             s2 -= math.e**(col.w * (b - a) / n)
149:         return s1 / n < s2 / n
150:
151:     def dist(i, j, the):
152:         d = n = 1E-32
153:         for col in i.rows.cols[the.COLS]:
154:             n += 1
155:             x, y = i.cells[at], j.cells[at]
156:             d += 1 if x == "?" and y == "?" else col.dist(x, y) ** the.P
157:         return (d/n) ** (1/the.P)
158:
159:     def far(i, rows, the):
160:         tmp = [(dist(i, j), j) for _ in range(the.SAMPLE)]
161:         return per(sorted(tmp, key=first), the.FAR)
162:
163:     def ys(i): return [i.cells[col.at] for col in i.rows.cols.y]
164: #
165: class Rows(o):
166:     def __init__(i, inits=[]):
167:         i.rows = []
168:         i.cols = o(all=[], names=[], x=[], y=[], klass=None)
169:         [i.add(x) for x in inits]
170:
171:     def add(i, a): i.data(a) if i.cols.names else i.header(a)
172:
173:     def best(i, the):
174:         i.rows.sort()
175:         ds = [the.IOTA*y.sd() for y in i.cols.y]
176:         best, rest = i.clone(), i.clone()
177:         for n, row in enumerate(i.rows):
178:             bestp = False
179:             for n1, n2, d in zip(i.rows[0].ys(), row.ys(), ds):
180:                 bestp |= abs(n1-n2) <= d
181:             (best if bestp else rest).add(row)
182:         return best, rest
183:
184:     def clone(i, inits=[]): return Rows([i.cols.names] + inits)
185:
186:     def data(i, a):
187:         a = a.cells if type(a) == Row else a
188:         i.rows += [Row([col.add(a[col.at]) for col in i.cols.all],

```

```

189:         rows=i)]
190:
191: def display(i, s, rule):
192: def n(x): return int(x) if x == int(x) else x
193: return f"{s:>.5f}: " + (' and ' .join(
194:     (i.cols.names[k])+" " = "+(
195:     str(n(lo)) if lo == hi else f"{n(lo)}..{n(hi)}")
196:     for k, (lo, hi) in sorted(rule.items()))))
197:
198: def header(i, a):
199: i.cols.names = a
200: for at, x in enumerate(a):
201:     new = Skip if i.skipp(x) else (Num if i.nump(x) else Sym)
202:     new = new(at=at, txt=x)
203:     i.cols.all += [new]
204:     if not i.skipp(x):
205:         i.cols["y" if i.y(x) else "x"] += [new]
206:     if i.klassp(x):
207:         i.cols.klass = new
208:
209: def klassp(i, x): return "!" in x
210: def nump(i, x): return x[0].isupper()
211: def skipp(i, x): return "?" in x
212: def yp(i, x): return "-" in x or "+" in x or i.klassp(x)
213: def ys(i): return [col.mid() for col in i.cols.y]
214: def ysd(i): return [col.sd() for col in i.cols.y]
215: #
216: def stratify(src):
217:     all, klass = None, {}
218:     for n, row in enumerate(src):
219:         if all:
220:             kl = row[all.cols.klass.at]
221:             here = klass[kl] = klass.get(kl, None) or all.clone()
222:             here.add(row)
223:             all.add(row)
224:         else:
225:             all = Rows([row])
226:     return o(all=all, klass=klass)
227:
228: # Discretizations
229: # Use 'bins' to divide numeric data into ranges.
230: def bins(xy, iota=0, size=30):
231: def merge(b4):
232:     j, tmp, n = 0, [], len(b4)
233:     while j < n:
234:         ((lo, _), ay) = a = b4[j]
235:         if j < n - 1:
236:             ((_, hi), by) = b4[j + 1]
237:             if cy := ay.merged(by):
238:                 a = ((lo, hi), cy)
239:                 j += 1
240:             tmp += [a]
241:             j += 1
242:     return merge(tmp) if len(tmp) < len(b4) else b4
243:
244: def divide(xy):
245:     bin = o(x=Num(), y=Sym())
246:     bins = [bin]
247:     for i, (x, y) in enumerate(xy):
248:         if bin.x.n >= size and x != b4:
249:             if i < len(xy)-size and bin.x.wide(iota):
250:                 bin = o(x=Num(), y=Sym())
251:                 bins += [bin]
252:             bin.x.add(x)
253:             bin.y.add(y)
254:             b4 = x
255:     return bins
256:
257: xy = sorted(xy, key=first)
258: return merge([(bin.x.span(), bin.y) for bin in divide(xy)])
259:
260: # Learn class deltas
261: def contrasts(here, there, the):
262:     goal = ('best' : lambda b, r: b**2/(b+r),
263:            'rest' : lambda b, r: r**2/(b+r),
264:            'other': lambda b, r: 1/(b+r) ) [the.GOAL]
265:
266: def like(d, kl):
267:     out = prior = (hs[kl] + the.K) / (n + the.K*2)
268:     for at, span in d.items():
269:         f = has.get((kl, at, span), 0)
270:         out *= (f + the.M*prior) / (hs[kl] + the.M)
271:     return out
272:
273: def top(a): return sorted(a, reverse=True, key=first)[:the.TOP]
274:
275: def val(d):
276:     for k, v in list(d.items()):
277:         if v == None:
278:             del d[k]
279:     s = goal(like(d, True), like(d, False)) if d else 0
280:     return s, d
281:
282: has = {(klass, at, (lo, hi)): f

```

```

283:         for col1, col2 in zip(here.cols.x, there.cols.x)
284:         for klass, f, at, (lo, hi) in col1.bins(col2, the))
285:     n = len(here.rows) + len(there.rows)
286:     hs = {True: len(here.rows), False: len(there.rows)}
287:     all= set([(at, span) for (_, at, span) in has])
288:     for x,ys in sorted(all):
289:         print(f"\n{here.cols.names[x]}")
290:         [print(f"\t{y}") for y in ys if y]
291:     solos = [val({at: span}) for (_, at, span) in all]
292:     ranges = {}
293:     for _, d in top(solos):
294:         for k in d:
295:             if k not in ranges:
296:                 ranges[k] = set([None])
297:                 ranges[k].add(d[k])
298:             print("\nrules:")
299:             for s, rule in top([val(d) for d in dict_product(ranges)]):
300:                 if s and rule:
301:                     print("\t"+here.display(s, rule))
302:
303: # Misc utils
304: # string stuff
305: def color(end="\n", **kw):
306:     s, a, z = "", "\u001b[", ";1m"
307:     c = dict(black=30, red=31, green=32, yellow=33,
308:             purple=34, pink=35, blue=36, white=37)
309:     for col, txt in kw.items():
310:         s = s+a + str(c[col]) + z+txt+"033[0m"
311:     print(s, end=end)
312:
313: def mline(m): m += [["-"*len(str(x)) for x in m[-1]]]
314:
315: def mlines(matrix):
316:     s = [[str(e) for e in row] for row in matrix]
317:     lens = [max(map(len, col)) for col in zip(*s)]
318:     fmt = ' | '.join('{{: >{} }}'.format(x) for x in lens)
319:     for row in [fmt.format(*row) for row in s]:
320:         print(row)
321:
322: # maths stuff
323: def r3(a): return [round(x, 3) for x in a]
324:
325: # dictionary stuff
326: def has(d, k): return d.get(k, 0)
327: def inc(d, k, n=1): tmp = d[k] = n + d.get(k, 0); return tmp
328:
329: def dict_product(d):
330:     keys = d.keys()
331:     for p in itertools.product(*d.values()):
332:         yield dict(zip(keys, p))
333:
334: # list stuff
335: def first(a): return a[0]
336: def last(a): return a[-1] # $Vlabel(comment)$
337: def per(a, p=.5): return a[int(p*len(a))]
338:
339: # file stuff
340: def csv(src=None, sep=",", zap=r'([\n\t\r ]|#.*)':
341:     if lines := (not src and sys.stdin or
342:                 type(src) == list and src or
343:                 "\n" in src and src.splitlines() or
344:                 None):
345:         for line in lines:
346:             if line := re.sub(zap, "", line):
347:                 yield line.split(sep)
348:         else:
349:             with open(src) as lines:
350:                 for line in lines:
351:                     if line := re.sub(zap, "", line):
352:                         yield line.split(sep)
353:
354: # command-line stuff
355: def cli(use, txt, config):
356:     fmt = argparse.RawTextHelpFormatter
357:     used, p = {}, argparse.ArgumentParser(prog=use, description=txt,
358:                                         formatter_class=fmt)
359:     for k, (_, b4, h) in sorted(config.items()):
360:         k0 = k[0]
361:         used[k0] = c = k0 if k0 in used else k0.lower()
362:         if b4 == False:
363:             p.add_argument("-"+c, dest=k, default=False,
364:                           help=h,
365:                           action="store_true")
366:         else:
367:             p.add_argument("-"+c, dest=k, default=b4,
368:                           help=h + " [" + str(b4) + "]",
369:                           type=type(b4), metavar=k)
370:     return o(**p.parse_args().__dict__)
371:
372: # Unit tests
373: class Eg:
374:     def ls(the):
375:         "list all examples."
376:         print("\nexamples:")

```

```

377: for k, f in vars(Eg).items():
378:     if k[0] != " ":
379:         print(f" {k:<13} {f.__doc__}")
380:
381: def _fail(the):
382:     "testing failure"
383:     assert False, "failing"
384:
385: def csv(the, file="./data/vote.csv"):
386:     "simple load of data into a table"
387:     for r in csv(file): print(r)
388:
389: def data(the, file="./data/vote.csv"):
390:     "simple load of data into a table"
391:     r = Rows(csv(file))
392:     assert 435 == len(r.rows)
393:     assert 195 == r.cols.all[1].has['y']
394:
395: def nclasses(the, file="./data/diabetes.csv", kl="positive"):
396:     "ead data with nclasses"
397:     rs = stratify(csv(file))
398:     assert 2 == len(rs.klass)
399:     assert 268 == len(rs.klass[kl].rows)
400:     assert 768 == len(rs.all.rows)
401:     assert 3.90625 == rs.klass[kl].cols.all[0].sd()
402:
403: def bins(the, file="./data/diabetes.csv",
404:          k1="positive", k2="negative"):
405:     "discretize some data"
406:     rs = stratify(csv(file))
407:     bins1(rs.klass[k1], rs.klass[k2], the)
408:
409: def bestrest(the, file="./data/auto93.csv"):
410:     "discretize some multi-goal data"
411:     r = Rows(csv(file))
412:     goods, bads = r.best(the)
413:     bins1(goods, bads, the)
414:
415: def con1(the, file="./data/auto93.csv"):
416:     "discretize some multi-goal data"
417:     r = Rows(csv(file))
418:     goods, bads = r.best(the)
419:     contrasts(goods, bads, the)
420:
421: def con2(the):
422:     "discretize china"
423:     Eg.con1(the, "./data/china.csv")
424:
425: def con3(the, file="./data/diabetes.csv",
426:          k1="positive", k2="negative"):
427:     "contrast sets from diabetes"
428:     rs = stratify(csv(file))
429:     contrasts(rs.klass[k1], rs.klass[k2], the)
430:
431: def bins1(goods, bads, the):
432:     for good, bad in zip(goods.cols.x, bads.cols.x):
433:         bins = sorted(good.bins(bad, the))
434:         if len(bins) > 1:
435:             print(f"\n{good.txt}")
436:             for bin in bins:
437:                 print("\t", bin)
438:
439: # Main program
440: def main(the):
441:     def run(fun, fails, the):
442:         s = f" {fun.__name__:<12}"
443:         if the.UNSAFE:
444:             print("unsafe mode:")
445:             fun(copy.deepcopy(the))
446:             sys.exit()
447:         try:
448:             fun(copy.deepcopy(the))
449:             random.seed(the.SEED)
450:             color(green=(chr(10003) + s), white=fun.__doc__)
451:         except Exception as err:
452:             fails = fails + 1
453:             color(red=(chr(10007) + s), white=str(err))
454:         return fails
455:     #
456:     fails = 0
457:     if the.XAMPLE == "all":
458:         for k, f in vars(Eg).items():
459:             if k[0] != " " and k != "ls":
460:                 fails = run(f, fails, the)
461:     else:
462:         if the.XAMPLE and the.XAMPLE in vars(Eg):
463:             f = vars(Eg)[the.XAMPLE]
464:             if the.XAMPLE == "ls":
465:                 f(the)
466:             else:
467:                 fails = run(f, fails, the)
468:     sys.exit(fails)
469:
470:

```

```

471: if __name__ == "__main__":
472:     main(cli("./keys", __doc__, config()))

```