

## keys

```

1: #!/usr/bin/env python3.9
2: # vim: ts=2 sw=2 sts=2 et :
3: # autopen8 --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: #
51: class Skip(Col):
52:     def add1(i, x, n=1): return x
53:
54: #
55: class Sym(Col):
56:     def __init__(i, **kw):
57:         i.has, i.mode, i.most = {}, None, 0
58:         super().__init__(**kw)
59:
60:     def add1(i, x, n=1):
61:         new = inc(i.has, x, n)
62:         if new > i.most:
63:             i.most, i.mode = new, x
64:         return x
65:
66:     def bins(i, j, _):
67:         for k in (i.has | j.has):
68:             yield True, i.has.get(k, 0), i.at, (k, k)
69:             yield False, j.has.get(k, 0), j.at, (k, k)
70:
71:     def dist(i, x, y): return 0 if x == y else 1
72:
73:     def ent(i):
74:         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: #
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), y.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:
134:     def span(i): return (first(i.all()), last(i.all()))
135:
136:     def wide(i, n=0): return last(i.all()) - first(i.all()) >= n
137:
138: # Row and Rows
139: class Row(o):
140:     def __init__(i, lst, rows=None): i.rows, i.cells = rows, lst
141:
142:     def __lt__(i, j):
143:         goals = i.rows.cols.y
144:         s1, s2, n = 0, 0, len(goals)
145:         for col in goals:
146:             a = col.norm(i.cells[col.at])
147:             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: #
164: class Rows(o):
165:     def __init__(i, inits=[]):
166:         i.rows = []
167:         i.cols = o(all=[], names=[], x=[], y=[], klass=None)
168:         [i.add(x) for x in inits]
169:
170:     def add(i, a): i.data(a) if i.cols.names else i.header(a)
171:
172:     def best(i, the):
173:         i.rows.sort()
174:         ds = [the.IOTA*y.sd() for y in i.cols.y]
175:         best, rest = i.clone(), i.clone()
176:         for n, row in enumerate(i.rows):
177:             bestp = False
178:             for n1, n2, d in zip(i.rows[0].ys(), row.ys(), ds):
179:                 bestp = abs(n1-n2) <= d
180:                 (best if bestp else rest).add(row)
181:             return best, rest
182:
183:     def clone(i, inits=[]): return Rows([i.cols.names] + inits)
184:
185:     def data(i, a):
186:         a = a.cells if type(a) == Row else a
187:         i.rows += [Row([col.add(a[col.at]) for col in i.cols.all],
188:             rows=i)]
189:
190:     def display(i, s, rule):
191:         def n(x): return int(x) if x==int(x) else x
192:         return f'{s:>.5f}: ' + (' and '.join(
193:             (i.cols.names[k])+":"+str(n(i)) if lo==hi else f'{n(lo)}...{n(hi)}'
194:             for k, (lo, hi) in sorted(rule.items())))
195:
196:     def header(i, a):
197:         i.cols.names = a
198:         for at, x in enumerate(a):
199:             new = Skip if i.skipp(x) else (Num if i.nump(x) else Sym)
200:             new = new(at=at, txt=x)
201:             i.cols.all += [new]
202:             if not i.skipp(x):
203:                 i.cols["y"] if i.y(x) else "x" += [new]
204:             if i.klassp(x):
205:                 i.cols.klass = new
206:
207:     def klassp(i, x): return "!" in x
208:
209:     def nump(i, x): return x[0].isupper()
210:
211:     def skipp(i, x): return "?" in x
212:
213:     def yp(i, x): return "-" in x or "+" in x or i.klassp(x)
214:
215:     def ys(i): return [col.mid() for col in i.cols.y]
216:
217:     def ysd(i): return [col.sd() for col in i.cols.y]
218:
219: #
220: def stratify(src):
221:     all, klass = None, {}
222:     for n, row in enumerate(src):

```

## keys

```

217: if all:
218:     kl = row[all.cols.klass.at]
219:     here = klass[kl] = klass.get(kl, None) or all.clone()
220:     here.add(row)
221:     all.add(row)
222: else:
223:     all = Rows([row])
224:     return o(all=all, klass=klass)
225:
226: # Discretizations
227: # Use 'bins' to divide numeric data into ranges.
228: def bins(xy, iota=0, size=30):
229:     def merge(b4):
230:         j, tmp, n = 0, [], len(b4)
231:         while j < n:
232:             ((lo, _), ay) = a = b4[j]
233:             if j < n - 1:
234:                 ((_, hi), by) = b4[j + 1]
235:                 if cy := ay.merged(by):
236:                     a = ((lo, hi), cy)
237:                     j += 1
238:             tmp += [a]
239:             j += 1
240:         return merge(tmp) if len(tmp) < len(b4) else b4
241:
242: def divide(xy):
243:     bin = o(x=Num(), y=Sym())
244:     bins = [bin]
245:     for i, (x, y) in enumerate(xy):
246:         if bin.x.n >= size and x != b4:
247:             if i < len(xy)-size and bin.x.wide(iota):
248:                 bin = o(x=Num(), y=Sym())
249:                 bins += [bin]
250:                 bin.x.add(x)
251:                 bin.y.add(y)
252:                 b4 = x
253:             return bins
254:
255: xy = sorted(xy, key=first)
256: return merge([(bin.x.span(), bin.y) for bin in divide(xy)])
257:
258: # Learn class deltas
259: def contrasts(here, there, the):
260:     goal = {'best': lambda b, r: b**2/(b+r),
261:            'rest': lambda b, r: r**2/(b+r),
262:            'other': lambda b, r: 1/(b+r)} [the.GOAL]
263:
264: def like(d, kl):
265:     out = prior = (hs[kl] + the.K) / (n + the.K*2)
266:     for at, span in d.items():
267:         f = has.get((kl, at, span), 0)
268:         out *= (f + the.M*prior) / (hs[kl] + the.M)
269:     return out
270:
271: def top(a): return sorted(a, reverse=True, key=first)[the.TOP]
272: def val(d):
273:     for k, v in list(d.items()):
274:         if v == None: del d[k]
275:     return (goal(like(d, True), like(d, False)), d)
276:
277: has = {(klass, at, (lo, hi)): f
278:        for col1, col2 in zip(here.cols.x, there.cols.x)
279:        for klass, f, at, (lo, hi) in col1.bins(col2, the)}
280: n = len(here.rows) + len(there.rows)
281: hs = {True: len(here.rows), False: len(there.rows)}
282: solos = [val((at, span)) for (at, span) in
283:           set([(at, span) for (_, at, span) in has])]
284: ranges = {}
285: for _, d in top(solos):
286:     for k in d:
287:         if k not in ranges: ranges[k] = set([None])
288:         ranges[k].add(d[k])

```

```

289: for s, rule in top([(val(d) for d in dict_product(ranges))]):
290:     if rule:
291:         print(here.display(s, rule))
292:
293: # Misc utils
294: # string stuff
295: def color(end="\n", **kw):
296:     s, a, z = "", "\u001b[" + ";1m"
297:     c = dict(black=30, red=31, green=32, yellow=33,
298:             purple=34, pink=35, blue=36, white=37)
299:     for col, txt in kw.items():
300:         s = s+a + str(c[col]) + z+txt+"\033[0m"
301:     print(s, end=end)
302:
303: def mline(m): m += [["-"*len(str(x)) for x in m[-1]]]
304:
305: def printm(matrix):
306:     s = [str(e) for e in row] for row in matrix]
307:     lens = [max(map(len, col)) for col in zip('s')]
308:     fmt = ' | '.join('{{{:>{}}}}'.format(x) for x in lens)
309:     for row in [fmt.format(*row) for row in s]:
310:         print(row)
311:
312: # maths stuff
313: def r3(a): return [round(x, 3) for x in a]
314:
315: # dictionary stuff
316: def has(d, k): return d.get(k, 0)
317: def inc(d, k, n=1): tmp = d[k] = n + d.get(k, 0); return tmp
318:
319: def dict_product(d):
320:     keys = d.keys()
321:     for p in itertools.product(*d.values()):
322:         yield dict(zip(keys, p))
323:
324: # list stuff
325: def first(a): return a[0]
326: def last(a): return a[-1] # $Label(comment)$
327: def per(a, p=.5): return a[int(p*len(a))]
328:
329: # file stuff
330: def csv(f=None, sep=","):
331:     def prep(s): return re.sub(r'([\\n\\t\\r \\#\\.\\*\\', ', ', s)
332:     if f:
333:         with open(f) as fp:
334:             for s in fp:
335:                 if s := prep(s):
336:                     yield s.split(sep)
337:     else:
338:         for s in sys.stdin:
339:             if s := prep(s):
340:                 yield s.split(sep)
341:
342: # command-line stuff
343: def cli(use, txt, config):
344:     fmt = argparse.RawTextHelpFormatter
345:     used, p = {}, argparse.ArgumentParser(prog=use, description=txt,
346:                                           formatter_class=fmt)
347:     for k, (_, b4, h) in sorted(config.items()):
348:         k0 = k[0]
349:         used[k0] = c = k0 if k0 in used else k0.lower()
350:         if b4 == False:
351:             p.add_argument("-"+c, dest=k, default=False,
352:                           help=h,
353:                           action="store_true")
354:         else:
355:             p.add_argument("-"+c, dest=k, default=b4,
356:                           help=h + " [" + str(b4) + "]",
357:                           type=type(b4), metavar=k)
358:     return o(**p.parse_args().__dict__)
359:
360: # Unit tests

```

```

361: class Eg:
362:     def ls(the):
363:         "list all examples."
364:         print("\nexamples:")
365:         for k, f in vars(Eg).items():
366:             if k[0] != "_":
367:                 print(f" {k:<13} {f.__doc__}")
368:
369:     def _fail(the):
370:         "testing failure"
371:         assert False, "failing"
372:
373:     def data(the, file="../data/vote.csv"):
374:         "simple load of data into a table"
375:         r = Rows(csv(file))
376:         assert 435 == len(r.rows)
377:         assert 195 == r.cols.all[1].has['y']
378:
379:     def nclasses(the, file="../data/diabetes.csv", kl="positive"):
380:         "read data with nclasses"
381:         rs = stratify(csv(file))
382:         assert 2 == len(rs.klass)
383:         assert 268 == len(rs.klass[kl].rows)
384:         assert 768 == len(rs.all.rows)
385:         assert 3.90625 == rs.klass[kl].cols.all[0].sd()
386:
387:     def bins(the, file="../data/diabetes.csv",
388:             k1="positive", k2="negative"):
389:         "discretize some data"
390:         rs = stratify(csv(file))
391:         bins1(rs.klass[k1], rs.klass[k2], the)
392:
393:     def bestrest(the, file="../data/auto93.csv"):
394:         "discretize some multi-goal data"
395:         r = Rows(csv(file))
396:         goods, bads = r.best(the)
397:         bins1(goods, bads, the)
398:
399:     def con1(the, file="../data/auto93.csv"):
400:         "discretize some multi-goal data"
401:         r = Rows(csv(file))
402:         goods, bads = r.best(the)
403:         contrasts(goods, bads, the)
404:
405:     def con2(the, file="../data/china.csv"):
406:         "discretize china "
407:         Eg.con1(the, file)
408:
409:     def bins1(goods, bads, the):
410:         for good, bad in zip(goods.cols.x, bads.cols.x):
411:             bins = sorted(good.bins(bad, the))
412:             if len(bins) > 1:
413:                 print(f"\n{good.txt}")
414:                 for bin in bins:
415:                     print("\t", bin)
416:
417: # Main program
418: def main(the):
419:     def run(fun, fails, the):
420:         s = f" {fun.__name__:<12}"
421:         if the.UNSAFE:
422:             print("unsafe mode:")
423:             fun(copy.deepcopy(the))
424:             sys.exit()
425:         try:
426:             fun(copy.deepcopy(the))
427:             random.seed(the.SEED)
428:             color(green=(chr(10003) + s), white=fun.__doc__)
429:         except Exception as err:
430:             fails = fails + 1
431:             color(red=(chr(10007) + s), white=str(err))
432:     return fails

```

```
433: # _____
434: fails = 0
435: if the.XAMPLE == "all":
436:     for k, f in vars(Eg).items():
437:         if k[0] != "_" and k != "ls":
438:             fails = run(f, fails, the)
439: else:
440:     if the.XAMPLE and the.XAMPLE in vars(Eg):
441:         f = vars(Eg)[the.XAMPLE]
442:         if the.XAMPLE == "ls":
443:             f(the)
444:         else :
445:             fails = run(f, fails, the)
446:     sys.exit(fails)
447:
448:
449: if __name__ == "__main__":
450:     main(cli("/keys", __doc__, config()))
```