```python
5:   #!/usr/bin/env python3.9
6:   # vim: ts=2 sw=2 sts=2 et :
7:   # autopep8: ignore E20,E401,E226,E302,E41
8:   import re, sys, math, argparse, itertools
9:   from argparse import ArgumentParser as parse
10:  from argparse import RawTextHelpFormatter as textual
11:  Float = Str = Int = Bool = lambda *l: l[0]
12:
13:  def keys(
14:    BINS   : Float("bins are of size n**BINS") = .5,
15:    COLS   : Str("columns to use for inference") = "x",
16:    DATA   : Str("where to read data") = "../data/auto2.csv",
17:    EPSILON: Float("small = sd**EPSILON") = .3,
18:    FAR    : Float("where to look for far things") = .9,
19:    GOAL   : Str("learning goals: best|rest|other") = "best",
20:    K      : Int("bayes low class frequency hack") = 2,
21:    M      : Int("bayes low range frequency hack") = 1,
22:    P      : Int("distance calculation exponent") = 2,
23:    SAMPLE : Int("#samples to find far things?") = 20,
24:    VERBOSE: Bool("set verbose") = False,
25:    TOP    : Int("focus on this many") = 20,
26:    XAMPLE : Str("egs: '-x ls' lists all, '-x all' runs all") = "" ):
27:    """
28:       ,-_|\    (c) Tim Menzies, 2021, unlicense.org.
29:      /    \   The delta between things is
30:      \_,-._*   simpler than the things.
31:          v """
32:
33:  GOAL = {'best' : lambda b, r: b**2/b+r,
34:          'rest' : lambda b, r: r**2/(b+r),
35:          'other': lambda b, r: 1/(b+r)   }[GOAL]
36:
37:  # -----------------------------------------------------------
38:  # Storing columns in 'Col', 'Skip', 'Sym', 'Num'.
39:  class Col(o):
40:    def __init__(i, at=0, txt="", inits=[]):
41:      i.n, i.at, i.txt = 0, at, txt
42:      i.w = -1 if "-" in txt else 1
43:      [i.add(x) for x in inits]
44:
45:    def add(i, x, n=1):
46:      if x != "?": i.n += 1; x = i.add1(x, n)
47:      return x
48:  # ----------------
49:  class Skip(Col):
50:    def add1(i, x, n=1): return x
51:  # ----------------
52:  class Sym(Col):
53:    def __init__(i, **kw): i.has = {}; super().__init__(**kw)
54:
55:    def add1(i, x, n=1): inc(i.has, x, n); return x
56:
57:    def bins(i, j):
58:      for k in (i.has | j.has):
59:        yield i.has.get(k, 0), True, (i.at, (k, k))
60:        yield j.has.get(k, 0), False, (j.at, (k, k))
61:
62:    def dist(i, x, y): return 0 if x == y else 1
63:
64:    def ent(i):
65:      return sum(-v/i.n * math.log(v/i.n) for v in i.has.values())
66:
67:    def merge(i, j):
68:      k = Sym(at=i.at, txt=i.txt)
69:      [k.add(x, n) for has in (i.has, j.has) for x, n in has.items()]
70:      return k
71:
72:    def merged(i, j):
73:      k = i.merge(j)
74:      e1, n1, e2, n2, e, n = i.ent(), i.n, j.ent(), j.n, k.ent(), k.n
75:      if e1 + e2 < 0.01 or e * .95 < n1 / n * e1 + n2 / n * e2:
76:        return k
77:  # ----------------
78:  class Num(Col):
79:    def __init__(i, **kw):
80:      i._all, i.ok = [], False
81:      super().__init__(**kw)
82:
83:    def add1(i, x, n):
84:      x, i.ok = float(x), False
85:      for _ in range(n): i._all += [x]
86:      return x
87:
88:    def all(i):
89:      if not i.ok: i.ok = True; i._all = sorted(i._all)
90:      return i._all

91:
92:    def bins(i, j):
93:      xy = [(z, True) for z in i._all]+[(z, False) for z in j._all]
94:      eps = EPSILON * (i.n*i.sd() + j.n*j.sd()) / (i.n + j.n)
95:      for ((lo, hi),s) in bins(xy,epsilon=eps,size=len(xy)**BINS):
96:        for klass, n in s.has.items():
97:          yield n, klass, (i.at, (lo, hi))
98:
99:    def dist(i, x, y):
100:     if   x == "?": y = i.norm(y); x = 1 if y < 0.5 else 0
101:     elif y == "?": x = i.norm(x);y = 1 if x < 0.5 else 0
102:     else         : x, y = i.norm(x), y.norm(y)
103:     return abs(x-y)
104:
105:    def norm(i, x):
106:      if x == "?": return x
107:      a = i.all()
108:      return max(0, min(1, (x-first(a)/(last(a)-first(a)+1E-32)))
109:
110:    def sd(i) : return (per(i.all(), .9) - per(i.all(), .1))/2.56
111:    def span(i) : return (first(i.all()), last(i.all()))
112:    def wide(i, n=0): return last(i.all()) - first(i.all()) >= n
113:
114:  #-----------------------------------------------------------
115:  # Data is in 'Row's which, in turn, are stored  in 'Table's.
116:  class Row(o):
117:    def __init__(i, lst, tab=None): i.tab, i.cells = tab, lst
118:
119:    def dist(i, j):
120:      d = n = 1E-32
121:      for col in i.tab.cols[COLS]:
122:        n += 1
123:        x, y = i.cells[at], j.cells[at]
124:        d += 1 if x == "?" and y == "?" else col.dist(x, y) ** P
125:      return (d/n) ** (1/P)
126:
127:    def far(i, rows):
128:      tmp = [(dist(i, j), j) for _ in range(SAMPLE)]
129:      return per(sorted(tmp, key=first), FAR)
130:
131:  # ----------------
132:  class Table(o):
133:    def __init__(i, inits=[]):
134:      i.rows = []
135:      i.cols = o(all=[], names=[], x=[], y=[], klass=None)
136:      [i.add(x) for x in inits]
137:
138:    def add(i, a): i.data(a) if i.cols.names else i.header(a)
139:    def clone(i, inits=[]): return Table([i.cols.names] + inits)
140:
141:    def data(i, a):
142:      a = a.cells if type(a) == Row else a
143:      a = [col.add(a[col.at]) for col in i.cols.all]
144:      i.rows += [Row(a, tab=i)]
145:
146:    def header(i, a):
147:      i.cols.names = a
148:      for at, x in enumerate(a):
149:        new = Skip if i.skipp(x) else (Num if i.nump(x) else Sym)
150:        new = new(at=at, txt=x)
151:        i.cols.all += [new]
152:        if not i.skipp(x):
153:          i.cols["y" if i.yp(x) else "x"] += [new]
154:          if i.klassp(x):
155:            i.cols.klass = new
156:
157:    def klassp(i, x): return "!" in x
158:    def nump(i, x): return x[0].isupper()
159:    def skipp(i, x): return "?" in x
160:    def yp(i, x): return "-" in x or "+" in x or i.klassp(x)
161:
162:  # ----------------
163:  def stratify(src):
164:    all, klass = None,{}
165:    for n,row in enumerate(src):
166:      if all:
167:        kl   = row[all.cols.klass.at]
168:        here = klass[kl] = klass.get(kl,None) or all.clone()
169:        here.add(row)
170:        all.add(row)
171:      else:
172:        all = Table([row])
173:    return o(all=all, klass=klass)
174:
175:  #-----------------------------------------------------------
176:  # Use 'bins' to divide numeric data into ranges.
```

```
177:   def bins(xy, epsilon=0, size=30):
178:     def merge(b4):
179:       j, tmp, n = 0, [], len(b4)
180:       while j < n:
181:         a = b4[j]
182:         if j < n - 1:
183:           b = b4[j + 1]
184:           print("\na",a[1])
185:           print("b",b[1])
186:           if cy := a[1].merged(b[1]):
187:             print("c",cy)
188:             a = ((a[0][0], b[0][1]), cy)
189:             j += 1
190:         tmp += [a]
191:         j += 1
192:       return merge(tmp) if len(tmp) < len(b4) else b4
193:
194:     def divide(xy):
195:       bin = o(x=Num(), y=Sym())
196:       bins = [bin]
197:       for i, (x, y) in enumerate(xy):
198:         if bin.x.n >= size:
199:           if x != b4 and i < len(xy)-size and bin.x.wide(epsilon):
200:             bin = o(x=Num(), y=Sym())
201:             bins += [bin]
202:         bin.x.add(x)
203:         bin.y.add(y)
204:         b4 = x
205:       return bins
206:
207:     return merge([(bin.x.span(), bin.y)
208:                   for bin in divide(sorted(xy, key=first))])
209:
210:   #---------------------------------------------------------------
211:   # Report ranges that are most different in two classes.
212:   def contrasts(here, there, t):
213:     def like(d, kl):
214:       out = prior = (hs[kl] + K) / (n + K*2)
215:       for at, span in d.items():
216:         f = has.get((kl, (at, span)), 0)
217:         out *= (f + M*prior) / (hs[kl] + M)
218:       return out
219:
220:     def val(d): return GOAL(like(d, True), like(d, False)), d
221:     def top(a): return sorted(a, reversed=True, key=first)[:TOP]
222:
223:     has = {(kl, (at, (lo, hi))): f
224:            for col1, col2 in zip(here.cols.x, there.cols.x)
225:            for f, kl, (at, (lo, hi)) in col1.bins(col2)}
226:     n = len(here.rows, there.rows)
227:     hs = {True: len(here.rows), False: len(there.rows)}
228:     solos = [val(dict(at=x)) for at, x in set([z for _, z in has])]
229:     ranges = {}
230:     for _, d in top(solos):
231:       for k in d:
232:         ranges[k] = ranges.get(k, set()).add(d[k])
233:     for rule in top([val(d) for d in dict_product(ranges)]):
234:       print(rule)
235:
236:   #---------------------------------------------------------------
237:   # Unit tests.
238:   class Eg:
239:     def ls():
240:       "list  all examples."
241:       print("\nexamples:")
242:       for k, f in vars(Eg).items():
243:         if k[0] != "_":
244:           print(f"  {k:<13} {f.__doc__}")
245:
246:     def data(file="../data/vote.csv"):
247:       "simple load of data into  a table"
248:       t = Table(csv(file))
249:       assert(435 == len(t.rows))
250:       assert(195 == t.cols.all[1].has['y'])
251:
252:     def nclasses(file="../data/diabetes.csv", kl="positive"):
253:       ts = stratify(csv(file))
254:       assert(2   == len(ts.klass))
255:       assert(268 == len(ts.klass[kl].rows))
256:       assert(768 == len(ts.all.rows))
257:
258:     def bins(file="../data/diabetes.csv",
259:              k1= "positive", k2= "negative"):
260:       ts = stratify(csv(file))
261:       goods, bads = ts.klass[k1], ts.klass[k2]
262:       for good,bad in zip(goods.cols.all, bads.cols.all):
263:         print(f"\n{good.at}")
264:         [print(f"\t{x}") for x in good.bins(bad)]
265:
266:   # ------------------------------------------------------------
267:   # main program for keys
268:   if XAMPLE == "all":
269:     for k, f in vars(Eg).items():
270:       if k[0] != "_": print("\n"+k); f()
271:   else:
272:     if XAMPLE and XAMPLE in vars(Eg): vars(Eg)[XAMPLE]()
273:
274:   ###############################################
275:   # things that don't use the config vars
276:   # dictionaries
277:   def has(d, k): return d.get(k, 0)
278:   def inc(d, k, n=1): tmp = d[k] = n + d.get(k, 0); return tmp
279:   def public(d): return {k:v for k, v in d.items() if k[0] != "_"}
280:
281:   def dict_product(d):
282:     keys = d.keys()
283:     for p in itertools.product(*d.values()):
284:       yield dict(zip(keys, p))
285:
286:   # lists
287:   def first(a): return a[0]
288:   def last(a): return a[-1]  ####
289:   def per(a, p=.5): return a[int(p*len(a))]
290:
291:   # objects
292:   class o(object):
293:     def __init__(i, **k): i.__dict__.update(**k)
294:     def __getitem__(i, k): return i.__dict__[k]
295:     def __repr__(i):
296:       return i.__class__.__name__+str(public(i.__dict__))
297:     def __setitem__(i, k, v): i.__dict__[k] = v
298:
299:   # read csv files
300:   def csv(f=None, sep=","):
301:     def prep(s): return re.sub(r'([\n\t\r ]|#.*)', '', s)
302:     if f:
303:       with open(f) as fp:
304:         for s in fp:
305:           if s := prep(s): yield s.split(sep)
306:     else:
307:       for s in sys.stdin:
308:         if s := prep(s): yield s.split(sep)
309:
310:   # Drive command line flags  from function annocations.
311:   def cli(f):
312:     p = parse(prog="./"+f.__name__, description=f.__doc__,
313:           formatter_class=textual)
314:     for (k, h),b4 in zip(
315:             list(f.__annotations__.items()),f.__defaults__):
316:       if b4 == False:
317:         p.add_argument("-"+(k[0].lower()), dest=k, help=h,
318:               default=False, action="store_true")
319:       else:
320:         p.add_argument("-"+(k[0].lower()), dest=k, default=b4,
321:               help=h+" ["+str(b4)+"]", type=type(b4),
322:               metavar=k)
323:     f(**p.parse_args().__dict__)
324:
325:   # Start up.
326:   if __name__ == "__main__": cli(keys)
327:
328:   # gs -dBATCH -dNOPAUSE -q -sDEVICE=pdfwrite -dAutoRotate
       Pages=/None -sOutputFile=finished.pdf  file1.pdf file2.pdf
329:
```