

keys

```

5: #!/usr/bin/env python3.9
6: # vim: ts=2 sw=2 sts=2 et :
7: # autoupep8: ignore E20,E401,E226,E302,E41
8: import re, sys, argparse, itertools
9: from argparse import ArgumentParser as parse
10: from argparse import RawTextHelpFormatter as textual
11: Float = Str = Int = Bool = lambda "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)
36:            }[GOAL]
37:
38:     # -----
39:     class Col(o):
40:         "Storing columns in 'Col', 'Skip', 'Sym', 'Num'."
41:         def __init__(i, at=0, txt="", inits=[]):
42:             i.n, i.at, i.txt = 0, at, txt
43:             i.w = -1 if "-" in txt else 1
44:             [i.add(x) for x in inits]
45:
46:         def add(i, x, n=1):
47:             if x != "?": i.n += 1; 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): i.has = {}; super().__init__(**kw)
55:
56:             def add1(i, x, n=1): inc(i.has, x, n); return x
57:
58:         def bins(i, j):
59:             for k in (i.has | j.has):
60:                 yield i.has.get(k, 0), True, (i.at, (k, k))
61:                 yield j.has.get(k, 0), False, (j.at, (k, k))
62:
63:         def dist(i, x, y): return 0 if x == y else 1
64:
65:         def ent(i):
66:             return sum(-v/i.n * math.log(v/i.n) for v in i.has.values())
67:

```

```

68: def merge(i, j):
69:     k = Sym(at=i.at, txt=i.txt)
70:     [k.add(x, n) for has in (i.has, j.has) for x, n in has.items()]
71:     return k
72:
73: def merged(i, j):
74:     k = i.merge(j)
75:     e1, n1, e2, n2, e, n = i.ent(), i.n, j.ent(), j.n, k.ent(), k.n
76:     if e1 + e2 < 0.01 or e * .95 < n1 / n * e1 + n2 / n * e2:
77:         return k
78:     # -----
79:     class Num(Col):
80:         def __init__(i, **kw):
81:             i._all, i.ok = [], False
82:             super().__init__(**kw)
83:
84:         def add1(i, x, n):
85:             x, i.ok = float(x), False
86:             for _ in range(n): i._all += [x]
87:             return x
88:
89:         def all(i):
90:             if not i.ok: i.ok = True; i._all = sorted(i._all)
91:             return i._all
92:
93:         def bins(i, j):
94:             xy = [(z, True) for z in i._all] + [(z, False) for z in j._all]
95:             eps = EPSILON * (i.n*i.sd() + j.n*j.sd()) / (i.n + j.n)
96:             for (lo, hi), sym in bins(xy, epsilon=eps, enough=len(xy)**BINS):
97:                 for klass, n in sym.has.items():
98:                     yield n, klass, (i.at, (lo, hi))
99:
100:         def dist(i, x, y):
101:             if x == "?": y = i.norm(y); x = 1 if y < 0.5 else 0
102:             elif y == "?": x = i.norm(x); y = 1 if x < 0.5 else 0
103:             else: x, y = i.norm(x), y.norm(y)
104:             return abs(x-y)
105:
106:         def norm(i, x):
107:             if x == "?": return x
108:             a = i.all()
109:             return max(0, min(1, (x-first(a))/(last(a)-first(a)+1E-32)))
110:
111:         def sd(i): return (per(i.all(), .9) - per(i.all(), .1))/2.56
112:         def span(i): return (first(i.all()), last(i.all()))
113:         def wide(i, n=0): return last(i.all()) - first(i.all()) >= n
114:
115:     # -----
116:     class Row(o):
117:         "Data is held in 'Row's which, in turn, are stored in 'Table's"
118:         def __init__(i, lst, tab=None): i.tab, i.cells = tab, lst
119:
120:         def dist(i, j):
121:             d = n = 1E-32
122:             for col in i.tab.cols[COLS]:
123:                 n += 1
124:                 x, y = i.cells[at], j.cells[at]
125:                 d += 1 if x == "?" and y == "?" else col.dist(x, y) ** P
126:             return (d/n) ** (1/P)
127:
128:         def far(i, rows):
129:             tmp = [(dist(i, j), j) for _ in range(SAMPLE)]
130:             return per(sorted(tmp, key=first), FAR)

```

```

131:
132: # -----
133: class Table(o):
134:     def __init__(i, inits=[]):
135:         i.rows = []
136:         i.cols = o(all=[], names=[], x=[], y=[], klass=None)
137:         [i.add(x) for x in inits]
138:
139:     def add(i, a): i.data(a) if i.cols.names else i.header(a)
140:     def clone(i, inits=[]): return Table([i.cols.names] + inits)
141:
142:     def data(i, a):
143:         a = a.cells if type(a) == Row else a
144:         a = [col.add(a[col.at]) for col in i.cols.all]
145:         i.rows += [Row(a, tab=i)]
146:
147:     def header(i, a):
148:         i.cols.names = a
149:         for at, x in enumerate(a):
150:             new = Skip if i.skipp(x) else (Num if i.nump(x) else Sym)
151:             new = new(at=at, txt=x)
152:             i.cols.all += [new]
153:             if not i.skipp(x):
154:                 i.cols["y" if i.yp(x) else "x"] += [new]
155:             if i.klassp(x):
156:                 i.cols.klass = new
157:
158:     def klassp(i, x): return "!" in x
159:     def nump(i, x): return x[0].isupper()
160:     def skipp(i, x): return "?" in x
161:     def yp(i, x): return "-" in x or "+" in x or i.klassp(x)
162:
163:     # -----
164:     def bins(xy, epsilon=0, enough=30):
165:         "Use 'bins' to divide numeric data into ranges."
166:         def merge(b4):
167:             j, tmp, n = 0, [], len(b4)
168:             while j < n:
169:                 a = b4[j]
170:                 if j < n - 1:
171:                     b = b4[j + 1]
172:                     if cy := a.y.merged(b.y):
173:                         a = o(x=(a.x[0], b.x[1]), y=cy)
174:                     j += 1
175:                 tmp += [a]
176:                 j += 1
177:             return merge(tmp) if len(tmp) < len(b4) else b4
178:
179:     def divide(xy):
180:         xy = sorted(xy, key=first)
181:         bin = o(x=Num(), y=Sym())
182:         bins = [bin]
183:         for i, (x, y) in enumerate(xy):
184:             if bin.x.n >= enough:
185:                 if x != b4 and i < len(xy)-enough and bin.x.wide(epsilon):
186:                     bin = o(x=Num(), y=Sym())
187:                     bins += [bin]
188:                 bin.x.add(x)
189:                 bin.y.add(y)
190:                 b4 = x
191:             return bins
192:
193:     return merge(divide([o(bin.x.span(), y=bin.y) for bin in bins]))

```

keys

```

194:
195: # -----
196: def contrasts(here, there, t):
197:     "Report the ranges that are most different in two classes."
198:     def like(d, kl):
199:         out = prior = (hs[kl] + K) / (n + K*2)
200:         for at, span in d.items():
201:             f = has.get(kl, (at, span)), 0)
202:             out *= (f + M*prior) / (hs[kl] + M)
203:         return out
204:
205: def val(d): return GOAL(like(d, True), like(d, False)), d
206: def top(a): return sorted(a, reversed=True, key=first)[:TOP]
207:
208: has = {(kl, (at, (lo, hi))): f
209:         for col1, col2 in zip(here.cols.x, there.cols.x)
210:         for f, kl, (at, (lo, hi)) in col1.bins(col2)}
211: n = len(here.rows, there.rows)
212: hs = {True: len(here.rows), False: len(there.rows)}
213: solos = [val(dict(at=x)) for at, x in set([z for _, z in has])]
214: ranges = {}
215: for _, d in top(solos):
216:     for k in d:
217:         ranges[k] = ranges.get(k, set()).add(d[k])
218: for rule in top([val(d) for d in dict_product(ranges)]):
219:     print(rule)
220:
221: # -----
222: class Eg:
223:     "Unit tests."
224:     def ls():
225:         "list all examples."
226:         print("\nexamples:")
227:         for k, f in vars(Eg).items():
228:             if k[0] != "_":
229:                 print(f" {k:<13} {f.__doc__}")
230:
231: def data(file="./data/vote.csv", goal="democrat"):
232:     "simple load of data into a table"
233:     t = Table(csv(file))
234:     assert(435 == len(t.rows))
235:     assert(195 == t.cols.all[1].has['y'])
236:
237: def clone(file="./data/diabetes.csv", k="positive"):
238:     "discretize test"
239:     t = Table(csv(file))
240:     kl = t.cols.klass.at
241:     u, v = t.clone(), t.clone()
242:     [(u if k == row.cells[kl] else v).add(row) for row in t.rows]
243:     good, bad = u.cols.x[1], v.cols.x[1]
244:     #print(good.all()[::25], len(good.all()), good.sd())
245:     print(bad.all()[::25], len(bad.all()), bad.sd())
246:
247: def bins(file="./data/diabetes.csv", k="positive"):
248:     "discretize test"
249:     t = Table(csv(file))
250:     kl = t.cols.klass.at
251:     u, v = t.clone(), t.clone()
252:     print("=====")
253:     for row in t.rows:
254:         (u if k == row.cells[kl] else v).add(row)
255:     good, bad = u.cols.x[1], v.cols.x[1]
256:     for x in good.bins(bad):

```

```

257:     print(good.at, x)
258:
259: # -----
260: # main program for keys
261: if XAMPLE == "all":
262:     for k, f in vars(Eg).items():
263:         if k[0] != "_": print("\n"+k); f()
264: else:
265:     if XAMPLE and XAMPLE in vars(Eg): vars(Eg)[XAMPLE]()
266:
267: #####
268: # things that don't use the config vars
269: # dictionaries
270: def has(d, k): return d.get(k, 0)
271: def inc(d, k, n=1): tmp = d[k] = n + d.get(k, 0); return tmp
272: def public(d): return {k: v for k, v in d.items() if k[0] != "_"}
273:
274: def dict_product(d):
275:     keys = d.keys()
276:     for p in itertools.product(*d.values()):
277:         yield dict(zip(keys, p))
278:
279: # -----
280: # lists
281: def first(a): return a[0]
282: def last(a): return a[-1]
283: def per(a, p=.5): return a[int(p*len(a))]
284:
285: # -----
286: class o(object):
287:     "objects"
288:     def __init__(i, **k): i.__dict__.update(**k)
289:     def __getitem__(i, k): return i.__dict__[k]
290:     def __repr__(i): return i.__class__.__name__+str(public(i.__dict__))
291:     def __setitem__(i, k, v): i.__dict__[k] = v
292:
293: # -----
294: def csv(f=None, sep=","):
295:     "read csv files"
296:     def prep(s): return re.sub(r'([\n\r ]|#.*)', '', s)
297:     if f:
298:         with open(f) as fp:
299:             for s in fp:
300:                 if s := prep(s): yield s.split(sep)
301:     else:
302:         for s in sys.stdin:
303:             if s := prep(s): yield s.split(sep)
304:
305: # -----
306: def cli(f):
307:     "Drive command line flags from function annotations."
308:     p = parse(prog="./"+f.__name__, description=f.__doc__,
309:              formatter_class=Textual)
310:     for (k, h), b4 in zip(
311:         list(f.__annotations__.items()), f.__defaults__):
312:         if b4 == False:
313:             p.add_argument("-"+k[0].lower(), dest=k, help=h,
314:                           default=False, action="store_true")
315:         else:
316:             p.add_argument("-"+k[0].lower(), dest=k, default=b4,
317:                           help=h+" ["+str(b4)+"]", type=type(b4),
318:                           metavar=k)
319:     f(**p.parse_args().__dict__)

```

```

320:
321: # -----
322: if __name__ == "__main__": cli(keys)

```