

keys

```

5: #!/usr/bin/env python3.9
6: # vim: ts=2 sw=2 sts=2 et :
7: # autpep8: ignore E20,E401,E226,E301, 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: [l]
12:
13: #####
14: def keys(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:     """
29:     ,_/_\ (c) Tim Menzies, 2021, unlicense.org.
30:     /_/_\ The delta between things is
31:     \_/_\_ simpler than the things.
32:     v """"
33:
34:     GOAL = {'best': lambda b, r: b**2/b+r,
35:            'rest': lambda b, r: r**2/(b+r),
36:            'other': lambda b, r: 1/(b+r)
37:            }[GOAL]
38:
39:     # -----
40:     # Modeling columns with Col, Skip, Sym, and Num
41:     class Col(o):
42:         def __init__(i, at=0, txt="", inits=[]):
43:             i.n, i.at, i.txt = 0, at, txt
44:             i.w = -1 if "-" in txt else 1
45:             [i.add(x) for x in inits]
46:
47:         def add(i, x, n=1):
48:             if x != "?": i.n += 1; x = i.add1(x, n)
49:             return x
50:
51:     # -----
52:     class Skip(Col):
53:         def add1(i, x, n=1): return x
54:
55:     # -----
56:     class Sym(Col):
57:         def __init__(i, **kw): i.has = {}; super().__init__(**kw)
58:
59:         def add1(i, x, n=1): inc(i.has, x, n); return x
60:
61:         def bins(i, j):
62:             for k in (i.has | j.has):
63:                 yield i.has.get(k, 0), True, (i.at, (k, k))
64:                 yield j.has.get(k, 0), False, (j.at, (k, k))
65:
66:         def dist(i, x, y): return 0 if x == y else 1
67:
68:         def ent(i):
69:             return sum(-v/i.n*math.log(v/i.n) for v in i.has.values())
70:
71:

```

```

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: # -----
80: class Num(Col):
81:     def __init__(i, **kw):
82:         i._all, i.ok = [], False
83:         super().__init__(**kw)
84:
85:     def add1(i, x, n):
86:         x, i.ok = float(x), False
87:         for _ in range(n): i._all += [x]
88:         return x
89:
90:     def all(i):
91:         if not i.ok: i.ok = True; i._all = sorted(i._all)
92:         return i._all
93:
94:     def bins(i, j):
95:         xy = [(z, True) for z in i._all]+[(z, False) for z in j._all]
96:         eps = EPSILON * (i.n*i.sd() + j.n*j.sd()) / (i.n + j.n)
97:         for ((lo, hi), sym) in bins(xy, epsilon=eps, enough=len(xy)**BINS):
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), i.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: # Modeling tables with Row and Table and (as defined above) Col
117: class Row(o):
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: 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.y(x) else "x"] += [new]
154:             if i.klass(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: # Discretization of numeric columns with merge and divide.
164: def bins(xy, epsilon=0, enough=30):
165:     j, tmp, n = 0, [], len(b4)
166:     while j < n:
167:         a = b4[j]
168:         if j < n - 1:
169:             b = b4[j + 1]
170:             if cy := a.y.merged(b.y):
171:                 a = o(x=(a.x[0], b.x[1]), y=cy)
172:             j += 1
173:             tmp += [a]
174:             j += 1
175:         return merge(tmp) if len(tmp) < len(b4) else b4
176:
177:     def divide(xy):
178:         xy = sorted(xy, key=first)
179:         bin = o(x=Num(), y=Sym())
180:         bins = [bin]
181:         for i, (x, y) in enumerate(xy):
182:             if bin.x.n >= enough:
183:                 if x != b4 and i < len(xy)-enough and bin.x.wide(epsilon):
184:                     bin = o(x=Num(), y=Sym())
185:                     bins += [bin]
186:                 bin.x.add(x)
187:                 bin.y.add(y)
188:                 b4 = x
189:             return bins
190:
191:     return merge(
192:         divide([o(bin.x.span(), y=bin.y) for bin in bins])) #

```

keys

```

193: # -----
194: # Learning the delta between two classes.
195: def contrasts(here, there, t):
196:     def like(d, kl):
197:         out = prior = (hs[kl] + K) / (n + K*2)
198:         for at, span in d.items():
199:             f = has.get(kl, (at, span), 0)
200:             out *= (f + M*prior) / (hs[kl] + M)
201:         return out
202:
203:     def val(d): return GOAL(like(d, True), like(d, False)), d
204:     def top(a): return sorted(a, reversed=True, key=first)[:TOP]
205:
206:     has = {(kl, (at, (lo, hi))): f
207:             for col1, col2 in zip(here.cols.x, there.cols.x)
208:             for f, kl, (at, (lo, hi)) in col1.bins(col2)}
209:     n = len(here.rows, there.rows)
210:     hs = {True: len(here.rows), False: len(there.rows)}
211:     solos = [val(dict(at=x)) for at, x in set([z for _, z in has])]
212:     ranges = {}
213:     for _, d in top(solos):
214:         for k in d:
215:             ranges[k] = ranges.get(k, set()).add(d[k])
216:     for rule in top([val(d) for d in dict_product(ranges)]):
217:         print(rule)
218: # -----
219: # Unit and system tests
220: class Eg:
221:     def ls():
222:         "list all examples."
223:         print("\nexamples:")
224:         for k, f in vars(Eg).items():
225:             if k[0] != "_": print(f" {k:<13} {f.__doc__}")
226:
227:     def data(file="./data/vote.csv", goal="democrat"):
228:         "simple load of data into a table"
229:         t = Table(csv(file)); print(1)
230:         assert(435 == len(t.rows))
231:         assert(195 == t.cols.all[1].has['y'])
232:
233:     def clone(file="./data/diabetes.csv", k="positive"):
234:         "discretize test"
235:         t = Table(csv(file))
236:         kl = t.cols.klass.at
237:         u, v = t.clone(), t.clone()
238:         [(u if k == row.cells[kl] else v).add(row) for row in t.rows]
239:         good, bad = u.cols.x[1], v.cols.x[1]
240:         #print(good.all()[::25], len(good.all()), good.sd())
241:         print(bad.all()[::25], len(bad.all()), bad.sd())
242:
243:     def bins(file="./data/diabetes.csv", k="positive"):
244:         "discretize test"
245:         t = Table(csv(file))
246:         kl = t.cols.klass.at
247:         u, v = t.clone(), t.clone()
248:         print("=====")
249:         for row in t.rows:
250:             (u if k == row.cells[kl] else v).add(row)
251:         good, bad = u.cols.x[1], v.cols.x[1]
252:         for x in good.bins(bad):
253:             print(good.at, x)
254: # -----

```

```

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

```