

```

1 #!/usr/bin/env python3.8
2 # vim: ts=2 sw=2 sts=2 et :
3 import re, sys, copy, math, random
4 from typing import Any
5 #
6 # 000b
7 #
8
9 big = sys.maxsize
10
11 is_klass = lambda x: "!" in x
12 is_less = lambda x: x[-1] == "-"
13 is_skip = lambda x: x[-1] == "."
14 is_num = lambda x: x[0].isupper()
15 is_goal = lambda x: "+" in x or "-" in x or is_klass(x)
16
17 def atom(x):
18     try: return float(x)
19     except: return x
20
21
22 def rows(f):
23     with open(f) as fp:
24         for line in fp:
25             line = re.sub(r'([\n\r\v\ ]|#|*)', ' ', line)
26             if line: yield [atom(cell.strip()) for cell in line.split(",")]
27
28 class o(object):
29     def __init__(i, **d): i.__dict__.update(**d)
30     def repr__(i):
31         pre = i.__class__.__name__ if isinstance(i,o) else ""
32         return pre+'['+(' '.join([f"{k}:{v}" for k, v in
33             sorted(i.__dict__.items()) if str(k)[0] != "_"]))+"]"
34
35 def merge(b4:list) -> list:
36     j,n,now = -1,len(b4),[]
37     while j < n-1:
38         j += 1
39         a = b4[j]
40         if j < n-2:
41             if merged := a.merge(b4[j+1]):
42                 a = merged
43                 j += 1
44             now += [a]
45     return b4 if len(now)==len(b4) else merge(now)
46 #
47 # eol
48 #
49 #
50
51 class Col(o):
52     def __init__(i, at=0, txt=""): i.n,i.at,i.txt = 0,at,txt
53     def add(i,x) : return x
54     def dist(i,x,y): return 1 if x=="?" and y=="?" else i.dist1(x,y)
55 #
56 # skip
57 #
58 #
59 #
60
61 class Skip(Col): ...
62
63 # eyn
64 #
65 #
66
67 class Sym(o):
68     def __init__(i,**kw): super().__init__(**kw); i.has = {}
69
70     def add(i, x, inc=1):
71         if x!="?": i.has[x] = inc + i.has.get(x,0)
72
73     def dist1(i,x,y): return 0 if x==y else 1
74
75     def div(i):
76         p = lambda x: x / (1E-31 + i.n)
77         return sum( -p(x)*math.log(p(x),2) for x in i.has.values() )
78
79     def merge(i,j):
80         k = Sym(at=i.at, txt=i.txt)
81         for x,n in i.has.items(): k.add(x,n)
82         for x,n in j.has.items(): k.add(x,n)
83         return k
84
85     def spans(i,j, _bins, out):
86         xys = [(x,"this",n) for x,n in i.has.items()] + [
87             (x,"that",n) for x,n in j.has.items()]
88         one, last = None,None
89         all = []
90         for x,y,n in sorted(xys, key=first):
91             if x != last:
92                 last = x
93                 one = Span(i, x, x)
94                 all += [one]
95             one.add(x,y,n)
96         if len(all) > 1 : out += all
97 #
98 # 0000
99 #
100
101 class Num(o):
102     def __init__(i,**kw):
103         super().__init__(**kw)
104         i.w, i.lo, i.hi = 1, big, -big
105         if i.txt and is_less(i.txt): i.w = -1
106
107     def add(i,x):
108         if x!="?": i.hi = max(i.hi,x); i.lo = min(i.lo,x)
109
110     def dist(i,x,y):
111         if x=="?" and y=="?": return 1
112         if x=="?": y = i.norm(y); x = 1 if y < .5 else 0
113         elif y=="?": x = i.norm(x); y = 1 if x < .5 else 0
114         else : x,y = i.norm(x), i.norm(y)
115         return abs(x - y)
116
117     def norm(i,x): return 0 if (i.hi-i.lo) < 1E-9 else (x - i.lo) / (i.hi-i.lo)
118
119     def spans(i,j, bins, out):
120         lo = min(i.lo, j.lo)
121         hi = max(i.hi, j.hi)
122         gap = (hi-lo) / bins
123         xys = [(x,"this",1) for x in i._all] + [
124             (x,"that",1) for x in j._all]
125         one = Span(i,lo,lo)
126         all = [one]
127         for x,y,n in sorted(xys, key=first):
128             if one.hi - one.lo > gap:
129                 one = Span(i, one.hi,x)
130                 all += [one]
131             one.add(x,y,n)
132         all = merge(all)
133         all[0].lo = -big
134         all[-1].hi = big
135         if len(all) > 1: out += all

```

```

136 #
137 # eeyen
138 #
139 #
140 #
141 class Span(o):
142     def __init__(i,col, lo, hi, ys=None):
143         i.col, i.lo, i.hi, i.ys = col, lo, hi, ys or Sym()
144
145     def add(i, x:float, y:Any, inc=1) -> None:
146         i.lo = min(x, i.lo)
147         i.hi = max(x, i.hi)
148         i.ys.add(y,inc)
149
150     def merge(i, j): # -> Span|None
151         a, b, c = i.ys, j.ys, i.ys.merge(j.ys)
152         if (i.ys.n==0 or j.ys.n==0 or
153             c.div()*0.99 <= (a.n*a.div() + b.n*b.div())/(a.n + b.n)):
154             return Span(i.col, min(i.lo,j.lo),max(i.hi,j.hi), ys=c)
155
156     def selects(i,row:list) -> bool:
157         x = row[i.col.at]; return x=="?" or i.lo<=x and x<i.hi
158
159     def show(i)-> None:
160         txt = i.col.txt
161         if i.lo == i.hi: return f"[txt] == {i.lo}"
162         elif i.lo == -big: return f"[txt] < {i.hi}"
163         elif i.hi == big: return f"[txt] >= {i.lo}"
164         else : return f"[i.lo] <= [txt] < {i.hi}"
165
166     def support(i) -> float:
167         return i.ys.n / i.col.n
168
169     @staticmethod
170     def sort(spans : list) -> list:
171         "Good spans have large support and low diversity."
172         divs, supports = Num(), Num()
173         sn = lambda s: supports.norm(s.support())
174         dn = lambda s: divs.norm(s.ys.div())
175         f = lambda s: ((1 - sn(s))*2 + dn(s)*2)**.5/2**.5
176         for s in spans:
177             divs.add(s.ys.div())
178             supports.add(s.support())
179         return sorted(spans, key=f)
180 #
181 # eols
182 #
183 #
184 #
185 class Cols(o):
186     def __init__(i,names):
187         i.x,i.y = [],[]
188         i.all = [i.head(at,txt) for at,txt in enumerate(names)]
189
190     def add(i,lst):
191         [col.add(x) for col,x in zip(i.all, lst)]
192         return lst
193
194     def head(i,at,txt):
195         if is_skip(txt): return Skip(at=at, txt=txt)
196         now = (Num if is_num(txt) else Sym)(at=at, txt=txt)
197         if is_klass(txt): i.klass=now
198         (i.y if is_goal(txt) else i.x).append(now)
199         return now
200 #
201 #
202 #
203 #
204 #
205
206 class Sample(o):
207     def __init__(i,the,init=None):
208         i.rows, i.cols, i.cache = [], None, []
209         i.the = copy.deepcopy(the)
210         i.lefts, i.rights, i.left, i.right, c = [],[],None,None,0
211         if type(inits) == list: [i.add(x) for x in inits]
212         if type(inits) == str : [i.add(x) for x in rows(inits)]
213
214     def add(i, sample, row):
215         if not i.left: i.left = row
216         elif not i.right: i.right = row
217         if i.right and i.left:
218             d1 = sample.dist(row, i.left)
219             d2 = sample.dist(row, i.right)
220             if d1 < d2:
221                 if d2 > c: i.left, d1,i.c = row,0,d2
222             else:
223                 if d1 > c: i.right,d2,i.c = row,0,d1
224             (left if d1 < d2 else right).append(row)
225
226     def clone(i, sample, inits=[]):
227         now = i.clone(the)
228         now.add(sample, [col.txt for col in i.cols])
229         [now.add(sample, row) for row in inits]
230         return now
231
232     def dist(i,j,k):
233         cols, p = i.cols.x, i.the.p
234         d = sum(col.dist(j[col.at], k[col.at])**p for col in cols)
235         return (d/len(cols)) ** (1/p)
236
237 #-----
238 s=Sample(o(p=2),"data/auto93.csv")

```