```python
1 '''
2 Created on Mar 11, 2012
3
4 @author: Masum
5 '''
6 import sys
7 from HTMLParser import HTMLParser
8 from shared import parse
9
10
11 class BiGramHTMLParser(HTMLParser):
12     text, N, newline ='', 0,{'br':'\n','BR':'\n'}
13
14     def __init__(self, config, ht):
15         HTMLParser.__init__(self)
16         self.config = config
17         self.ht = ht
18
19
20     def handle_data(self, raw):
21         self.text = self.text+ raw+ self.newline.get(self.lasttag,'')
22
23     #format is (token =>value, here val)
24     def feed(self,data, did):
25         #extracting html text
26         try:
27             HTMLParser.feed(self,data)
28         except:
29             sys.exc_clear()
30
31         #tokenizing extracted data
32         for line in self.text.splitlines():
33             tokens= parse(line, self.config)
34             if not tokens: continue
35             for i in xrange(len(tokens)):
36                 if i > 0 :
37                     self.ht[(tokens[i], tokens[i-1])] +=1;
38                     self.N +=1
39
40         #cleaning for next feed
41         self.text  = ''; self.N=0
42         HTMLParser.reset(self)
43
44
```