

shared

```
1 '''
2 Created on Mar 11, 2012
3
4 @author: Masum
5 '''
6
7 import re
8 from collections import defaultdict
9
10 #===== global regex =====
11 is_word = re.compile("[a-zA-Z_]+$")
12 delim = re.compile("\t|\.|-|:|@|\\|/|,|;|\"|!|\\*|+")
13
14 #=====path =====
15 dst_dir = ''
16 src_dir = ''
17
18 #===== global hash variables =====
19 doc_id= defaultdict(int)
20 term_count= defaultdict(int)
21
22 def parse(line, config):
23     line =line.strip()
24     if not line: return []
25
26     tokens= []
27     for word in delim.split(line):
28         word = word.strip("~`$%^&()[|<>+_-/_").lower()
29         if word in config['str_stop_list']: continue
30         if config['min_token_len']<= len(word) <= config['max_token_len'] and
is_word.match(word):
31             tokens.append(word)
32     return tokens
```