

(<http://xiaorui.cc/2015/06/13/python%E4%BD%BF%E7%94%A8esmr%E4%BB%A3%E6%9B%BFahocorasick%E5%AE%9E%E7%8E%B0ac%t>

为什么会用AC自动机？如果你想知道一篇文章有没有你要过滤的敏感词，怎么办？不可能用正则一个个的匹配吧？敏感词超过300个之后，用Trie来构建模式树(字典树)的速度优势相当的明显...

特别说下，trie图也是一种DFA，可以由trie树为基础构造出来，对于插入的每个模式串，其插入过程中使用的最后一个节点都作为DFA的一个终止节点。

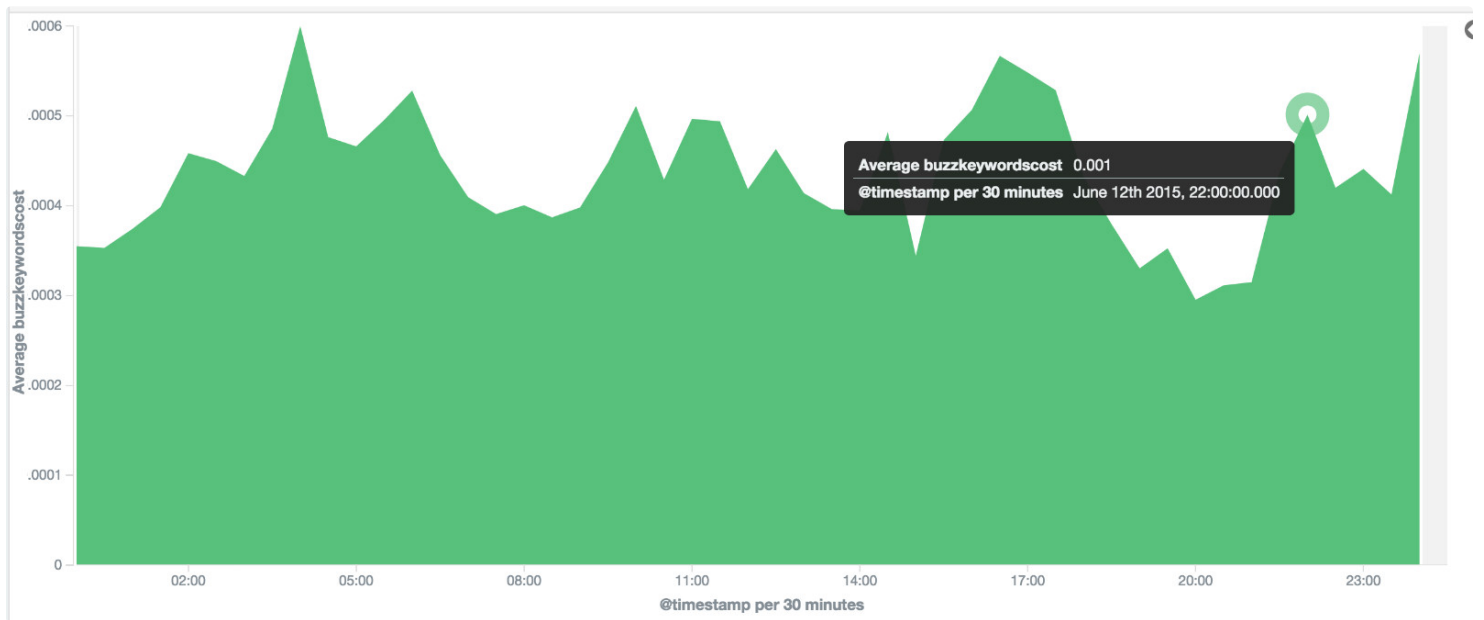
ps: AC自动机是Trie的一种实现，也就是说AC自动机是构造Trie图的DFA的一种方法。还有别的构造DFA的方法...

安装简单，直接pip install esmre

```
In [1]: import esm
>>> index.query('this here in history')
[[{'(1, 4)', 'his'}, {'(5, 7)', 'he'}, {'(13, 16)', 'his'}]]
>>> index.query('Those are his sheep!')
[[{'(8, 13)', 'his'}, {'(14, 17)', 'she'}, {'(15, 17)', 'he'}]]
>>>
In [2]: index = esm.Index()
>>>
In [3]: index.enter("宝马")
You can see more usage examples in the tests.
In [4]: index.enter("马")
<ipython.GeneratorTester>
In [5]: index.enter("奔驰")
In [6]: index.enter("保时捷")
In [7]: index.fix()

In [8]: index.query("哎呀，今天在楼下看到了宝马，我老家倒是有养马的，以前的邻居有个奔驰，不对是保时捷，大爷的，都是马")
Out[8]:
[[{'(33, 39)', '\xe5\xae\x9d\xe9\xa9\xac'),
({'(36, 39)', '\xe9\xa9\xac'),
({'(63, 66)', '\xe9\xa9\xac'),
({'(93, 99)', '\xe5\xa5\x94\xe9\xa9\xb0'),
({'(111, 120)', '\xe4\xbf\x9d\xe6\x97\xb6\xe6\x8d\xb7'),
({'(141, 144)', '\xe9\xa9\xac'})]]
```

测了几天，性能还是可以的，500KB的文章，6000多个关键词，消耗的时间在0.002左右，相比ahocorasick一点都不差的。观察了下，esmr是发现没有发现内存异常泄露等问题。



```
Python
1
2 [2015-06-12 23:34:01,043] INFO extractor "Get keywords takes 0.0003 seconds"
3 [2015-06-12 23:34:01,069] INFO extractor "Get keywords takes 0.0002 seconds"
4 [2015-06-12 23:34:01,178] INFO extractor "Get keywords takes 0.0002 seconds"
5 [2015-06-12 23:34:02,372] INFO extractor "Get keywords takes 0.0002 seconds"
6 [2015-06-12 23:34:02,386] INFO extractor "Get keywords takes 0.0012 seconds"
7 [2015-06-12 23:34:02,631] INFO extractor "Get keywords takes 0.0002 seconds"
8 [2015-06-12 23:34:03,656] INFO extractor "Get keywords takes 0.0021 seconds"
9 [2015-06-12 23:34:03,744] INFO extractor "Get keywords takes 0.0001 seconds"
10 [2015-06-12 23:34:03,785] INFO extractor "Get keywords takes 0.0001 seconds"
11 [2015-06-12 23:34:03,910] INFO extractor "Get keywords takes 0.0002 seconds"
12 [2015-06-12 23:34:04,031] INFO extractor "Get keywords takes 0.0002 seconds"
13 [2015-06-12 23:34:05,004] INFO extractor "Get keywords takes 0.0035 seconds"
14 [2015-06-12 23:34:05,579] INFO extractor "Get keywords takes 0.0055 seconds"
15 [2015-06-12 23:34:05,602] INFO extractor "Get keywords takes 0.0005 seconds"
16 [2015-06-12 23:34:05,662] INFO extractor "Get keywords takes 0.0010 seconds"
17 [2015-06-12 23:34:06,125] INFO extractor "Get keywords takes 0.0002 seconds"
18 [2015-06-12 23:34:06,299] INFO extractor "Get keywords takes 0.0002 seconds"
19 [2015-06-12 23:34:06,404] INFO extractor "Get keywords takes 0.0003 seconds"
20 [2015-06-12 23:34:07,396] INFO extractor "Get keywords takes 0.0002 seconds"
21 [2015-06-12 23:34:07,595] INFO extractor "Get keywords takes 0.0004 seconds"
22 [2015-06-12 23:34:08,725] INFO extractor "Get keywords takes 0.0015 seconds"
23 [2015-06-12 23:34:09,504] INFO extractor "Get keywords takes 0.0004 seconds"
24 [2015-06-12 23:34:09,515] INFO extractor "Get keywords takes 0.0005 seconds"
25 [2015-06-12 23:34:10,650] INFO extractor "Get keywords takes 0.0002 seconds"
26 [2015-06-12 23:34:11,206] INFO extractor "Get keywords takes 0.0003 seconds"
27 [2015-06-12 23:34:12,298] INFO extractor "Get keywords takes 0.0002 seconds"
28 [2015-06-12 23:34:12,319] INFO extractor "Get keywords takes 0.0001 seconds"
29 [2015-06-12 23:34:13,547] INFO extractor "Get keywords takes 0.0006 seconds"
30 [2015-06-12 23:34:13,853] INFO extractor "Get keywords takes 0.0005 seconds"
31
```



python开发及运维研发
多步到 <http://xiaorui.cc>