

D06 使用結巴進行中文斷詞



- >
- 重要知識點 >
- 中文斷詞利器: 結巴(Jieba) >
- 系統環境安裝 >
- 以jieba進行斷詞 >



重要知識點



- 如何使用結巴來進行實際的斷詞操作。



補充資料：簡單易用的中英文斷詞和詞性標註

中文斷詞利器: 結巴(Jieba)

- 要對非結構化的文字資料進行分析，第一件事情是對文字資料抽取結構化的量化數值特徵，除了用「文字探勘分析器」簡單分析字數、句數之外，最常見的分析方式就是斷詞和詞性分析。
- 作為最廣為人使用的斷詞器，結巴(jieba)不只可搭配「非結構化資料分析：文本分類、等機器學習來使用，更可以用在質性研究的內容分析、文本分析或敘說分析上。先用jieba找出特定詞性的文本內容來分析。

系統環境安裝

- 安裝jieba, 只需在終端機上使用

```
pip install jieba
```

- 在.py的腳本中, 引入jieba模組

```
import jieba
```

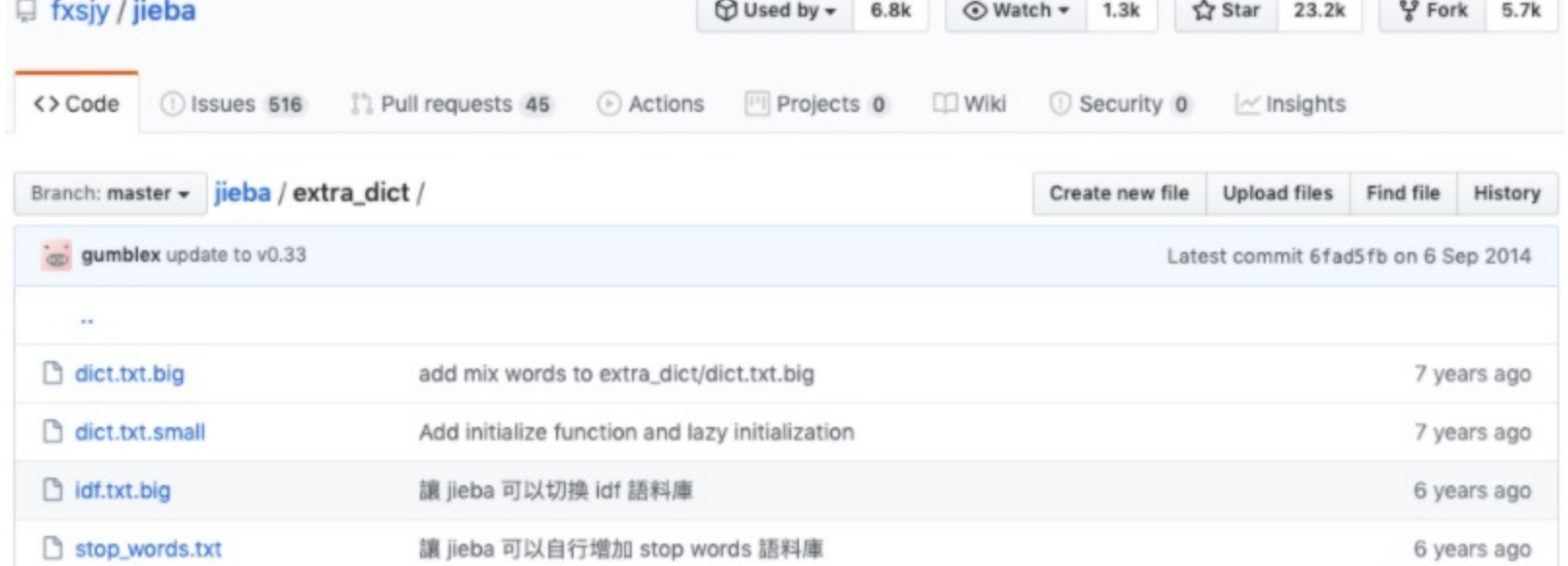
以jieba進行斷詞

- 以結巴進行斷詞
- 結巴斷詞的輸出為generator, 需要使用遊標將值取出



更改使用字典

- #將結巴使用的字典更改為對繁體中文表現較好的字典
- 此字典可從結巴專案github下載



- 設定使用字典



辨識新字詞

- 啟用HMM已辨識新字詞
- 預設HMM功能即為啟用, 可以不用特地設為True



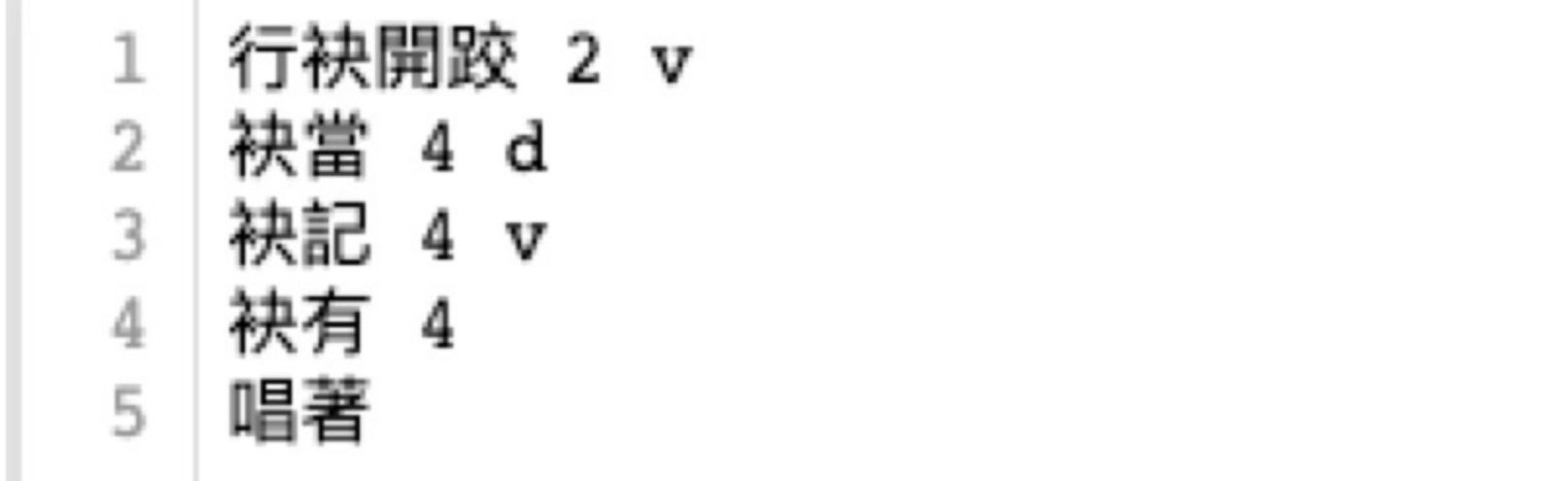
補充資料：github

新增自定義詞庫

在既有使用的字典下新增自定義字詞



userdict.txt 格式如下



動態調整字典

動態加入字典



動態調整詞頻



詞性標注

進行詞性標注(PoS Tagging)



n	名詞	名詞	取英語名詞noun的第1個字樣。	人, 姓, 國家
ng	名詞	名詞數	名詞性語素 - 名詞代綴如n, 建築代綴g前置量詞n。	子, 身, 姓, 堂
nr	名詞	人名	名詞代綴nr(A/en)的釋母放在一起。	羅, 王, 楊

學員可以參考這個網站查看完整對照表，也可以查看結巴官方網站

取出斷詞位置



知識點回顧

這章節我們學習到了

- 使用結巴來進行各項任務(如斷詞)

詳細使用操作

請參閱使用結巴進行斷詞.ipynb檔進行更詳細的使用操作