

[Code-Korean] Tokenize Sentences

Code Snippet

```
# package import
import pandas as pd
import re
from konlpy.tag import Mecab # Korean Tokenizer, Do not import at other languages

# Basic Cleaning Text Function
def CleanText(readData):

    # Remove Retweets
    text = re.sub('RT @[\\w_]+: ', '', readData)

    # Remove Mentions
    text = re.sub('@[\\w_]+', '', text)

    # Remove or Replace URL
    # text = url_re.sub('URL', text)
    text = re.sub(r"http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&]|![!'\(\)\,\;]|(?:%[0-9a-fA-F][0-9a-fA-F]))+", ' ', text) # start with http
    text = re.sub(r"[-a-zA-Z0-9@:%_\+~#={1,256}\. [a-zA-Z0-9()]{2,6}\b([-a-zA-Z0-9()@:%_\+~#?&//=]*)", ' ', text) # Don't start wi

    # Remove Hashtag
    text = re.sub('[#]+[0-9a-zA-Z_]+', ' ', text)

    # Remove Garbage Words (ex. < >, etc)
    text = re.sub('&[a-z]+', ' ', text)

    # Remove Special Characters
    text = re.sub('[^0-9a-zA-Zㄱ-힣]', ' ', text)

    # Remove Numbers (If you want, activate the code)
    # text = re.sub(r'\d+', ' ', text)

    # Remove English (If you want, activate the code)
    # text = re.sub('[a-zA-Z]', ' ', text)

    # Remove newline
    text = text.replace('\n', ' ')

    # Remove multi spacing & Reform sentence
    text = ' '.join(text.split())

    # If you want to normalize Korean text, activate code below:
    # from konlpy.tag import Okt # Must use Konlpy ver 0.5.2 above
    # okt = Okt()
```

