

# 01 tokenization

Kazuki Yabe

September 2019

## 1 Segmenter

### 1.1 Segmenters

#### 1.1.1 Segmenter from the slide

This segmenter simply add new line after "?", "!", or ".".

#### 1.1.2 NLTK segmenter

This segmenter uses `nltk.tokenize.sent_tokenizer` for sentence segmentation.

### 1.2 Quantitative

I manually looked at the first 200 examples, and calculated the accuracy for both segmenters. NLTK segmenter performed slightly better than the segmenter from the slide by 1%, where NLTK got 100% accuracy and The slide's segmenter got 99.5% accuracy. The only difference was that the slide's segmenter wrongly segmented a quote with a period inside quotation marks.

### 1.3 Qualitative

NLTK segmenter performs better in segmenting URLs, quotations, abbreviations, and decimals.

Types	NLTK	Slide
URL	<code>http://c2.com/cgi/wiki?WikiHistory</code>	<code>http://c2.</code>
Quatations	<code>"mangingibig ng karunungan."</code>	<code>"mangingibig ng karunungan.</code>
Abbreviations	<code>Osmeña Jr."</code> at lolo nina ...	<code>Osmeña Jr.</code>
Decimal Numbers	<code>hanggang 26.5%</code> ...	<code>hanggang 26.</code>

URL, quotation with period eg "hello." Jr. 26.5 (decimal number)

## 2 Tokenizer

### 2.1 Implementation

/01.Tokenisation/maxmatch.py

### 2.2 Instruction

```
python maxmatch.py
```

The above command run maxmatch on test data which is extracted from test.conllu. Use ">file\_name" to save the result to a desired file.

```
python maxmatch_evaluation.py
```

The above command run WER on results of maxmatch.py. Use ">file\_name" to save the result to a desired file.

### 2.3 Evaluation

WER score is calculated for evaluation of the maxmatch algorithm. The score is 8.4%. In general, if a word is not in dictionary, then the tokenizer separates that word into some miscellaneous words that the tokenizer recognizes.

For example, the algorithm had a problem in identifying names of organizations, like "幸福の科学". It parsed it as "幸福" "の" "科学," by treating the name as noun-particle-noun.

Since there are many homograph