

Tokenization

```
In [1]: 1 import nltk
        2 nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\Xiaowei\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping corpora\wordnet.zip.
```

Out[1]: True

```
In [2]: 1 text = "This is Andrew's text, isn't it?"
```

```
In [3]: 1 tokenizer = nltk.tokenize.WhitespaceTokenizer()
        2 tokenizer.tokenize(text)
```

Out[3]: ['This', 'is', "Andrew's", 'text,', "isn't", 'it?']

```
In [4]: 1 tokenizer = nltk.tokenize.TreebankWordTokenizer()
        2 tokenizer.tokenize(text)
```

Out[4]: ['This', 'is', 'Andrew', "'s", 'text', ',', 'is', "n't", 'it', '?']

```
In [5]: 1 tokenizer = nltk.tokenize.WordPunctTokenizer()
        2 tokenizer.tokenize(text)
```

Out[5]: ['This', 'is', 'Andrew', "'", 's', 'text', ',', 'isn', '"', 't', 'it', '?']

Stemming (further in the video)

```
In [8]: 1 text = "feet wolves cats talked"
        2 tokenizer = nltk.tokenize.TreebankWordTokenizer()
        3 tokens = tokenizer.tokenize(text)
```

```
In [9]: 1 stemmer = nltk.stem.PorterStemmer()
        2 " ".join(stemmer.stem(token) for token in tokens)
```

Out[9]: 'feet wolv cat talk'

```
In [10]: 1 stemmer = nltk.stem.WordNetLemmatizer()
         2 " ".join(stemmer.lemmatize(token) for token in tokens)
```

Out[10]: 'foot wolf cat talked'

```
In [ ]: 1
```