

Tokenization in NLTK

Bird, Loper and Klein (2009), *Natural Language Processing with Python*. O'Reilly

```
>>> text = 'That U.S.A. poster-print costs $12.40...'
>>> pattern = r'''(?x)          # set flag to allow verbose regexps
...     ([A-Z]\.)+              # abbreviations, e.g. U.S.A.
...     | \w+(-\w+)*            # words with optional internal hyphens
...     | \$?\d+(\.\d+)?%?      # currency and percentages, e.g. $12.40, 82%
...     | \.\.\.               # ellipsis
...     | [][.,;"'()?:_-']     # these are separate tokens; includes ], [
...     '''
>>> nltk.regexp_tokenize(text, pattern)
['That', 'U.S.A.', 'poster-print', 'costs', '$12.40', '...']
```