

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization
Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

Natural Language ToolKit

November 27, 2017

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text preprocessing

Text normalization

Segmentation

Tokenization

Stemming

Lemmatization

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

Overview

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization
Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

Natural Language ToolKit (NLTK)

- ▶ Python-based
- ▶ Comprehensive
- ▶ Useful for research and prototyping

NLTK Modules

Language processing task	NLTK modules	Functionality
Accessing corpora	corpus	standardized interfaces to corpora and lexicons
String processing	tokenize, stem	tokenizers, sentence tokenizers, stemmers
Collocation discovery	collocations	t-test, chi-squared, PMI
Part-of-speech tagging	tag	n-gram, backoff, Brill, HMM, TnT
Machine learning	classify, cluster, tbl	decision tree, max entropy, naive Bayes, EM, k-means
Chunking	chunk	regex, n-gram, NER
Parsing	parse, ccg	chart, feature-based, unification, probabilistic, dependency
Semantic interpretation	sem, inference	lambda calculus, first-order logic, model checking
Evaluation metrics	metrics	precision, recall, agreement coefficients
Probability and estimation	probability	frequency and smoothed probability distributions
Applications	app, chat	concordancer, parsers, WordNet browser, chatbots
Linguistic fieldwork	toolbox	manipulate data in SIL Toolbox format

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text normalization

Segmentation

Tokenization

Stemming

Lemmaatization

Stop word removal

Rare word removal

References

References

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

- ▶ Simplicity
- ▶ Consistency
- ▶ Extensibility
- ▶ Modularity

- ▶ Installation
- ▶ Theoretical background
- ▶ NLU
- ▶ Accessing corpora
- ▶ FreqDist, collocations, n-grams
- ▶ Stemming
- ▶ Lemmatization
- ▶ Segmentation, Tokenization
- ▶ POS-Taggers
- ▶ Chunking
- ▶ Information extraction
- ▶ Classifiers
- ▶ NER
- ▶ Parsers

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Installation

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

- ▶ Python
<https://wiki.python.org/moin/BeginnersGuide/Download>
- ▶ NLTK
<https://pypi.python.org/pypi/nltk>
- ▶ PyCharm
<https://www.jetbrains.com/pycharm/download/>

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text
preprocessing

Text
normalization
Segmentation
Tokenization
Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

Text preprocessing

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

**Text
normalization**

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

- ▶ tokenization
- ▶ stemming
- ▶ lemmatization
- ▶ stop word removal
- ▶ rare word removal

Tokenization

- ▶ Segmentation =
- ▶ Sentence tokenization =
- ▶ Sentence splitting =
- ▶ Sentence boundary detection, etc.

From splitting the string on (.) to a predictive classifier

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text
preprocessing

Text
normalization
Segmentation
Tokenization
Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

Sample rules

1. Period?
2. The preceding token is in the list of abbreviations?
3. The next token is capitalized?
4. etc.

Better: maximum entropy models, NNs

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

On-demand segmentation

```
1 from nltk.tokenize import sent_tokenize
2 inputstring = 'This is an example sent.
   Will the sentence tokenizer split on
   sent markers?!'
3 all_sent = sent_tokenize(inputstring)
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text
preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

On-demand segmentation

```
1 from nltk.tokenize import sent_tokenize
2 inputstring = 'This is an example sent.
   Will the sentence tokenizer split on
   sent markers?!'
3 all_sent = sent_tokenize(inputstring)
```

*[' This is an example sent', 'Will the sentence tokenizer split
on sent markers?!']*

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text
preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

PunktSentenceTokenizer instantiation

```
1 import nltk.tokenize.punkt
2 tokenizer = nltk.tokenize.punkt.  
    PunktSentenceTokenizer()
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

PunktSentenceTokenizer instantiation

```
1 import nltk.tokenize.punkt
2 tokenizer = nltk.tokenize.punkt.
    PunktSentenceTokenizer()
```

*[' This is an example sent', 'Will the sentence tokenizer split
on sent markers?!']*

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Tokenizing sentences in other languages

```
1 import nltk.data
2 tokenizer = nltk.data.load('tokenizers/
    punkt/PY3/english.pickle')
3 tokenizer.tokenize(inputstring)
4 spanish_tokenizer = nltk.data.load('
    tokenizers/punkt/PY3/spanish.pickle')
5 spanish_tokenizer.tokenize('Hola amigo.
    Estoy bien.')
```

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization

Segmentation

- Tokenization

- Stemming

- Lemmatization

- Stop word removal

- Rare word removal

- References

References

Tokenizing sentences in other languages

```
1 import nltk.data
2 tokenizer = nltk.data.load('tokenizers/
    punkt/PY3/english.pickle')
3 tokenizer.tokenize(inputstring)
4 spanish_tokenizer = nltk.data.load('
    tokenizers/punkt/PY3/spanish.pickle')
5 spanish_tokenizer.tokenize('Hola amigo.
    Estoy bien.')
```

['Hola amigo.', 'Estoy bien.']

You can see a list of all the available language tokenizers in
/usr/share/nltk_data/ tokenizers/punkt/PY3 (or
C:\nltk_data\tokenizers\punkt\PY3).

[Overview](#)[NLTK Modules](#)
[NLTK Principles](#)
[NLTK Features](#)[Installation](#)[Installation](#)[Text
preprocessing](#)[Text
normalization](#)**[Segmentation](#)**[Tokenization](#)[Stemming](#)[Lemmatization](#)[Stop word removal](#)[Rare word removal](#)[References](#)[References](#)

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization

Segmentation

- Tokenization

- Stemming

- Lemmatization

- Stop word removal

- Rare word removal

- References

References

Training a sentence tokenizer

White guy: So, do you have any plans for this evening?

Asian girl: Yeah, being angry!

White guy: Oh, that sounds good.

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization

Segmentation

- Tokenization

- Stemming

- Lemmatization

- Stop word removal

- Rare word removal

- References

References

```
1 from nltk.tokenize import
   PunktSentenceTokenizer
2 from nltk.corpus import webtext
3 text = webtext.raw('overheard.txt')
4 sent_tokenizer = PunktSentenceTokenizer(
   text)
5 sents1 = sent_tokenizer.tokenize(text)
6 sents1[0]
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

```
1 from nltk.tokenize import
   PunktSentenceTokenizer
2 from nltk.corpus import webtext
3 text = webtext.raw('overheard.txt')
4 sent_tokenizer = PunktSentenceTokenizer(
   text)
5 sents1 = sent_tokenizer.tokenize(text)
6 sents1[0]
```

'White guy: So, do you have any plans for this evening?'

Compare:

```
1 from nltk.tokenize import sent_tokenize
2 sents2 = sent_tokenize(text)
3 sents1[678]
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Compare:

```
1 from nltk.tokenize import sent_tokenize
2 sents2 = sent_tokenize(text)
3 sents1[678]
```

'Girl: But you already have a Big Mac...'

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Compare:

```
1 from nltk.tokenize import sent_tokenize
2 sents2 = sent_tokenize(text)
3 sents1[678]
```

'Girl: But you already have a Big Mac...'

```
1 sents2[678]
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Compare:

```
1 from nltk.tokenize import sent_tokenize
2 sents2 = sent_tokenize(text)
3 sents1[678]
```

'Girl: But you already have a Big Mac...'

```
1 sents2[678]
```

'Girl: But you already have a Big Mac... \nHobo: Oh, this is all theatrical.'

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

PunktSentenceTokenizer

Heuristics:

- ▶ Orthography
- ▶ Frequent sentence starters
- ▶ Collocations
- ▶ Initials

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation

Tokenization

- Stemming
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

The simplest tokenizer: the `split()` method of Python strings.

```
1 s = "Hi Everyone ! hola gr8"  
2 s.split()
```

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation

Tokenization

- Stemming
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

The simplest tokenizer: the `split()` method of Python strings.

```
1 s = "Hi Everyone ! hola gr8"  
2 s.split()
```

['Hi', 'Everyone', '!', 'hola', 'gr8']

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization**
- Stemming
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

The word_tokenize method

```
1 from nltk.tokenize import word_tokenize
2 word_tokenize(s)
```

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation

Tokenization

- Stemming
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

The word_tokenize method

```
1 from nltk.tokenize import word_tokenize
2 word_tokenize(s)
```

['Hi', 'Everyone', '!', 'hola', 'gr8']

The `regex_tokenize` method

```
1 from nltk.tokenize import regexp_tokenize
2 regexp_tokenize(s, pattern='\\w+')
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

The `regex_tokenize` method

```
1 from nltk.tokenize import regexp_tokenize
2 regexp_tokenize(s, pattern='\\w+')

```

['Hi', 'Everyone', 'hola', 'gr8']

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

The `regex_tokenize` method

```
1 from nltk.tokenize import regexp_tokenize
2 regexp_tokenize(s, pattern='\w+')

```

['Hi', 'Everyone', 'hola', 'gr8']

```
1 regexp_tokenize(s, pattern='\d+')

```

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

The `regex_tokenize` method

```
1 from nltk.tokenize import regexp_tokenize
2 regexp_tokenize(s, pattern='\w+')

```

['Hi', 'Everyone', 'hola', 'gr8']

```
1 regexp_tokenize(s, pattern='\d+')

```

['8']

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

The wordpunct_tokenize and blankline_tokenize methods

```
1 from nltk.tokenize import  
   wordpunct_tokenize, blankline_tokenize  
2 wordpunct_tokenize(s)
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text
preprocessing

Text
normalization
Segmentation

Tokenization

Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

The wordpunct_tokenize and blankline_tokenize methods

```
1 from nltk.tokenize import  
   wordpunct_tokenize, blankline_tokenize  
2 wordpunct_tokenize(s)
```

```
['Hi', ',', 'Everyone', '!!', 'hola', 'gr8']
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text
preprocessing

Text
normalization
Segmentation

Tokenization

Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

The wordpunct_tokenize and blankline_tokenize methods

```
1 from nltk.tokenize import  
   wordpunct_tokenize, blankline_tokenize  
2 wordpunct_tokenize(s)
```

['Hi', ',', 'Everyone', '!!', 'hola', 'gr8']

```
1 blankline_tokenize(s)
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text
preprocessing

Text
normalization
Segmentation

Tokenization

Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

The wordpunct_tokenize and blankline_tokenize methods

```
1 from nltk.tokenize import  
   wordpunct_tokenize, blankline_tokenize  
2 wordpunct_tokenize(s)
```

['Hi', ',', 'Everyone', '!!', 'hola', 'gr8']

```
1 blankline_tokenize(s)
```

['Hi, Everyone !! hola gr8']

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation

Tokenization

- Stemming
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

Experiment

```
1 from nltk.tokenize import word_tokenize,  
   regexp_tokenize, wordpunct_tokenize,  
   blankline_tokenize
```


Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation

Tokenization

- Stemming
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

Experiment

```
1 from nltk.tokenize import word_tokenize,  
   regexp_tokenize, wordpunct_tokenize,  
   blankline_tokenize
```

<http://text-processing.com/demo>

More on the word_tokenize method

```
1 from nltk.tokenize import  
   TreebankWordTokenizer  
2 tokenizer = TreebankWordTokenizer()  
3 tokenizer.tokenize('Hello World.')
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

More on the word_tokenize method

```
1 from nltk.tokenize import  
   TreebankWordTokenizer  
2 tokenizer = TreebankWordTokenizer()  
3 tokenizer.tokenize('Hello World.')
```

```
['Hello', 'World', '.']
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Other subclasses of Tokenizerl:

- ▶ TreebankWordTokenizer
- ▶ WhitespaceTokenizer
- ▶ SpaceTokenizer
- ▶ PunktWordTokenizer
- ▶ WordPunctTokenizer
- ▶ BlanklineTokenizer
- ▶ etc.

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text
preprocessing

Text
normalization
Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Separating contractions

```
1 from nltk.tokenize import  
   TreebankWordTokenizer  
2 tokenizer = TreebankWordTokenizer()  
3 tokenizer.tokenize("can't")
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Separating contractions

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

```
1 from nltk.tokenize import  
   TreebankWordTokenizer  
2 tokenizer = TreebankWordTokenizer()  
3 tokenizer.tokenize("can't")
```

['ca', 'n't']

Separating contractions

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation

Tokenization

- Stemming
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.tokenize import  
   PunktWordTokenizer  
2 tokenizer = PunktWordTokenizer()  
3 tokenizer.tokenize("Can't is a  
   contraction.")
```

Separating contractions

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text

- normalization

- Segmentation

- Tokenization**

- Stemming

- Lemmatization

- Stop word removal

- Rare word removal

- References

References

```
1 from nltk.tokenize import  
   PunktWordTokenizer  
2 tokenizer = PunktWordTokenizer()  
3 tokenizer.tokenize("Can't is a  
   contraction.")
```

['Can', "'t", 'is', 'a', 'contraction.']

Separating contractions

```
1 from nltk.tokenize import  
   WordPunctTokenizer  
2 tokenizer = WordPunctTokenizer()  
3 tokenizer.tokenize("Can't is a  
   contraction.")
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

Separating contractions

```
1 from nltk.tokenize import  
   WordPunctTokenizer  
2 tokenizer = WordPunctTokenizer()  
3 tokenizer.tokenize("Can't is a  
   contraction.")
```

```
['Can', "'", 't', 'is', 'a', 'contraction', '.']
```

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

More on tokenization with regular expressions

Match on:

- ▶ tokens
- ▶ separators or gaps

Regex tokenization

```
1 from nltk.tokenize import RegexpTokenizer
2 tokenizer = RegexpTokenizer("[\w']+")
3 tokenizer.tokenize("Can't is a
   contraction.")
```

Equal to:

```
1 from nltk.tokenize import regexp_tokenize
2 regexp_tokenize("Can't is a contraction.",
   , "[\w']+")
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Regex tokenization

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

```
1 from nltk.tokenize import RegexpTokenizer
2 tokenizer = RegexpTokenizer("[\w']+")
3 tokenizer.tokenize("Can't is a
   contraction.")
```

Equal to:

```
1 from nltk.tokenize import regexp_tokenize
2 regexp_tokenize("Can't is a contraction.",
   , "[\w']+")
```

["Can't", 'is', 'a', 'contraction']

Simple whitespace tokenizer

```
1 tokenizer = RegexpTokenizer('\s+', gaps=  
    True)  
2 tokenizer.tokenize("Can't is a  
    contraction.")
```

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization**
- Stemming
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

Simple whitespace tokenizer

```
1 tokenizer = RegexpTokenizer('\s+', gaps=True)
2 tokenizer.tokenize("Can't is a contraction.")
```

["Can't", 'is', 'a', 'contraction.']

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming**
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

Example: *eat*

eating, eaten, eats

Stemming strategies

- ▶ Lookup table
- ▶ Rule-based affix stripping
- ▶ Hybrid approaches
- ▶ Suffix substitution
- ▶ Recursive rules

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Stemming

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization

Stemming

Lemmatization
Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import PorterStemmer
2 from nltk.stem.lancaster import
   LancasterStemmer
3 pst = PorterStemmer()
4 lst = LancasterStemmer()
5 lst.stem("eating")
```

Stemming

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization

Stemming

Lemmatization
Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import PorterStemmer
2 from nltk.stem.lancaster import
   LancasterStemmer
3 pst = PorterStemmer()
4 lst = LancasterStemmer()
5 lst.stem("eating")
```

'eat'

Stemming

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization

Stemming

Lemmatization
Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import PorterStemmer
2 from nltk.stem.lancaster import
    LancasterStemmer
3 pst = PorterStemmer()
4 lst = LancasterStemmer()
5 lst.stem("eating")

    'eat'

1 pst.stem("hopping")
```

Stemming

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization

Stemming

Lemmatization
Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import PorterStemmer
2 from nltk.stem.lancaster import
    LancasterStemmer
3 pst = PorterStemmer()
4 lst = LancasterStemmer()
5 lst.stem("eating")

'eat'

1 pst.stem("hopping")
```

'hop'

Subclasses of Stemmerl:

- ▶ PorterStemmer
- ▶ LancasterStemmer
- ▶ RegexStemmer
- ▶ SnowballStemmer

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming**
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import RegexpStemmer
2 stemmer = RegexpStemmer('ing')
3 stemmer.stem('cooking')
```

Stemming

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming**
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import RegexpStemmer
2 stemmer = RegexpStemmer('ing')
3 stemmer.stem('cooking')
```

'cook'

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization

Stemming

- Lemmatization
- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import RegexpStemmer
2 stemmer = RegexpStemmer('ing')
3 stemmer.stem('cooking')
```

'cook'

```
1 stemmer.stem('ingenious')
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization

Stemming

Lemmatization
Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import RegexpStemmer
2 stemmer = RegexpStemmer('ing')
3 stemmer.stem('cooking')
```

'cook'

```
1 stemmer.stem('ingenious')
```

'enious'

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization

Stemming

Lemmatization
Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import SnowballStemmer
2 SnowballStemmer.languages('danish', '
    dutch', 'english', 'finnish', 'french'
    , 'german', 'hungarian', 'italian', '
    norwegian', 'porter', 'portuguese', '
    romanian', 'russian', 'spanish', '
    swedish')
3 spanish_stemmer = SnowballStemmer('
    spanish')
4 spanish_stemmer.stem('hola')
```

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization

Stemming

Lemmatization
Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import SnowballStemmer
2 SnowballStemmer.languages('danish', '
    dutch', 'english', 'finnish', 'french'
    , 'german', 'hungarian', 'italian', '
    norwegian', 'porter', 'portuguese', '
    romanian', 'russian', 'spanish', '
    swedish')
3 spanish_stemmer = SnowballStemmer('
    spanish')
4 spanish_stemmer.stem('hola')
```

'hol'

Experiment

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

```
1 from nltk.stem import PorterStemmer,
   SnowballStemmer, RegexpStemmer
2 from nltk.stem.lancaster import
   LancasterStemmer
3 g_snst = SnowballStemmer('german')
4 pst = PorterStemmer()
5 lst = LancasterStemmer()
6 rst = RegexpStemmer()
7
8 stemer.stem()
```

Experiment

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming**
- Lemmatization
- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import PorterStemmer,
   SnowballStemmer, RegexpStemmer
2 from nltk.stem.lancaster import
   LancasterStemmer
3 g_snst = SnowballStemmer('german')
4 pst = PorterStemmer()
5 lst = LancasterStemmer()
6 rst = RegexpStemmer()
7
8 stemer.stem()
```

<http://text-processing.com/demo>

Lemmatisation algorithms

- ▶ Stochastic algorithms
- ▶ n-gram analysis
- ▶ Hybrid approaches

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming

Lemmatization

- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import WordNetLemmatizer
2 wlem = WordNetLemmatizer()
3 wlem.lemmatize("ate")
```


Lemmatization

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization
Stemming

Lemmatization

Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import WordNetLemmatizer
2 wlem = WordNetLemmatizer()
3 wlem.lemmatize("ate")
```

'eat'

Lemmatization

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming

Lemmatization

- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import WordNetLemmatizer
2 wlem = WordNetLemmatizer()
3 wlem.lemmatize('cooking')
```

Lemmatization

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming

Lemmatization

- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import WordNetLemmatizer
2 wlem = WordNetLemmatizer()
3 wlem.lemmatize('cooking')
```

'cooking'

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming

Lemmatization

- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import WordNetLemmatizer
2 wlem = WordNetLemmatizer()
3 wlem.lemmatize('cooking')
```

'cooking'

```
1 wlem.lemmatize('cooking', pos='v')
```

Lemmatization

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming

Lemmatization

- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import WordNetLemmatizer
2 wlem = WordNetLemmatizer()
3 wlem.lemmatize('cooking')
```

'cooking'

```
1 wlem.lemmatize('cooking', pos='v')
```

'cook'

Lemmatization

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization
Stemming

Lemmatization

Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import WordNetLemmatizer
2 wlem = WordNetLemmatizer()
3 wlem.lemmatize('cooking')
```

'cooking'

```
1 wlem.lemmatize('cooking', pos='v')
```

'cook'

```
1 wlem.lemmatize('cookbooks')
```

Lemmatization

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization
Stemming

Lemmatization

Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import WordNetLemmatizer
2 wlem = WordNetLemmatizer()
3 wlem.lemmatize('cooking')
```

'cooking'

```
1 wlem.lemmatize('cooking', pos='v')
```

'cook'

```
1 wlem.lemmatize('cookbooks')
```

'cookbook'

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming

Lemmatization

- Stop word removal
- Rare word removal
- References

References

```
1 from nltk.stem import PorterStemmer
2 from nltk.stem import WordNetLemmatizer
3 stemmer = PorterStemmer()
4 stemmer.stem('believes')
```


Lemmatization

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References

```
1 from nltk.stem import PorterStemmer
2 from nltk.stem import WordNetLemmatizer
3 stemmer = PorterStemmer()
4 stemmer.stem('believes')
```

'believ'

Lemmatization

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization
Stemming

Lemmatization

Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import PorterStemmer
2 from nltk.stem import WordNetLemmatizer
3 stemmer = PorterStemmer()
4 stemmer.stem('believes')
```

'believ'

```
1 lemmatizer = WordNetLemmatizer()
2 lemmatizer.lemmatize('believes')
```

Lemmatization

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization
Stemming

Lemmatization

Stop word removal
Rare word removal
References

References

```
1 from nltk.stem import PorterStemmer
2 from nltk.stem import WordNetLemmatizer
3 stemmer = PorterStemmer()
4 stemmer.stem('believes')
```

'believ'

```
1 lemmatizer = WordNetLemmatizer()
2 lemmatizer.lemmatize('believes')
```

'belief'

Stop word removal

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming
- Lemmatization

Stop word removal

- Rare word removal
- References

References

```
1 from nltk.corpus import stopwords
2 stoplist = stopwords.words('english')
3 text = "This is just a test"
4 cleanwordlist = [word for word in text.
                    split() if word not in stoplist]
```

Stop word removal

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming
- Lemmatization

Stop word removal

- Rare word removal
- References

References

```
1 from nltk.corpus import stopwords
2 stoplist = stopwords.words('english')
3 text = "This is just a test"
4 cleanwordlist = [word for word in text.
                    split() if word not in stoplist]
```

['test']

Stop word removal

Overview

- NLTK Modules
- NLTK Principles
- NLTK Features

Installation

- Installation

Text

preprocessing

- Text normalization
- Segmentation
- Tokenization
- Stemming
- Lemmatization

Stop word removal

- Rare word removal
- References

References

```
1 from nltk.corpus import stopwords
2 stoplist = stopwords.words('english')
3 text = "This is just a test"
4 cleanwordlist = [word for word in text.
                    split() if word not in stoplist]
```

['test']

```
1 stopwords.fileids()
```

Stop word removal

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization
Stemming
Lemmatization

Stop word removal

Rare word removal
References

References

```
1 from nltk.corpus import stopwords
2 stoplist = stopwords.words('english')
3 text = "This is just a test"
4 cleanwordlist = [word for word in text.
                    split() if word not in stoplist]
```

['test']

```
1 stopwords.fileids()
```

*['danish', 'dutch', 'english', 'finnish', 'french', 'german',
'hungarian', 'italian', 'norwegian', 'portuguese', 'russian',
'spanish', 'swedish', 'turkish']*

Rare word removal

Overview

NLTK Modules
NLTK Principles
NLTK Features

Installation

Installation

Text

preprocessing

Text
normalization
Segmentation
Tokenization
Stemming
Lemmatization
Stop word removal
Rare word removal
References

References

```
1 freq_dist = nltk.FreqDist(tokens)
2 rarewords = freq_dist.keys()[-50:]
3 nonrarewords = [ word for word in token
                   not in rarewords]
```


- Bird, Steven, Ewan Klein, and Edward Loper (2009). *Natural language processing with Python: analyzing text with the natural language toolkit*. " O'Reilly Media, Inc."
- Perkins, Jacob (2014). *Python 3 Text Processing with NLTK 3 Cookbook*. Packt Publishing Ltd.

Overview

NLTK Modules

NLTK Principles

NLTK Features

Installation

Installation

Text

preprocessing

Text

normalization

Segmentation

Tokenization

Stemming

Lemmatization

Stop word removal

Rare word removal

References

References