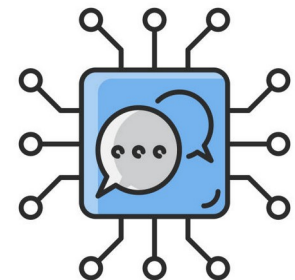


Python for NLP

Tushar B. Kute,
<http://tusharkute.com>



Pre-processing

- To prepare the text data for the model building we perform text preprocessing. It is the very first step of NLP projects. Some of the preprocessing steps are:
 - Removing punctuations like . , ! \$ () * % @
 - Removing URLs
 - Removing Stop words
 - Lower casing
 - Tokenization
 - Stemming
 - Lemmatization

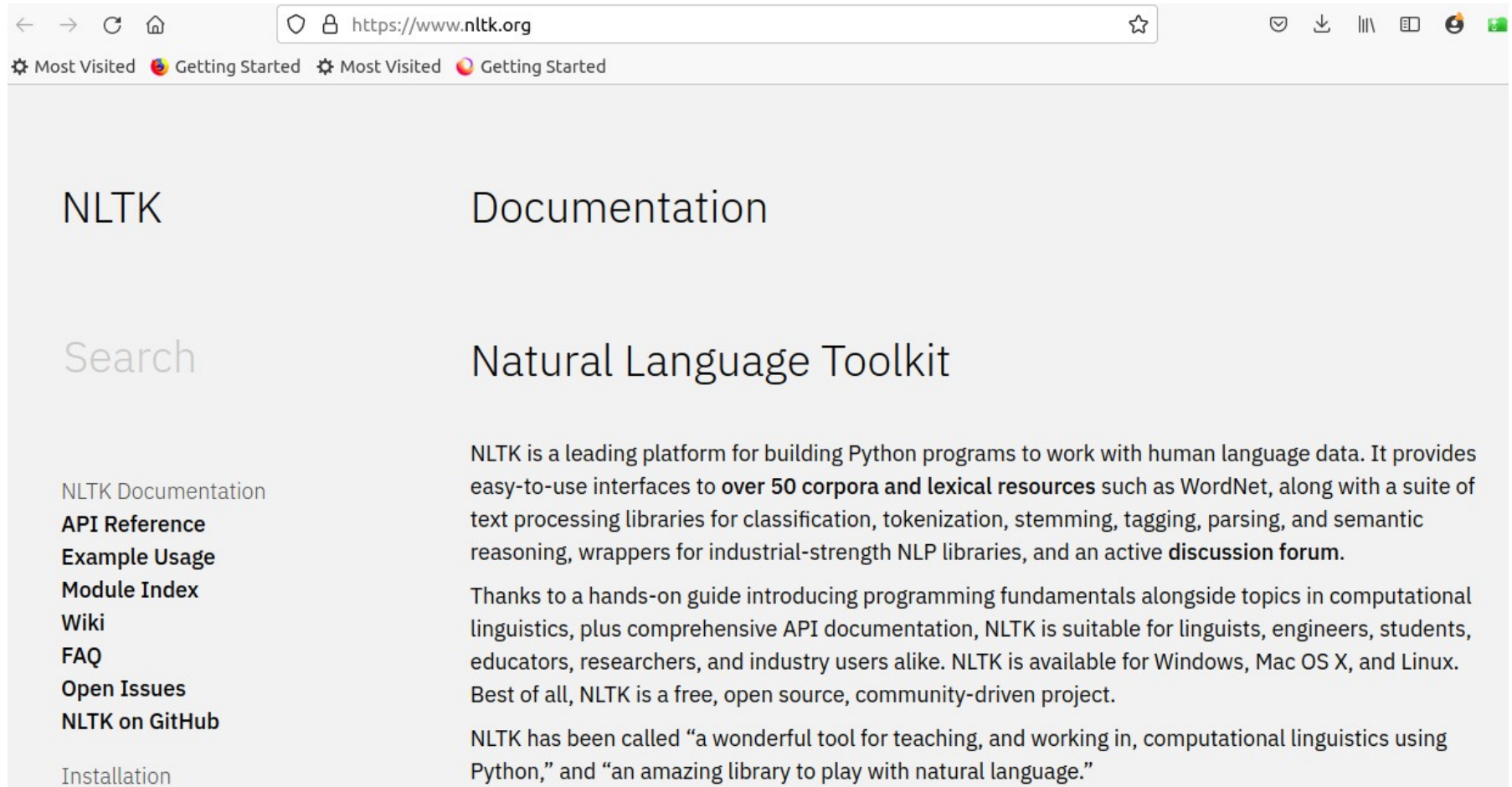
Why Pre-processing?

- Significance of text preprocessing in the performance of models.
- Data preprocessing is an essential step in building a Machine Learning model and depending on how well the data has been preprocessed; the results are seen.
- In NLP, text preprocessing is the first step in the process of building a model.

- The Natural Language Toolkit, or more commonly NLTK, is a suite of libraries and programs for symbolic and statistical natural language processing (NLP) for English written in the Python programming language.
- It was developed by Steven Bird and Edward Loper in the Department of Computer and Information Science at the University of Pennsylvania.
- NLTK includes graphical demonstrations and sample data. It is accompanied by a book that explains the underlying concepts behind the language processing tasks supported by the toolkit, plus a cookbook.

NLTK

- NLTK is intended to support research and teaching in NLP or closely related areas, including empirical linguistics, cognitive science, artificial intelligence, information retrieval, and machine learning.
- NLTK has been used successfully as a teaching tool, as an individual study tool, and as a platform for prototyping and building research systems.
- There are 32 universities in the US and 25 countries using NLTK in their courses.
- NLTK supports classification, tokenization, stemming, tagging, parsing, and semantic reasoning functionalities



The screenshot shows a web browser window with the URL <https://www.nltk.org>. The page layout includes a left sidebar with navigation links and a main content area. The sidebar links are: NLTK Documentation, API Reference, Example Usage, Module Index, Wiki, FAQ, Open Issues, NLTK on GitHub, and Installation. The main content area features the title 'NLTK Documentation' and 'Natural Language Toolkit'. Below the title, there is a paragraph describing NLTK as a leading platform for building Python programs to work with human language data, providing interfaces to over 50 corpora and lexical resources. It also mentions a hands-on guide, comprehensive API documentation, and availability for Windows, Mac OS X, and Linux. A quote at the bottom describes NLTK as a wonderful tool for teaching and working in computational linguistics.

NLTK Documentation

Natural Language Toolkit

NLTK is a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to **over 50 corpora and lexical resources** such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning, wrappers for industrial-strength NLP libraries, and an active **discussion forum**.

Thanks to a hands-on guide introducing programming fundamentals alongside topics in computational linguistics, plus comprehensive API documentation, NLTK is suitable for linguists, engineers, students, educators, researchers, and industry users alike. NLTK is available for Windows, Mac OS X, and Linux. Best of all, NLTK is a free, open source, community-driven project.

NLTK has been called “a wonderful tool for teaching, and working in, computational linguistics using Python,” and “an amazing library to play with natural language.”

- NLTK Documentation
- API Reference
- Example Usage
- Module Index
- Wiki
- FAQ
- Open Issues
- NLTK on GitHub
- Installation

Install nltk

- `!pip install nltk -U`
- Installing nltk packages
 - `import nltk`
 - `nltk.download('package-name')`

Using Python Scripts

```
>>> import nltk
>>> sentence = """At eight o'clock on Thursday morning
... Arthur didn't feel very good."""
>>> tokens = nltk.word_tokenize(sentence)
>>> tokens
['At', 'eight', "o'clock", 'on', 'Thursday', 'morning',
 'Arthur', 'did', "n't", 'feel', 'very', 'good', '.']
>>> tagged = nltk.pos_tag(tokens)
>>> tagged[0:6]
[('At', 'IN'), ('eight', 'CD'), ("o'clock", 'JJ'), ('on', 'IN'),
 ('Thursday', 'NNP'), ('morning', 'NN')]
```


Using Python Scripts

```
>>> entities = nltk.chunk.ne_chunk(tagged)
>>> entities
Tree('S', [(['At', 'IN'], ('eight', 'CD'), ("o'clock", 'JJ'),
            ('on', 'IN'), ('Thursday', 'NNP'), ('morning', 'NN'),
            Tree('PERSON', [(['Arthur', 'NNP'])],
                ('did', 'VBD'), ("n't", 'RB'), ('feel', 'VB'),
                ('very', 'RB'), ('good', 'JJ'), ('.', '.')]])])])
```

Thank you

This presentation is created using LibreOffice Impress 7.4.1.2, can be used freely as per GNU General Public License



@mitu_skillologies



@mITuSkillologies



@mitu_group



@mitu-skillologies



@MITUSkillologies

kaggle

@mituskillologies

Web Resources

<https://mitu.co.in>

<http://tusharkute.com>



@mituskillologies

contact@mitu.co.in
tushar@tusharkute.com