

[Download This Cheat Sheet \(PDF\)](#)[Comments](#)

Rating: ★★★★★ (6)



saeeddian

[Home](#)[Cheat Sheets](#)[Create](#)[Community](#)[Help](#)[Home](#) > [Programming](#) > [Python Cheat Sheets](#)

Natural Language Processing with Python & nltk Cheat Sheet by [murenei](#)

A quick reference guide for basic (and more advanced) natural language processing tasks in Python, using mostly nltk (the Natural Language Toolkit package), including POS tagging, lemmatizing, sentence parsing and text classification.

Handling Text

<code>text='Some words'</code>	assign string
<code>list(text)</code>	Split text into character tokens
<code>set(text)</code>	Unique tokens
<code>len(text)</code>	Number of characters

Accessing corpora and lexical resources

<code>from nltk.corpus import brown</code>	import CorpusReader object
<code>brown.words(text_id)</code>	Returns pretokenised document as list of words
<code>brown.fileids()</code>	Lists docs in Brown corpus
<code>brown.categories()</code>	Lists categories in Brown corpus

Tokenization

<code>text.split(" ")</code>	Split by space
<code>nltk.word_tokenizer(text)</code>	nltk in-built word tokenizer
<code>nltk.sent_tokenize(doc)</code>	nltk in-built sentence tokenizer

Sentence Parsing

<code>g=nltk.data.load('grammar.cfg')</code>	Load a grammar from a file
<code>g=nltk.CFG.fromstring("...")</code>	Manually define grammar
<code>parser=nltk.ChartParser(g)</code>	Create a parser out of the grammar

<code>trees=parser.parse_all(text)</code>	
<code>for tree in trees: ... print tree</code>	
<code>from nltk.corpus import treebank</code>	
<code>treebank.parsed_sents('wsj_0001.mrg')</code>	Treebank parsed sentences

Text Classification

<code>from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer</code>	
<code>vect=CountVectorizer().fit(X_train)</code>	Fit bag of words model to data
<code>vect.get_feature_names()</code>	Get features
<code>vect.transform(X_train)</code>	Convert to doc-term matrix

[Download This Cheat Sheet \(PDF\)](#)
[Comments](#) | Rating: ★★★★★ (6) | [❤](#)

```
input="List listed lists
listing listings"

words=input.lower().split(' ')

porter=nltk.PorterStemmer

[port.stem(t) for t in
words]

WNL=nltk.WordNetLemmatizer()

[WNL.lemmatize(t) for t in
words]
```

Different suffixes

Normalize (lower-case) words

Initialise Stemmer

Create list of stems

Initialise WordNet lemmatizer

Use the lemmatizer

Part of Speech (POS) Tagging

```
nltk.help.upenn_tagset('MD')

nltk.pos_tag(words)
```

Lookup definition for a POS tag

nltk in-built POS tagger

<use an alternative tagger to illustrate ambiguity>

```
g="NP: {<DT>?<JJ>*<-
NN>}"

cp=nltk.RegexpParser(g)

ch=cp.parse(pos_sent)

print(ch)

ch.draw()

cp.evaluate(test_sents)

sents=nltk.corpus.treebank.tagged_sents()

print(nltk.ne_chunk(sent))
```

Regex chunk grammar

Parse grammar

Parse tagged sent. using grammar

Show chunks

Show chunks in IOB tree

Evaluate against test doc

Print chunk tree

RegEx with Pandas & Named Groups

```
df=pd.DataFrame(time_sents, columns=['text'])

df['text'].str.split().str.len()

df['text'].str.contains('word')

df['text'].str.count(r'\d')

df['text'].str.findall(r'\d')

df['text'].str.replace(r'\w+day\b', '???')

df['text'].str.replace(r'(\w)', lambda x: x.groups()[0][:3])

df['text'].str.extract(r'(\d?\d):(\d\d)')

df['text'].str.extractall(r'((\d?\d):(\d\d) ?([ap]m))')

df['text'].str.extractall(r'(?<digits>\d)')
```



How's Your Readability?

Cheatography is sponsored by Readable.com. Check out Readable to make your content and copy more

Help Us Go Positive!

We offset our carbon usage with Ecologi.

[Download This Cheat Sheet \(PDF\)](#)
[Comments](#)

Rating: ★★★★★ (6)



Measure Your Readability
Now!

Ecology
always
positive
positive

344 trees



Download the Natural Language Processing with Python & nltk Cheat Sheet



PDF (recommended)

[PDF \(2 pages\)](#)

Alternative Downloads

[PDF \(black and white\)](#)

[LaTeX](#)

Comments

No comments yet. Add yours below!

Created By

murenei

<https://tutify.com.au>

Add a Comment

[Download This Cheat Sheet \(PDF\)](#)[Comments](#) | Rating: ★★☆☆☆ (6) | [♥](#)

Your Name

Your Email
Address

Your Comment

[Post Your Comment](#)

Metadata

Languages: [English](#)

Published: 28th May, 2018

Last Updated: 29th May, 2018

Rated: 5 stars based on 6 ratings

Favourited By

Related Cheat Sheets

[Python 3 Cheat Sheet](#)
by Finxter[NLP For Arabic Cheat Sheet](#)

More Cheat Sheets by murenei

★★★★☆

Network Analysis with
Python and NetworkX
Cheat Sheet

Latest Cheat Sheet

AP Physics Formulas (Kinematic) Cheat Sheet

Some physics formulas that will be useful in kinematics. Not a truly complete list of formulas though, as some things are missing. I can't think of any more formulas for this cheat sheet though, so suggestions on what to add would be helpful.

 1 Page

☆☆☆☆☆ (0)



ReSummit



23 Oct 20



physics

Random Cheat Sheet

FREQUENTLY USED DX CODES Cheat Sheet

List of diagnosis frequently used on admission to rehab and long term care facilities.

 1 Page

★★★★★ (2)



charlesnurse



9 May 13, updated 12 May 16



nursing, coding, nurse, admission, icd-9 and 3 more ...

[Download This Cheat Sheet \(PDF\)](#)[Comments](#) | Rating: ★★☆☆☆ (6) | [Heart](#)

references in 25 languages for
everything from travel to maths!

right here:

[DaveChild](#)

[SpaceDuck](#)

[Cheatography](#)

[Sheet.](#)

1 day 9 hours ago

[ReSummit published AP
Physics Formulas
\(Kinematic\).](#)

1 day 11 hours ago

[cjdvslee updated UTS.](#)

1 day 18 hours ago