

---

# **LexNLP Documentation**

***Release 2.2.0***

**ContraxSuite, LLC**

**Jul 12, 2022**



<b>1</b>	<b>Table of Contents</b>	<b>3</b>
1.1	About LexNLP	3
1.1.1	Purpose	3
1.1.2	ContraxSuite Projects	3
1.1.3	Citing for academic use	4
1.2	LexNLP package	4
1.2.1	<code>lexnlp.extract</code> : Extracting structured data from unstructured text	4
1.2.1.1	Pattern-based extraction methods	4
1.2.1.2	NLP-based extraction methods	6
1.2.2	<code>lexnlp.nlp</code> : Natural language processing	6
1.2.2.1	Tokenization and related methods	6
1.2.2.2	Segmentation and related methods for real-world text	7
1.2.2.3	Transforming text into features	7
1.3	Changelog	7
1.3.1	2.2.0 - July 7, 2022	7
1.3.2	2.1.0 - September 16, 2021	7
1.3.3	2.0.0 - May 10, 2021	7
1.3.4	1.8.0 - December 2, 2020	7
1.3.5	1.7.0 - August 27, 2020	8
1.3.6	1.6.0 - May 27, 2020	8
1.3.7	1.4.0 - December 20, 2019	8
1.3.8	1.3.0 - November 1, 2019	8
1.3.9	0.2.7 - August 1, 2019	8
1.3.10	0.2.6 - Jun 12, 2019	8
1.3.11	0.2.5 - Mar 1, 2019	9
1.3.12	0.2.4 - Feb 1, 2019	9
1.3.13	0.2.3 - Jan 10, 2019	9
1.3.14	0.2.2 - Sep 30, 2018	9
1.3.15	0.2.1 - Aug 24, 2018	9
1.3.16	0.2.0 - Aug 1, 2018	10
1.3.17	0.1.9 - Jul 1, 2018	10
1.3.18	0.1.8 - May 1, 2018	10
1.3.19	0.1.7 - Apr 1, 2018	10
1.3.20	0.1.6 - Mar 1, 2018	10
1.3.21	0.1.5 - Feb 1, 2018	10
1.3.22	0.1.4 - Jan 1, 2018	10

1.3.23	0.1.3 - Dec 1, 2017 . . . . .	10
1.3.24	0.1.2 - Nov 1, 2017 . . . . .	11
1.3.25	0.1.1 - Oct 1, 2017 . . . . .	11
1.3.26	0.1.0 - Sep 1, 2017 . . . . .	11
1.4	License . . . . .	12
1.4.1	AGPL License . . . . .	12
1.4.2	License Release . . . . .	12
<b>2</b>	<b>Indices and tables</b>	<b>13</b>





## 1.1 About LexNLP

### 1.1.1 Purpose

LexNLP is a library for working with real, unstructured legal text, including contracts, plans, policies, procedures, and other material. LexNLP provides functionality such as:

#### **Segmentation and tokenization, such as**

- A sentence parser that is aware of common legal abbreviations like LLC. or F.3d.
  - Pre-trained segmentation models for legal concepts such as pages or sections.
  - Pre-trained word embedding and topic models, broadly and for specific practice areas
- Pre-trained classifiers for document type and clause type
- Broad range of fact extraction, such as:
  - Monetary amounts, non-monetary amounts, percentages, ratios
  - Conditional statements and constraints, like “less than” or “later than”
  - Dates, recurring dates, and durations
  - Courts, regulations, and citations
- Tools for building new clustering and classification methods
- Hundreds of unit tests from real legal documents

### 1.1.2 ContraxSuite Projects

LexNLP is often used as part of ContraxSuite, an open source contract analytics and document exploration platform built by LexPredict. LexNLP and ContraxSuite are related through the following project structure:

- ContraxSuite web application: <https://github.com/LexPredict/lexpredict-contraxsuite>

- LexNLP library for extraction: <https://github.com/LexPredict/lexpredict-lexnlp>
- ContraxSuite pre-trained models and “knowledge sets”: <https://github.com/LexPredict/lexpredict-legal-dictionary>
- ContraxSuite agreement samples: <https://github.com/LexPredict/lexpredict-contraxsuite-samples>
- ContraxSuite deployment automation: <https://github.com/LexPredict/lexpredict-contraxsuite-deploy>

### 1.1.3 Citing for academic use

We are currently drafting and submitting a technical whitepaper describing the LexNLP library and documenting its performance on a large corpus of gold-standard contracts. Please contact us at [support@contraxsuite.com](mailto:support@contraxsuite.com) for more information on citing in academic use.

## 1.2 LexNLP package



### 1.2.1 `lexnlp.extract`: Extracting structured data from unstructured text

The `lexnlp.extract` module contains methods that allow for the extraction of structured data from unstructured textual sources. Supported data types include a wide range of facts relevant to contract or document analysis, including dates, amounts, proper noun types, and conditional statements.

**This module is structured along ISO 2-character language codes. Currently, the following languages are stable:**

- English: `lexnlp.extract.en`
- German: `lexnlp.extract.de`
- Spanish: `lexnlp.extract.es`

Extraction methods follow a simple *get\_X* pattern as demonstrated below:

```
>>> import lexnlp.extract.en.amounts
>>> text = "There are ten cows in the 2 acre pasture."
>>> print(list(lexnlp.extract.en.amounts.get_amounts(text)))
[10, 2.0]
```

#### 1.2.1.1 Pattern-based extraction methods

**The full list of supported pattern-based structured data types is below:**

- “EN” locale:
  - acts, e.g., “section 1 of the Advancing Hope Act, 1986”
  - amounts, e.g., “ten pounds” or “5.8 megawatts”



- citations, e.g., “10 U.S. 100” or “1998 S. Ct. 1”
- companies, e.g., “Lexpredict LLC”
- conditions, e.g., “subject to ...” or “unless and until ...”
- constraints, e.g., “no more than” or “
- copyright, e.g., “(C) Copyright 2000 Acme”
- courts, e.g., “Supreme Court of New York”
- CUSIP, e.g., “392690QT3”
- dates, e.g., “June 1, 2017” or “2018-01-01”
- definitions, e.g., “Term shall mean ...”
- distances, e.g., “fifteen miles”
- durations, e.g., “ten years” or “thirty days”
- geographic and geopolitical entities, e.g., “New York” or “Norway”
- money and currency usages, e.g., “\$5” or “10 Euro”
- percents and rates, e.g., “10%” or “50 bps”
- PII, e.g., “212-212-2121” or “999-999-9999”
- ratios, e.g., “3:1” or “four to three”
- regulations, e.g., “32 CFR 170”
- trademarks, e.g., “MyApp (TM)”
- URLs, e.g., “<http://acme.com/>”
- “DE” locale:
  - amounts, e.g., “1 tausend” or “eine halbe Million Dollar”
  - citations, e.g., “BGBI. I S. 434”
  - copyrights, e.g., “siemens.com globale Website Siemens © 1996 – 2019”
  - court citations, e.g., “BStBl I 2003, 240”
  - courts, e.g., “Amtsgerichte”
  - dates, e.g., “vom 29. März 2017”
  - definitions
  - durations, e.g., “14. Lebensjahr” or “fünfundzwanzig Jahren”
  - geographic and geopolitical entities, e.g., “Albanien”
  - percents, e.g., “15 Volumenprozent”
- “ES” locale:
  - copyrights, e.g., “”Website BBC Mundo © 1996 – 2019”
  - courts, e.g., “Tribunal Superior de Justicia de Madrid”
  - dates, e.g., “15 de febrero” or “1º de enero de 1999”
  - definitions, e.g., “”El ser humano”: una anatomía moderna humana”
  - regulations, e.g., “Comisión Nacional Bancaria y de Valores”

### 1.2.1.2 NLP-based extraction methods

In addition to pattern-based structured data types, the *lexnlp.extract* module also supports NLP methods based on tagged part-of-speech classifiers. These classifiers are based on NLTK and, optionally, Stanford NLP libraries. The list of these modules is below:

- named entity extraction with NLTK maximum entropy classifier
- named entity extraction with NLTK and regular expressions
- named entity extraction with Stanford Named Entity Recognition (NER) models

**These modules allow to extract data types like:**

- addresses, e.g., “1999 Mount Read Blvd, Rochester, NY, USA, 14615”
- companies, e.g., “Lexpredict LLC”
- persons, e.g., “John Doe”

## 1.2.2 `lexnlp.nlp`: Natural language processing

The `lexnlp.nlp` module contains methods that assist in natural language processing (NLP) tasks, especially in the context of developing unsupervised, semi-supervised, or supervised machine learning. Methods range from tokenizing, stemming, and lemmatizing to the creation of custom sentence segmentation or word embedding models.

**This module is structured along ISO 2-character language codes. Currently, the following languages are stable:**

- English: *lexnlp.nlp.en*

Extraction methods follow a simple *get\_X* pattern as demonstrated below:

```
>>> import lexnlp.nlp.en.tokens
>>> text = "There are ten cows in the 2 acre pasture."
>>> print(list(lexnlp.nlp.en.tokens.get_nouns(text)))
['cows', 'pasture']
```

The methods in this package are primarily built on the Natural Language Toolkit (NLTK), but some functionality from the Stanford NLP, gensim, and spaCy packages is available to users depending on their use case.

**Attention:** The sections below are a work in progress. Thank you for your patience while we continue to expand and improve our documentation coverage.

If you have any questions in the meantime, please feel free to log issues on GitHub at the URL below or contact us at the email below:

- GitHub issues: <https://github.com/LexPredict/lexpredict-lexnlp>
- Email: [support@contraxsuite.com](mailto:support@contraxsuite.com)

### 1.2.2.1 Tokenization and related methods

- Extracting tokens, stems, lemmas, and parts of speech

### 1.2.2.2 Segmentation and related methods for real-world text

- Sentences
- Paragraphs
- Sections
- Pages
- Titles
- Utilities

### 1.2.2.3 Transforming text into features

- Character Transforms
- Token Transforms, including n-grams and skip-grams

## 1.3 Changelog

### 1.3.1 2.2.0 - July 7, 2022

- Improved LexNLP handling for dates, durations and persons for the all locales.
- Added parameterizable contract classifiers.
- Improved LexNLP handling for ML models.
- Updated python requirements and tests, retrained ML models to use gensim-4.

### 1.3.2 2.1.0 - September 16, 2021

- Improved LexNLP handling for companies for the “EN” locale.
- Improved LexNLP handling for dates for all locales and dates parser accuracy for the “DE” locale.

### 1.3.3 2.0.0 - May 10, 2021

- Tune extracting facts from text for different locales.
- Updated regex patterns and tests for “DE” amount parser.
- Added support for delimiter inference in “EN” and “DE” amount parsers.

### 1.3.4 1.8.0 - December 2, 2020

- Improved LexNLP handling for definitions for the “EN” locale.
- Implemented rating OCR quality in texts.
- Migrated numeric data in parsers results to decimal format to avoid losing fraction digits.

### 1.3.5 1.7.0 - August 27, 2020

- Improved LexNLP handling for dates for the “EN” locale.
- Implemented lists of exceptions for entity extractors.
- Implemented strongly typed response for entity extractors.
- Updated third-party python requirements.

### 1.3.6 1.6.0 - May 27, 2020

- Update psutil package version from 5.4.0 to 5.6.6.

### 1.3.7 1.4.0 - December 20, 2019

- Improved accuracy of locating and converting date phrases into typed format.
- Introduced new text vectorizing and classifying models.
- Implemented ML-based definitions locator.

### 1.3.8 1.3.0 - November 1, 2019

- Made massive improvements to EN definitions and companies parsers.
- Updated EN dates parser to catch more date formats.
- Made company parsing strongly typed

### 1.3.9 0.2.7 - August 1, 2019

- Standardized LexNLP methods response to return a generator of Annotation objects or a generator of dictionaries (tuples)
- Improved LexNLP handling for definitions for the “EN” locale.
- Improved LexNLP handling for companies for the “EN” locale.
- Improved sentence splitting logic.
- Improved LexNLP unit test coverage.
- Updated python requirements in python-requirements\*.txt.
- Dropped support for python 3.4 and 3.5.

### 1.3.10 0.2.6 - Jun 12, 2019

- Improved LexNLP handling for dates for all locales.
- Improved LexNLP handling for currencies for “EN” locale.
- Updated documentation for ReadTheDocs.
- Improved LexNLP unit test coverage.

### **1.3.11 0.2.5 - Mar 1, 2019**

- Improved LexNLP handling for courts for “DE” and “ES” locales.
- Improved LexNLP handling for dates for “ES” locale.
- Improved LexNLP handling for amounts, acts, regulations and definitions for “EN” locale.
- Added CUSIP parser for “EN” locale.
- Improved LexNLP unit test coverage.

### **1.3.12 0.2.4 - Feb 1, 2019**

- Added universal courts parser, configured LexNLP handling for courts for “DE” locale.
- Added universal dates parser, configured LexNLP handling for dates for “DE” and “ES” locales.
- Added definitions, citations and dates parsers for “DE” locale.
- Added amounts, percents and durations parsers for “DE” locale.
- Added geo entities parser for “DE” locale.
- Added courts and definitions parsers for “ES” locale.
- Added acts parser for “EN” locale.
- Improved LexNLP unit test coverage.

### **1.3.13 0.2.3 - Jan 10, 2019**

- Updated python requirements.
- Improved LexNLP handling for definitions and paragraphs.
- Improved LexNLP unit test coverage.

### **1.3.14 0.2.2 - Sep 30, 2018**

- Improved LexNLP handling for different date formats.
- Improved LexNLP handling for titles.
- Improved LexNLP unit test coverage.

### **1.3.15 0.2.1 - Aug 24, 2018**

- Updated python requirements.
- Improved LexNLP handling for amounts.
- Optimized processing of sentences and titles.
- Improved LexNLP unit test coverage.

### **1.3.16 0.2.0 - Aug 1, 2018**

- Improved LexNLP handling for addresses and sentences.
- Improved LexNLP unit test coverage.

### **1.3.17 0.1.9 - Jul 1, 2018**

- Improved handling of TOC during sentence processing.
- Added contracts locator to LexNLP.
- Improved LexNLP handling for citations, titles and definitions.
- Improved LexNLP unit test coverage.

### **1.3.18 0.1.8 - May 1, 2018**

- Improved LexNLP handling for addresses and currencies.
- Improved LexNLP unit test coverage.

### **1.3.19 0.1.7 - Apr 1, 2018**

- Improved LexNLP handling for companies, organizations and dates.
- Implemented generating train/test dataset for addresses.
- Exclude common false positives for persons parser.

### **1.3.20 0.1.6 - Mar 1, 2018**

- Improved LexNLP unit test coverage.

### **1.3.21 0.1.5 - Feb 1, 2018**

- Improved LexNLP unit test coverage.

### **1.3.22 0.1.4 - Jan 1, 2018**

- Improved LexNLP unit test coverage.
- Implemented method to get sentence ranges in addition to sentence texts.

### **1.3.23 0.1.3 - Dec 1, 2017**

- Improved LexNLP unit test coverage.

### 1.3.24 0.1.2 - Nov 1, 2017

- Implemented LexNLP title locator.
- Implemented additional LexNLP transforms for skipgrams and n-grams.
- Improved LexNLP handling for parties with abbreviations and other cases.
- Improved LexNLP handling for amounts with mixed alpha and numeric characters.
- Improved LexNLP unit test coverage.

### 1.3.25 0.1.1 - Oct 1, 2017

- Improve unit test framework handling for language and locales.
- Implemented method and input-level CPU and memory benchmarking for unit tests.
- Migrated all unit tests to 60 separate CSV files.
- Added over 1,000 new unit tests for most LexNLP methods.
- Reduced memory usage for paragraph and section segmenters.
- Improved handling of brackets and parentheses within noun phrases.
- Added URL locator to LexNLP.
- Added trademark locator to LexNLP.
- Added copyright locator to LexNLP.
- Improved default Punkt sentence boundary detection.
- Added custom sentence boundary training methods.
- Improved handling of multilingual text, especially around geopolitical entities.
- Improved default handling of party names with non-standard characters.
- Enhanced metadata related to party type in LexNLP.
- Improved continuous integration for public repositories.

### 1.3.26 0.1.0 - Sep 1, 2017

- Refactored and integrate core extraction into separate LexNLP package.
- Released nearly 200 unit tests with over 500 real-world test cases in LexNLP.
- Improved definition, date, and financial amount locators for corner cases.
- Integrated PII locator for phone numbers, SSNs, and names from LexNLP.
- Integrated ratio locator from LexNLP.
- Integrated percent locator from LexNLP.
- Integrated regulatory locator from LexNLP.
- Integrated distance locator from LexNLP.
- Integrated case citation locator from LexNLP.
- Improved geopolitical locator to allow non-master-data entity location.

- Improved party locator to allow configuration and better handle corner cases

## 1.4 License

### 1.4.1 AGPL License

LexNLP is available by default under the terms of the GNU Affero General Public License v3.0. <https://github.com/LexPredict/lexpredict-lexnlp/blob/2.2.0/LICENSE>

### 1.4.2 License Release

If you would like to request a release from the terms of the default AGPLv3 license, please contact us at: ContraxSuite Licensing <[license@contraxsuite.com](mailto:license@contraxsuite.com)>.



## CHAPTER 2

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`