

Glossary

Tokenization

Segmenting text into words, punctuation etc.

Lemmatization

Assigning the base forms of words, for example: "was" → "be" or "rats" → "rat".

Sentence Boundary Detection

Finding and segmenting individual sentences.

Part-of-speech (POS) Tagging

Assigning word types to tokens like verb or noun.

Dependency Parsing

Assigning syntactic dependency labels, describing the relations between individual tokens, like subject or object.

Named Entity Recognition (NER)

Labeling named "real-world" objects, like persons, companies or locations.

Text Classification

Assigning categories or labels to a whole document, or parts of a document.

Statistical model

Process for making predictions based on examples.

Training

Updating a statistical model with new examples.