

Terms in text preprocessing:

1. Corpus \rightarrow paragraph.
2. Document \rightarrow sentences from corpus.
3. vocabulary \rightarrow set of unique words from document.
4. Words \rightarrow All the words from corpus

Tokenization :

Corpus \rightarrow $\left\{ \begin{array}{l} \text{This is DL class and we are covering NLP.} \\ \text{It is organized by learnbay. Taken by suresh} \end{array} \right\}$

\downarrow

sentence tokenization \Rightarrow corpus \Rightarrow document

\downarrow

document \rightarrow $\left\{ \begin{array}{l} \text{This is DL class and we are covering NLP.} \Rightarrow \text{document-1} \\ \text{It is organized by learnbay.} \Rightarrow \text{document-2} \\ \text{Taken by suresh.} \Rightarrow \text{document-3} \end{array} \right.$

\downarrow

word tokenization \rightarrow document \Rightarrow words.

\downarrow

words

$\left\{ \begin{array}{l} \text{The cat drink the milk. The dog chases the cat to} \\ \text{drink the milk.} \end{array} \right\} \Rightarrow \text{corpus}$

The cat drink the milk. \Rightarrow document-1 (d-1)

The dog chases the cat to drink the milk. \Rightarrow document-2 (d-2)

[The, cat, drink, the, milk] \Rightarrow words from d-1

[The, dog, chases, the, cat, to, drink, the, milk] \Rightarrow words from d-2

[The, cat, drink, the, milk, dog, chases, to] \Rightarrow vocabulary