

Basic Terminology in NLP

- ① Corpus
- ② Documents
- ③ Vocabulary
- ④ Words

Corpus

- ① Corpus is large collection of text documents
- ② Example: If we are analysing customer reviews for product

Corpus = ["The product is great",
"I love this product",
"product did.... expectation"]

Document

A document is an individual piece of text within the corpus

Corpus = [$\frac{R_1}{R_2}$, $\frac{R_3}{R_4}$]

Vocabulary

It is the set of unique words or tokens in the entire corpus without duplications

Vocabulary = {
"The", "product", "is", "great" . . .
- - - - -
}

Words (tokens)

It is the individual components of document