



What is Bag of Words (BoW) in NLP

- ① widely used technique in NLP to represent text data into vectors
- ② In this model document is represented as a collection (or bag) of its words without considering grammar or word order but keeping track of frequency of each word
- ③ The core idea is to create a vector that counts the number of occurrences of each word in the document

1. Cat eats fish
2. dog eats fish
3. dog chases cat

→ Remove stop words

## 1. Create Vocabulary

[Cat, eats, fish, dog, chases] → unique words

## 2. Create frequency vector for each sentence

Sentence 1 → "Cat eats fish"

Cat → 1, eats → 1, fish → 1, dog → 0, chases → 0

vector :- [1, 1, 1, 0, 0]

Sentence 2 → "dog eats fish"

vector → [0, 1, 1, 1, 0]

Sentence 3 → "dog chases cat"

vector → [1, 0, 0, 1, 1]

Binary Bow  
{ 1 and 0 }

Bow  
{ update based on frequency }

Dataset

~~the~~ ~~the~~ ~~the~~ good boy  
~~the~~ ~~the~~ ~~the~~ good girl  
Boy ~~the~~ girl ~~the~~ good

lower all  
sentence  
for stopword

s1 → good boy  
s2 → good girl boy girl  
s3 → boy girl good

<u>Vocabulary</u>	<u>frequency</u>
good	3
boy	2
girl	2

$\begin{bmatrix} \text{good} & \text{boy} & \text{girl} \end{bmatrix}$   
 $s1 \rightarrow \begin{bmatrix} 1 & 1 & 0 \end{bmatrix}$   
 $s2 \rightarrow \begin{bmatrix} 1 & 0 & 2 \end{bmatrix}$   
 $s3 \rightarrow \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$

$\begin{bmatrix} 1 & 2 & 2 \end{bmatrix}$

Advantages

- ① simple and intuitive
- ② Fixed sized I/p

Disadvantages

- ① sparse matrix & overfitting
- ② ordering of the word is getting changed
- ③ semantic meaning is still not captured
- ④ out of vocabulary (OOV)

## N-grams in Bag of Words

→ An  $n$ -gram is a contiguous sequence of  $n$  items (words, characters)  
from given sample of text

→  $n$  grams are typically sequence of words

unigram ( $n=1$ ) : each individual word is treated as feature

Bigram ( $n=2$ ) pair of consecutive words is treated as feature

Trigram ( $n=3$ ) triplets of consecutive words are treated as feature

① → Cats <sup>②</sup> eats fish.

② → dog eats fish

③ → dog chases cat

unigram ["cat", "eats", "fish"]

Bigram

["cat eats", "eats fish"]