

Bag of words Intuition

Binary

Normal

S₁ → ~~He~~ is a good boy.

stop keywords

S₁

good boy

S₂ → she is a good girl.

removing

S₂

good girl

S₃ → Boy and girl are good.

sentences

S₃

Boy girl good.

words frequency.

good 3
boy 2
girl 2

vectors

Bag of words

we need to sort it in descending order.

independent features
f₁ f₂ f₃
dependent features

| | good | boy | girl | o/p |
|----------------|------|-----|------|-----|
| S ₁ | 1 | 1 | 0 | |
| S ₂ | 1 | 0 | 1 | |
| S ₃ | 1 | 1 | 1 | |

histogram

```
import re
```

↳ regular expression.

```
from import nltk.stem import PorterStemmer
```

```
from nltk.stem import WordNetLemmatizer
```

```
from nltk.corpus import stopwords
```

cleaning the text.

```
ps = PorterStemmer()
```

```
lr = WordNetLemmatizer()
```

```
sentences = paragraph nltk.sent_tokenize(paragraph)
```

```
corpus = []
```

```
for i in range(len(sentences))
```

```
    review = re.sub('[^a-zA-Z]', '', sentences[i])
```

```
    review = review.lower()
```

```
    review = review.split()
```

```
    review = [p.stem(word) for word in review if  
               not word in set(stopwords.words('english'))]
```

```
    review = ' '.join(review)
```

```
    corpus.append(review)
```

Creating the bag of words model.

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
cv = CountVectorizer(max_features=1500)
```

```
x = cv.fit_transform(corpus).toarray()
```