

The background of the slide is a dark gray to black gradient, overlaid with a complex, white, abstract network of interconnected nodes and lines. The nodes are represented by small circles of varying sizes, and the lines are thin, creating a web-like structure that spans the entire frame. The density of the network is higher in the upper right and lower right areas, with more sparse connections in the upper left.

Introduction to NLP

Video 6: Methods of Text Preprocessing - Part 3

Preprocessing Techniques



Lowercasing



Removing Punctuation and Special Characters



Stop Words Removal

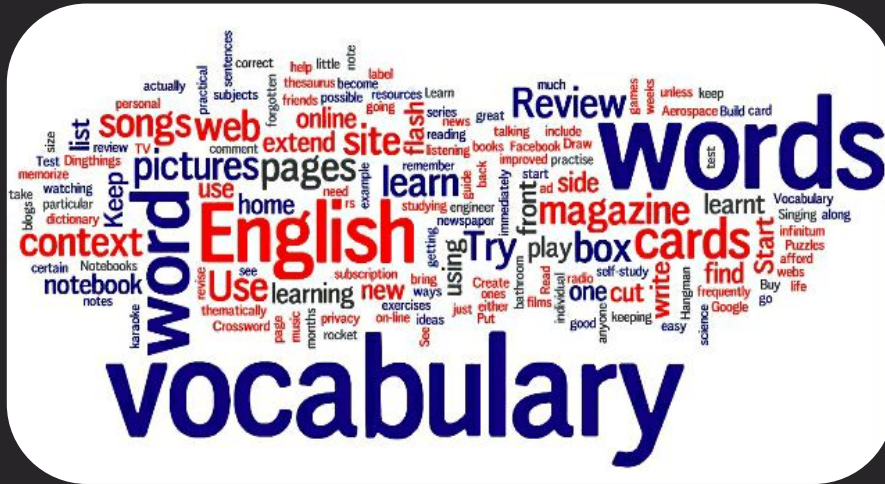


Stemming and Lemmatization

Vocabulary Building

Vectorization

Vocabulary Building



Select and define the set of unique words NLP model will analyse.

Tokens left after preprocessing
steps : Base for the model

Jupyter Notebook

Vectorization

Vectorization

```
of': 0.058554537981986046  
'member': 0.053946576476573944  
'popular': 0.03960049197053909  
'is': 0.038850940465688705  
'an': 0.03689520010328293  
'one': 0.021529447048783302  
'love': 0.012194252967238426  
'helps': -0.005911458611011982  
'important': -0.011899725027620792  
'spiderman': -0.03565903055644035
```

Convert words into numerical representations called vectors.

Vectorization



One Hot Encoding

Count Vectorizer

TF-IDF

Word Embedding

Vectorization



One Hot Encoding

Count Vectorizer

TF-IDF

Word Embedding

Vectorization



One Hot Encoding

Count Vectorizer

TF-IDF

Word Embedding

One-Hot Encoding

Binary vector to show whether a word appears in a document.

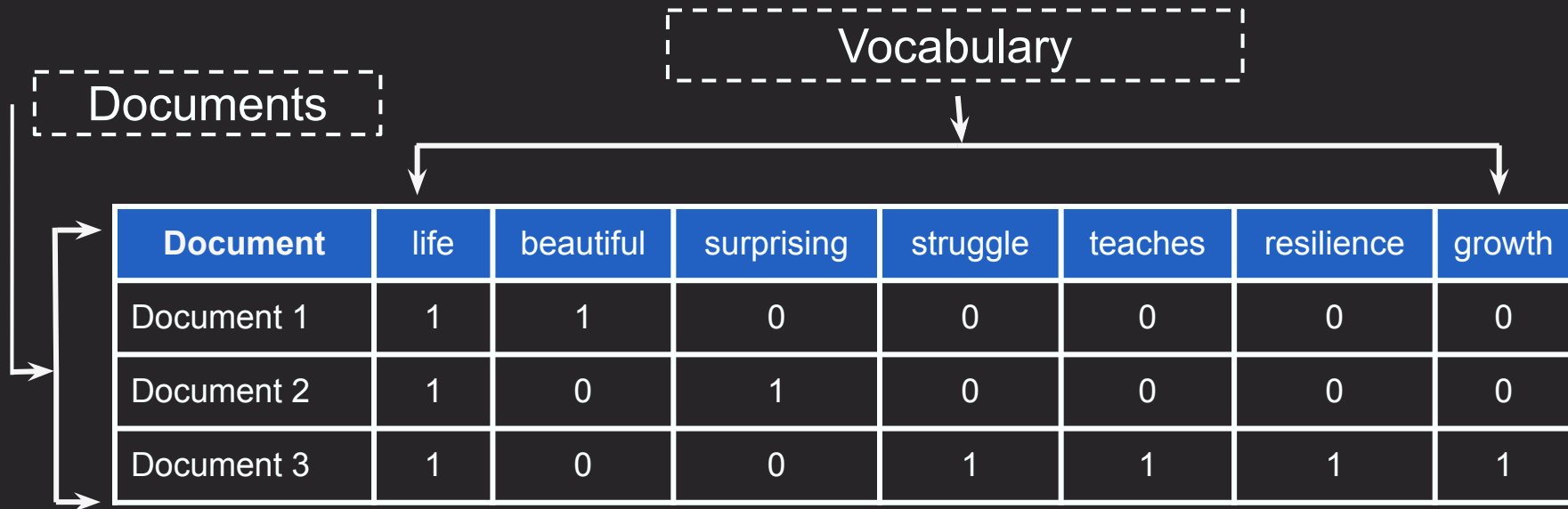
1

Word is present

0

Word is absent

One Hot Encoding



One Hot Encoding

Document 1 : “Life is beautiful.”

Document 2 : “ Life is surprising.”

Document 3 : “Struggle in life teaches resilience and growth.”

One Hot Encoding

Document 1 : “**Life** is **beautiful**.”

Document 2 : “ Life is surprising.”

Document 3 : “Struggle in life teaches resilience and growth.”

life

beautiful

One Hot Encoding

Document 1 : “**Life** is **beautiful**.”

Document 2 : “ Life is **surprising**.”

Document 3 : “Struggle in life teaches resilience and growth.”

life

beautiful

surprising

One Hot Encoding

Document 1 : “**Life** is **beautiful**.”

Document 2 : “ Life is **surprising**.”

Document 3 : “**Struggle** in life **teaches resilience** and **growth**.”

life

beautiful

surprising

struggle

teaches

resilience

growth

One Hot Encoding

Document	life	beautiful	surprising	struggle	teaches	resilience	growth
Document 1	1	1	0	0	0	0	0
Document 2	1	0	1	0	0	0	0
Document 3	1	0	0	1	1	1	1

life

beautiful

surprising

struggle

teaches

resilience

growth

Jupyter Notebook

Vectorization



One Hot Encoding

Count Vectorization

TF-IDF

Word Embedding

Count Vectorization

Represented as matrix of token counts.

Document	life	beautiful	surprising	struggle	turn	teaches	resilience	growth
Document 1	3	2	0	0	1	0	0	0
Document 2	2	0	1	0	0	0	0	0
Document 3	3	0	0	1	1	1	1	1

Count Vectorization

Document 1 : “Life is life, life is beautiful and full of beautiful turns.”

Document 2 : “ Life surprises us at every turn of life”

Document 3 : “Struggle in life teaches resilience, life, and growth in life.”

Count Vectorization

Document 1 : “**Life** is **life**, **life** is **beautiful** and full **beautiful turns**.”

Document 2 : “ **Life surprises** us at every **turn** of **life**”

Document 3 : “**Struggle** in **life teaches resilience**, **life**, and **growth** in **life**.”

life

beautiful

surprises

struggle

turn

teaches

resilience

growth

Count Vectorization

Document	life	beautiful	surprising	struggle	turn	teaches	resilience	growth
Document 1	3	2	0	0	1	0	0	0
Document 2	2	0	1	0	0	0	0	0
Document 3	3	0	0	1	1	1	1	1

life

beautiful

surprises

struggle

turn

teaches

resilience

growth

Jupyter Notebook