

Lesson 5.06

NLP PART 2

Vectorization Types

- **Count Vectorizer** creates a document-term matrix where the entry of each cell will be a count of the number of times that word occurred in that document.
- **TF-IDF (Term Frequency — Inverse Document Frequency)** is basically a count vectorizer that includes some consideration for the length of the document, and also how common the word is across other text messages.
- **N-Grams** is used to look for groups of adjacent words instead of just looking for single terms.
- They're all just slight modifications of each other, and typically you'll test different vectorization methods depending on your problem, and the results determine which one you will proceed with.

TF-IDF Formula

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

$tf_{i,j}$ = number of occurrences of i in j

df_i = number of documents containing i

N = total number of documents

	Document 1	Document 2	Document 3	Document 4	Document 5	Document 6	Document 7	Document 8
Term(s) 1	10	0	1	0	0	0	0	2
Term(s) 2	0	2	0	0	0	18	0	2
Term(s) 3	0	0	0	0	0	0	0	2
Term(s) 4	6	0	0	4	6	0	0	0
Term(s) 5	0	0	0	0	0	0	0	2
Term(s) 6	0	0	1	0	0	1	0	0
Term(s) 7	0	1	8	0	0	0	0	0
Term(s) 8	0	0	0	0	0	3	0	0

Word Vector (Passage Vector) ←

Document Vector ↗

Document Term Matrix

- A document-term matrix is a mathematical matrix that describes the frequency of terms that occur in a collection of documents.
- In a document-term matrix, rows correspond to documents in the collection and columns correspond to terms.

Sparse Matrix

- A matrix comprising of mostly zero values
- Only locations and values non-zero elements are stored to save space

1	[[1 0 0 1 0 0]		
2	[0 0 2 0 0 1]		
3	[0 0 0 2 0 0]]		
4			
5	(0, 0)	1	
6	(0, 3)	1	
7	(1, 2)	2	
8	(1, 5)	1	
9	(2, 3)	2	
10			

Normal Matrix
(All values)

Sparse Matrix
(Location and Values)