

Lesson 23 (Bag of Words) Term (word) Frequency (**TF**) is the word count of a word in a document, normalized by the total number of words in the document. Consider the corpus consisting of documents that are the sentences listed below. Fill in the tables below after removing stop words. Sort words in alphabetical order. Leave zero values blank.

Document		
4	0	The <u>car</u> <u>crashed</u> <u>long</u> ago.
2	1	The <u>car</u> has <u>rusted</u> .
3	2	A <u>rusted</u> <u>car</u> is <u>unsafe</u> .
5	3	<u>Spare</u> <u>car</u> <u>parts</u> are <u>needed</u> <u>urgently</u> .

unique words sort in alphabetical order

don't need to fill 0's for empty

(a) Document Term Matrix:

terms (words)	<u>ago</u>	<u>car</u>	<u>crashed</u>	<u>long</u>	<u>needed</u>	<u>parts</u>	<u>rusted</u>	<u>spare</u>	<u>unsafe</u>	<u>urgently</u>
Document 0	1	1	1	1						
1		1					1			
2		1					1		1	
3		1			1	1		1		1

(b) TF (Term Frequency):

terms (words)	<u>ago</u>	<u>car</u>	<u>crashed</u>	<u>long</u>	<u>needed</u>	<u>parts</u>	<u>rusted</u>	<u>spare</u>	<u>unsafe</u>	<u>urgently</u>
4 Document 0	1/4	1/4	1/4	1/4						
2 1		1/2					1/2			
3 2		1/3					1/3		1/3	
5 3		1/5			1/5	1/5		1/5		1/5