

Bag of Words

Sentence \rightarrow S1 \rightarrow He is a good boy
 S2 \rightarrow She is a good girl
 S3 \rightarrow Boy & Girl are good

Step 1:- Stop Keywords & lower the sentences :-
stop keywords like \rightarrow He, She, is, a are not much important
 the keywords \rightarrow boy, girl, good are important
Lower the sentences \rightarrow Convert uppercase to lowercase

S1 good boy
 S2 good girl
 S3 boy girl good

Step 2:- Create histogram \rightarrow write the ^(count) frequency for the important keywords

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

note: when we calculate frequency, it should be in sorted order.

Step 3:- now apply bag of words (BOW):

		f_1	f_2	f_3	o/p
		good	boy	girl	
Vectors \Rightarrow BOW	S ₁	1	1	0	
	S ₂	1	0	1	
	S ₃	1	1	0	

binary bag of words (0 & 1)

When we apply ML Model; f_1, f_2 & f_3 are independent features & o/p is dependent feature and train the ML Model.

Here, good-1 & boy-1; both words have same representation; but we don't know which word gives more importance. For this, we use the technique called TF-IDF (Term Frequency-Inverse Document Frequency)

TF-IDF

→ Term frequency (TF) = $\frac{\text{no. of repetition of words in sentence}}{\text{no. of words in sentence}}$

→ Inverse-Document frequency (IDF) = $\log\left(\frac{\text{no. of sentences}}{\text{no. of sentences containing words}}\right)$

→ Then calculate $TF * IDF$.

Given:-

		words	frequency
S1 →	good boy	good	3
S2 →	good girl	boy	2
S3 →	good boy girl	girl	2

Calculate TF:-

Calculate IDF:-

	S1	S2	S3	words	IDF
good	1/2	1/2	1/3	good	$\log(3/3) = 0$
boy	1/2	0	1/3	boy	$\log(3/2)$
girl	0	1/2	1/3	girl	$\log(3/2)$

now Calculate: Tf*IDF:-

		f_1	f_2	f_3	o/p
		good	boy	girl	
vectors \Rightarrow	S_1	0	$\frac{1}{2} \log(3/2)$	0	
Tf-IDF	S_2	0	0	$\frac{1}{2} \log(3/2)$	
	S_3	0	$\frac{1}{3} \log(3/2)$	$\frac{1}{3} \log(3/2)$	

Calculation for above table:-

	f_1	f_2	f_3
	good	boy	girl
S_1	$\frac{1}{2} * 0 = 0$	$\frac{1}{2} * \log(3/2) = \frac{1}{2} \log(3/2)$	$0 * \log(3/2) = 0$
S_2	$\frac{1}{2} * 0 = 0$	$0 * \log(3/2) = 0$	$\frac{1}{2} \log(3/2) = \frac{1}{2} * \log(3/2)$
S_3	$\frac{1}{3} \log(3/2) = 0$	$\frac{1}{3} * \log(3/2) = \frac{1}{3} \log(3/2)$	$\frac{1}{3} \log(3/2) = \frac{1}{3} \log(3/2)$

Here, we will get decimals which gives Schematic importance to boy, girl words.