



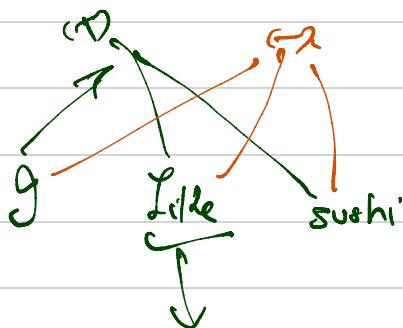
## Agenda

1. BOW recap + discussion
2. Use-case Medium Articles
3. Use BOW to generate recommendations
4. TF-IDF, use TF-IDF to generate recommendations

Stemming & Lemmatization) → Sent classification

use this or not

Problem statement



for each word  
Find out freq

across entire D → (1) v

(2) v

# Bag of Words

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

47K → daily → 200M Visitor  
↳ every month

→ Contractions

can't → can not  
don't → do not  
you'll → you will

Vocabulary → Set of all unique words

good, nice, great

problems collection of all documents

corpus

Step 2 - Each word in the sentence would be represented as below:

Word	Lecture	on	text	representation
Lecture	1	0	0	0
on	0	1	0	0
text	0	0	1	0
representation	0	0	0	1

- 1 Distance
- 2 doesn't Preve similar.
- 3 Order

### Quiz time!

Time Left: 13s

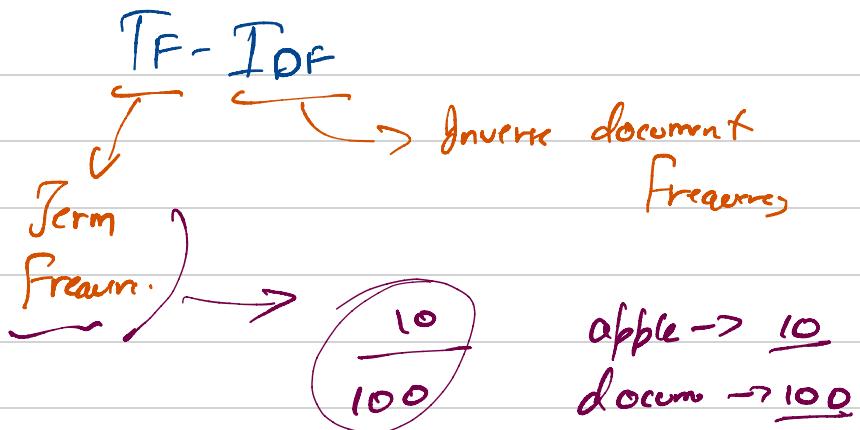
What is a potential limitation of the Bag of Words model?

21 users have participated

- A It is computationally expensive. 14%
- B It does not consider the order of words. 76% ✓
- C It requires labeled training data. 5%
- D It is only applicable to short documents. 5%

[End Quiz Now](#)

X  
X  
long ones



IDF  $\rightarrow$

$d_1$

$d_1 - d_{10}$

$= 2$

$D \rightarrow d_1 \rightarrow d_{10}$

"Harry Potter"  $\rightarrow 5 \rightarrow d_1$

$\rightarrow d_2, d_6, d_7, d_8, d_9, d_{10}$

$$\log\left(\frac{5}{8}\right) \approx 1$$

IDF  $\rightarrow$  rarity of word

$$D \rightarrow d_1 - \underbrace{d_{20}}$$

"mouse"  $\rightarrow$   $d_1$   $\rightarrow$  10

$$\begin{aligned} & d_2, d_4, d_5, \dots, d_{19} \\ & \underbrace{10}_{\text{other documents}} \quad \underbrace{\text{mouse}}_{\text{document}} = \log\left(\frac{10}{10}\right) \\ & = \cancel{\log(1)} \\ & = 0 \end{aligned}$$

If wasn't  $P_x = d_2 - d_p$

$$\log\left(\frac{10}{1}\right) =$$

1  
how  $\rightarrow$  rarity

Document matrix

$\hookrightarrow$  Replace Frequency  $\rightarrow$

TF  $\times$  IDF

In the TF-IDF formula, what does the Inverse Document Frequency (IDF) component measure? ?

22 users have participated

- |   |   |     |
|---|---|-----|
| A | The frequency of a term in a document                             | 0%  |
| B | The importance of a term in a specific document                   | 1%  |
| C | <b>The rarity of a term in the entire collection of documents</b> | 82% |
| D | The total number of documents in the collection                   | 0%  |

[End Quiz Now](#)

TF-IDF  
↓  
Sentiment

$d_1$

mulchond  $\rightarrow \frac{10}{100}$ ) Euclidian

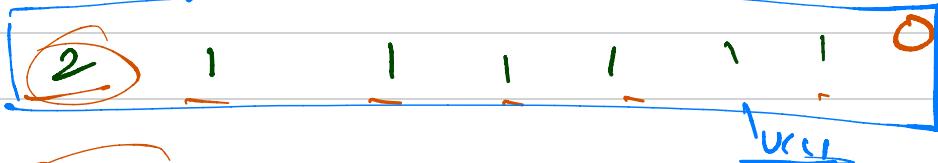
$d_2$

Recall  $\frac{100}{1000}$ ) Euclidean

NDCG  
MRR  
Jaccard - accuracy  
brown

A quick brown fox jumped over a lazy dog

Bow



TF-IDF

TFxIDF

Vec<sub>2</sub>