

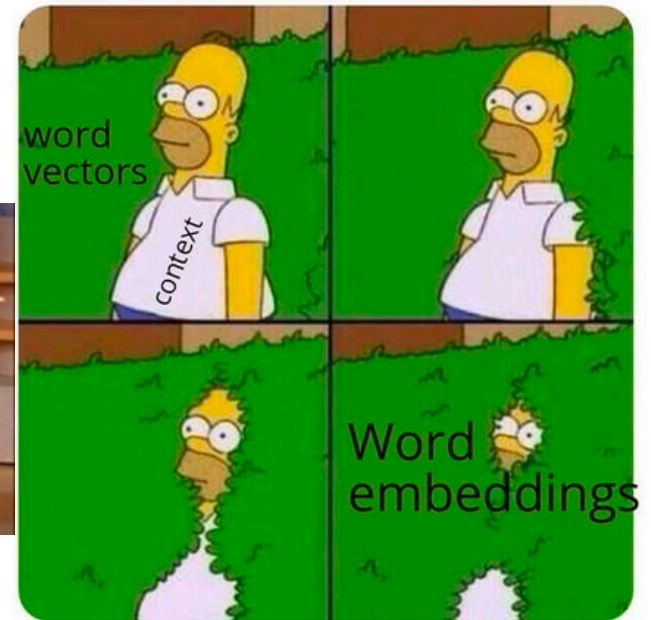
April 21, 2023 .

DSML : NLP module .

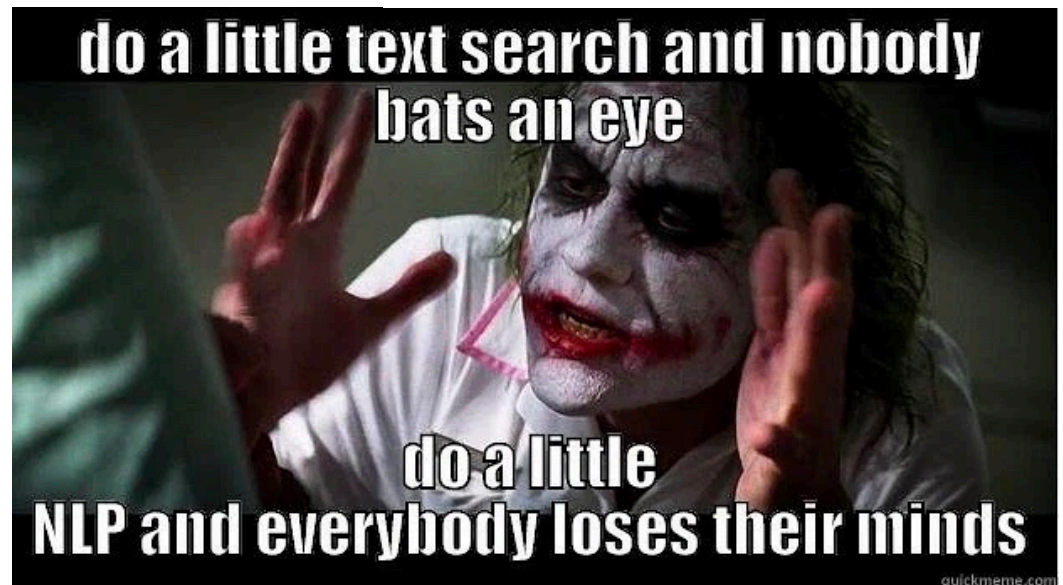
Word embeddings
in a nutshell

Text Representations.

Class starts
@ 9:05



When you penalize your Natural
Language Generation model for
large sentence lengths



Announcements :

Hey Everyone,

On 22nd April at 03:00 PM IST, we have organised a session with one of your peers where we will be discussing on the roadmap to crack Data Analyst Offers.

More details in the link to register.

Link to register - <https://bit.ly/3UXzE2P>

Regards,

Scaler Community

From learners: Could we have a similar event for Data Scientist roles?

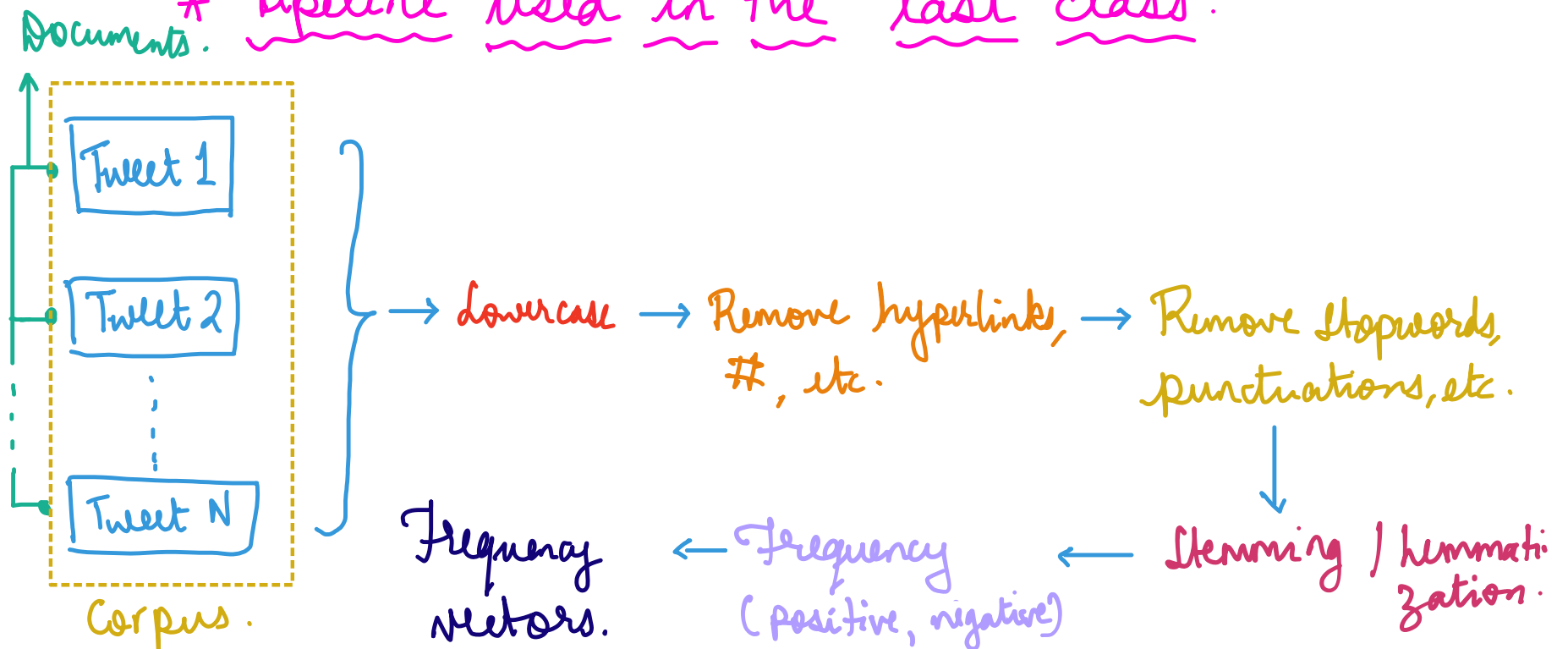
Recap:

* NLP → Natural language processing.

* Core Issue: Preprocessing.

"How to convert text to vectors?"

* Pipeline used in the last class:



Agenda: Text \rightarrow Numerical vectors.

- BoW: Bag of words.

- TF-IDF: Term frequency, Inverse Document Frequency.

- Call: Finding similar blogs.

- T-SNE for data insights.

- Case : Finding similar blogs.

Finding similar Medium articles

You are working as a Data Scientist at Medium

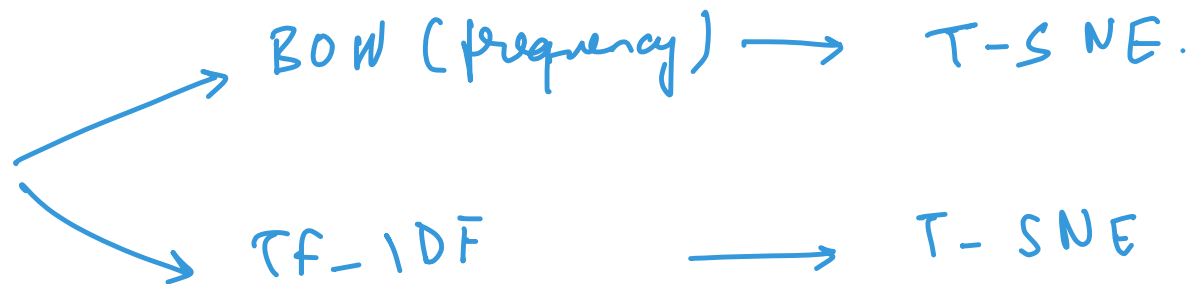
- [Medium](#) is an online publishing platform which hosts a hybrid collection of blog posts from both amateur and professional people and publications.
- In 2020, about 47,000 articles were published daily on the platform and it had about 200M visitors every month.

Problem Statement:

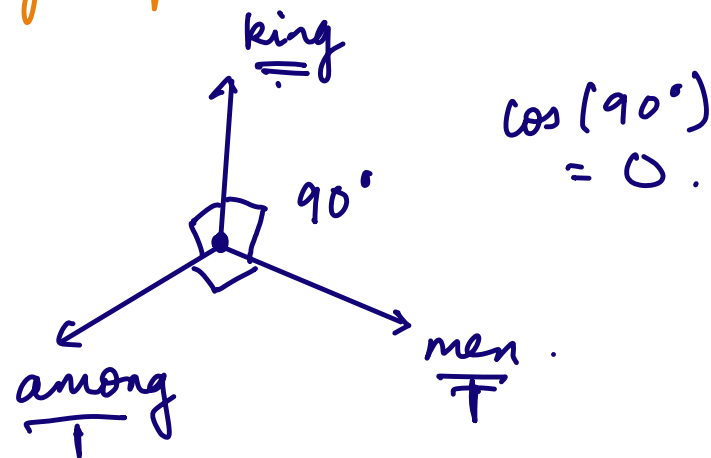
- You want to give readers a better reading experience at Medium. To do that, you want to recommend articles to the user on the basis of current article that the user is reading.
- More concretely, given a Medium article find a set of similar articles.

Corpus

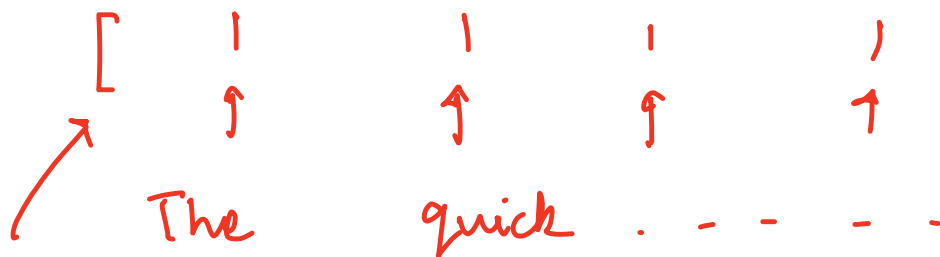
d_1
 d_2
 d_3
 \dots 209



The quick brown fox jumped over the lazy dog.

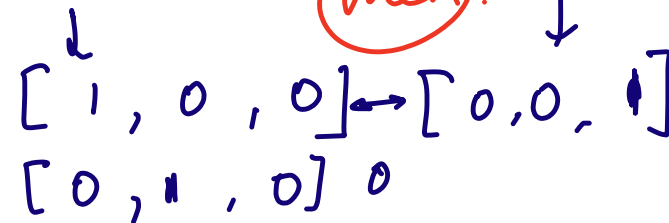


One-hot encoding.



- sparsity.
- All words have weightage!
- Similar words are not close by.

He is a king among men.



Bag-of-words representations.

Corpus - doc1
doc2
⋮
docn.

- ⑥ pre-processing pipeline - lowercase, stop-word removal, stem/lem etc.
- ① get all unique words from the corpus.
- ② Put these words as the column names.
- ③ For each doc, find frequency & store in matrix.

Unresolved issues with BOW:

- ① Weightage \rightarrow frequency \times
- ② cosine similarity.

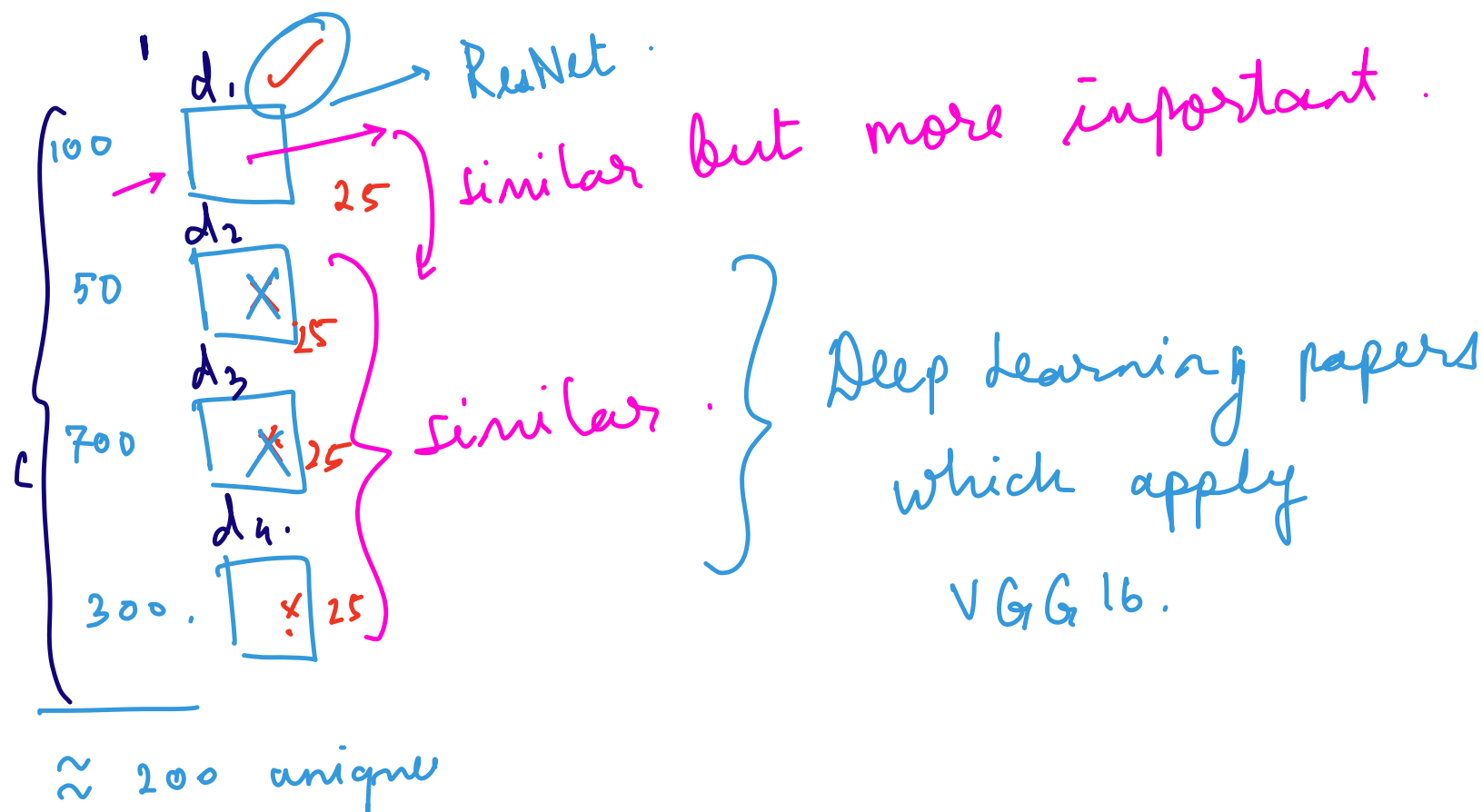
Document \rightarrow Deep learning paper ✓

neural — 300

network — 1000

\rightarrow residual — 50

\uparrow
Gives us some
important info.



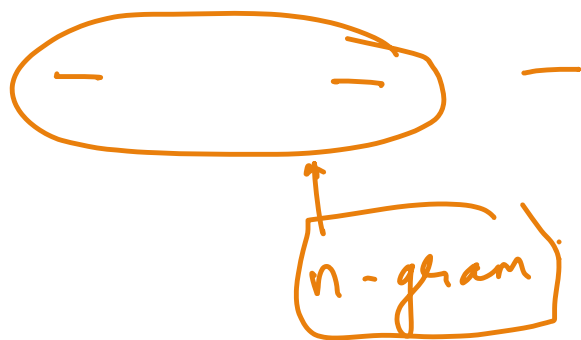
$$IDF(t, d, C) = \log\left(\frac{|C|}{|\{d \in C : t \in d\}|}\right)$$

$$IDF(\text{"residual"}, d_1, C) = \log\left(\frac{4}{1}\right) = \log(4)$$

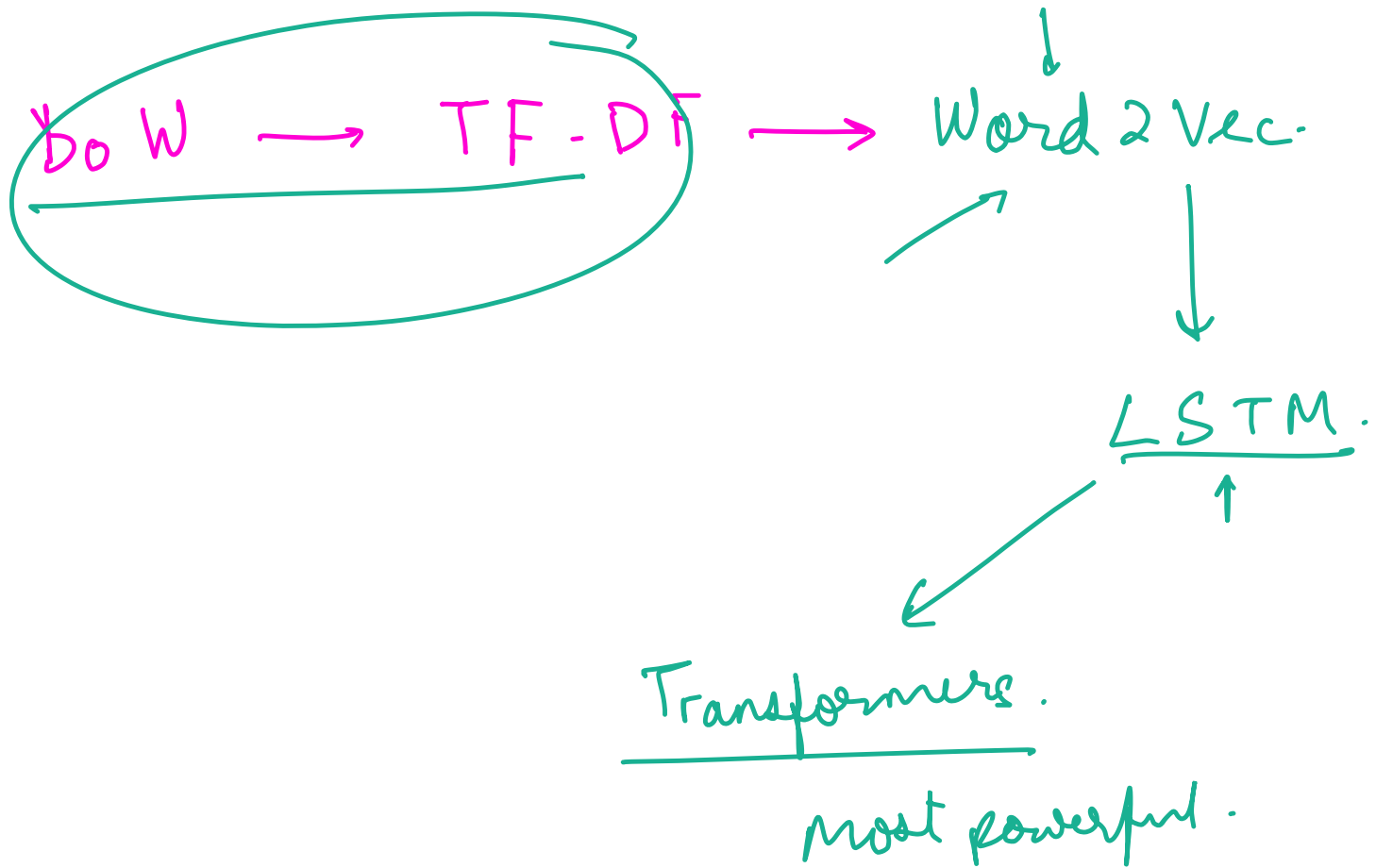
$$IDF(\text{"network"}, d_1, C) = \log\left(\frac{4}{4}\right) = \frac{\log(1)}{0.6}$$

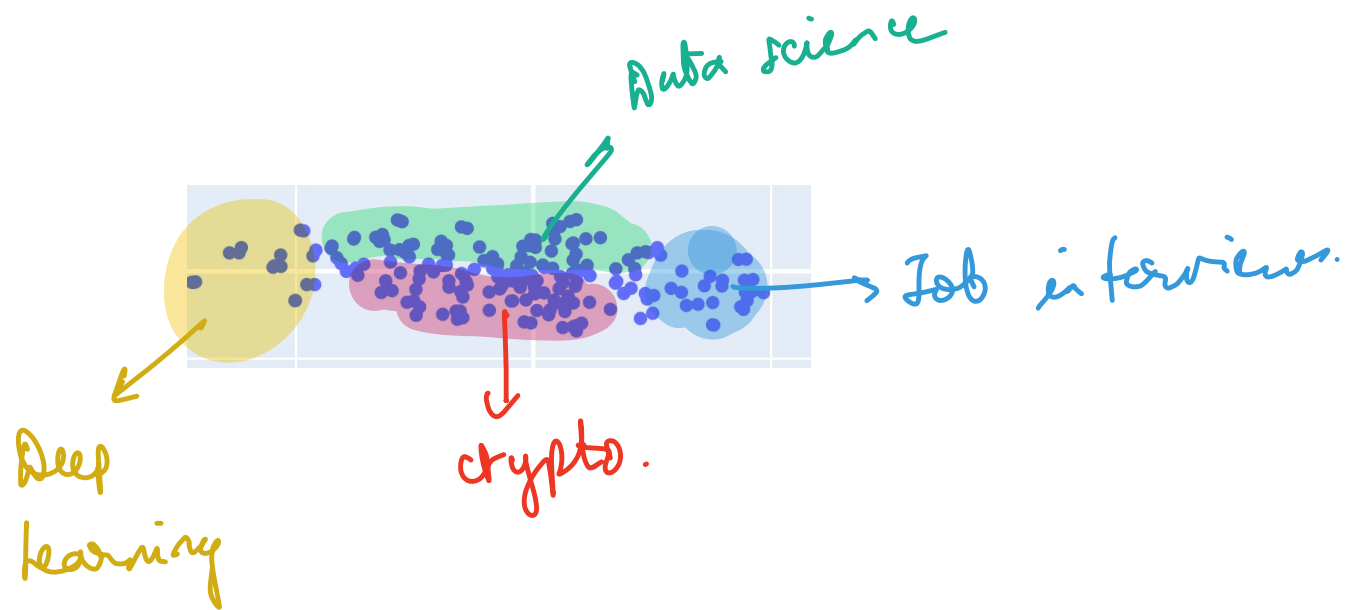
"The ^{new} iPad is not good compared to previous releases."

good x bad
not good ↗
n = 2



Instead of taking individual words, take pairs, triplet, quadruplet etc.

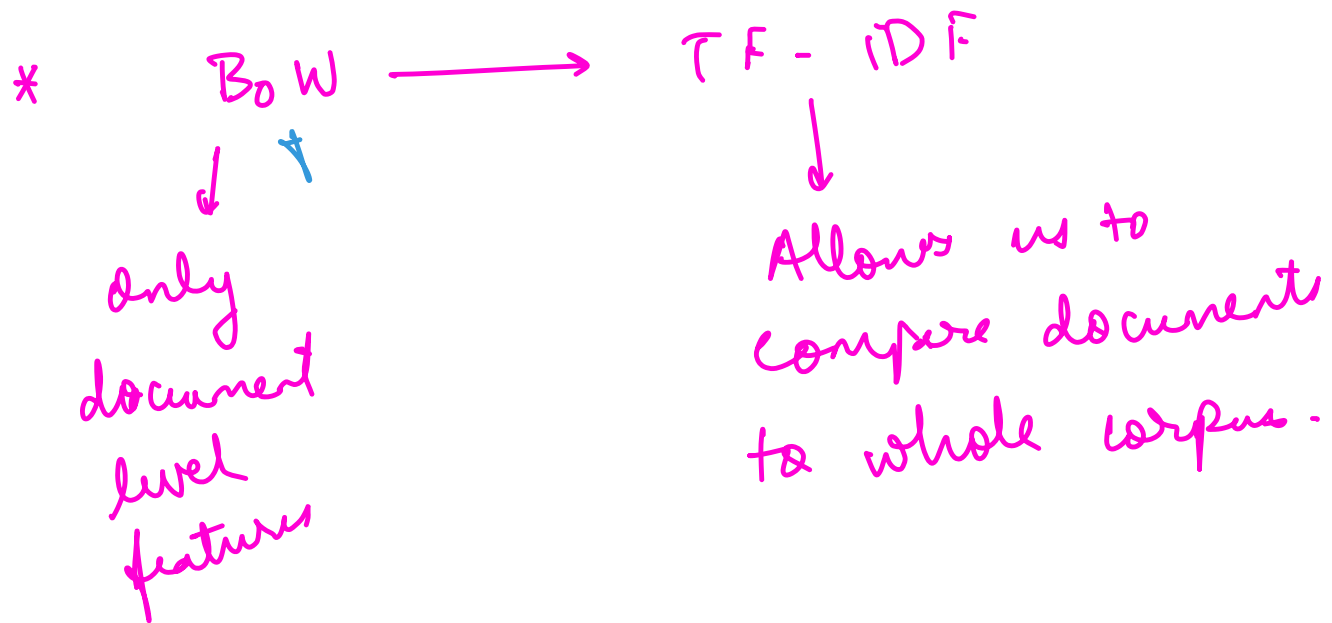




Concluding Remarks.

* How to convert text \rightarrow Vectors.

* We studied 2 methods - BoW, TF-IDF.



* Shortcomings:

- ① Content linkage missing.
- ② Size of the vector depends on unique words.

IDF \rightarrow $|C| \rightarrow$ # of docs in corpus.

2. Max IDF,

$$IDF(t, C) = \log\left(\frac{\max_{t' \in d} (|d \in C : t' \in d|)}{|d \in C : t \in d|}\right)$$

