

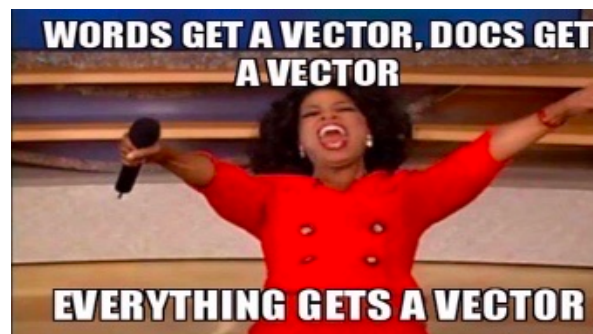
April 26, 2023.

DSML: NLP module.

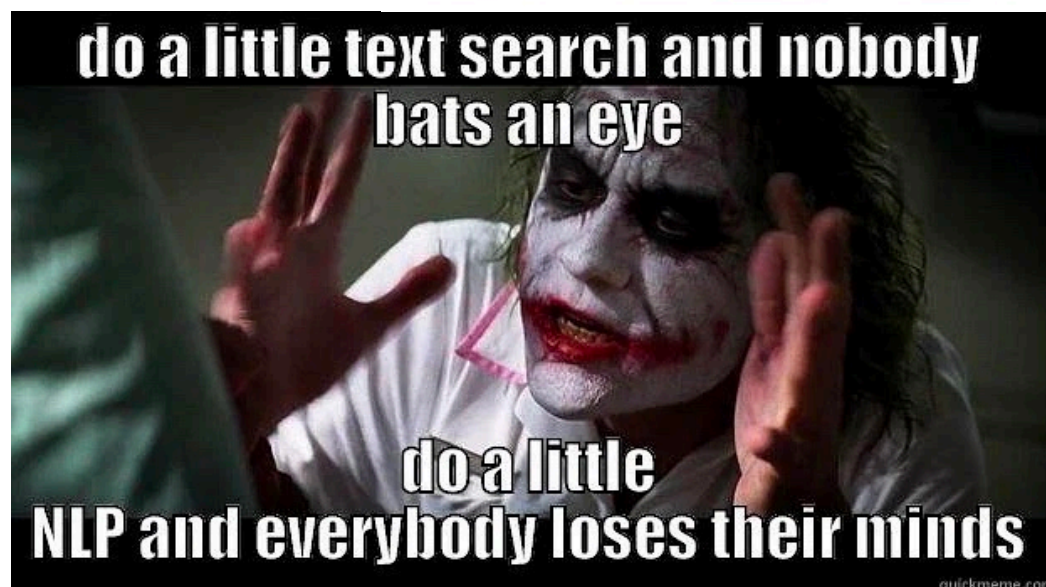
Word embeddings  
in a nutshell

Word Embedding: Word 2 Vec

Class starts  
@ 9:05



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



## Recap:

- \* Corpus: Collection of documents.
- \* Documents: A collection of sentences.
- Agenda (NLP): Find good vector representations for documents.
- Bag of Words: Document  $\rightarrow$  Vector.
- TF-IDF: Document  $\rightarrow$  Vector.

$$= A \cdot B$$

Agenda ↓

$n \times k$  \*  $k \times d$  Fundamental change: Words  $\rightarrow$  Vectors.

- \* Why? We want vector representations which capture semantics (meaning).
- \* Plan: Let's build/invent techniques together!!  
(We have all pre-requisites to figure this out)
- \* But first... The business case!!

Corpus.

$$\begin{bmatrix} d_1 - \\ d_2 - \\ \vdots \\ d_n - \end{bmatrix} \quad \uparrow$$

$d$  - dim vector.  
 $\uparrow$   
 the size of the vocabulary.

documents are paper abstracts.

Recap: Rec. Sys.

Movie Recos.

M. F

$$\underbrace{\begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}}_{n \times d} \xrightarrow{\text{pink arrow}} R = \begin{matrix} A & \cdot & B \\ \downarrow & & \downarrow \\ n \times \textcircled{k} & & \underbrace{k \times d.} \end{matrix}$$

B. o. W representation.

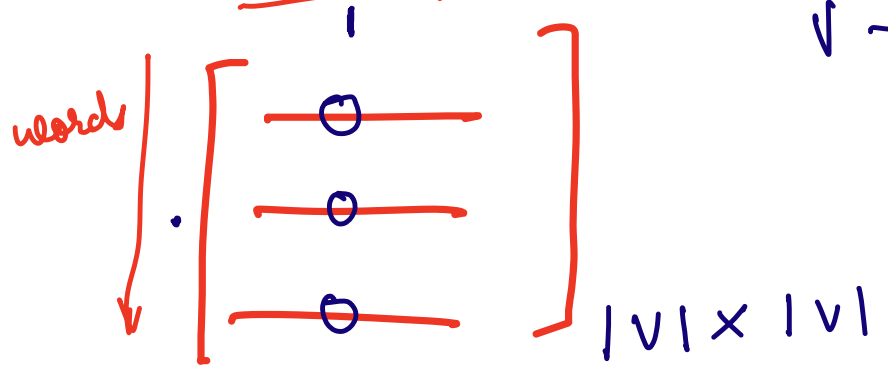
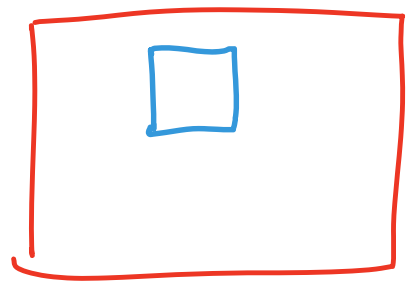
$$R = A \cdot B$$

$n \times k$        $k \times d$        $t-1$        $t$        $t+1$        $t+2$

How is object detection related to segmentation?

proximity.

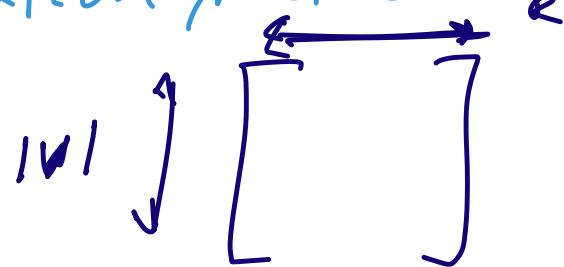
ngram  $\Rightarrow n = 2$ . words.



$V \rightarrow$  Vocabulary set.

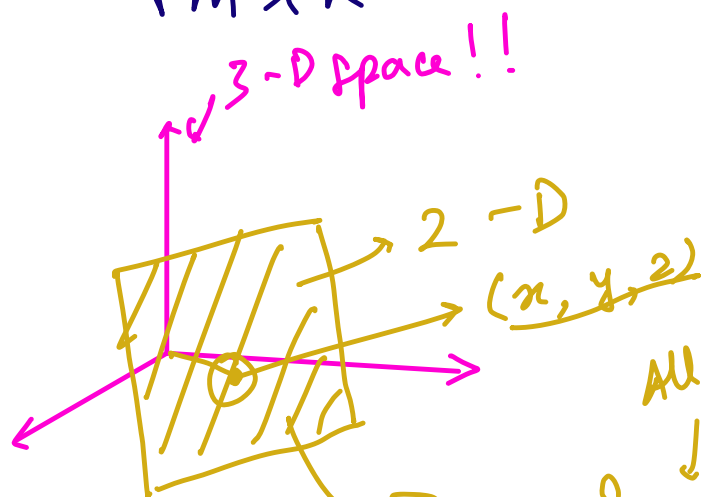
$$\boxed{R}_{|v| \times |v|} = \boxed{U}_{\uparrow} \Sigma \downarrow V.$$

Singular Value Decomposition method.  
SVD.



$$A_{m \times n} = U_{m \times k} \Sigma_{k \times k} V_{k \times n}$$

Inherent dimensionality  
of  $A$ .



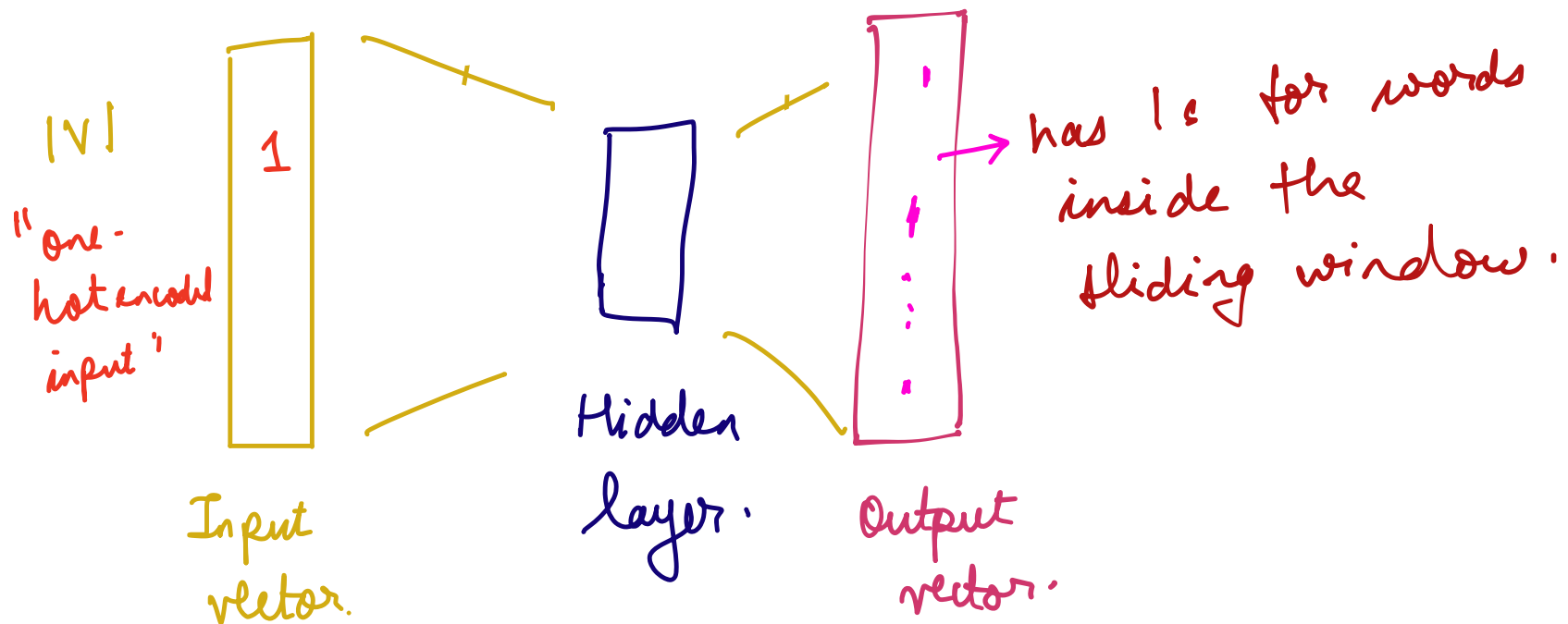
All vectors  
↓ lying on the plane.

$$P = \{ \cdot \}$$

# Neural networks solution.

↳ Autoencoder solution for P.R. !!

Word 2 Vec. → Name given to autoencoders which gives word-embeddings.



1) Fill - in - the blanks?

ResNet is a type of Convolutional Neural Network.

Continuous Bag of words strategy for Word 2 Vec.



2) Write an essay on ... skip-gram strategy for Word 2 Vec.

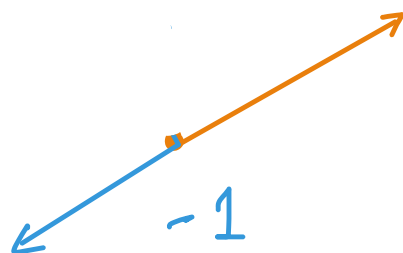
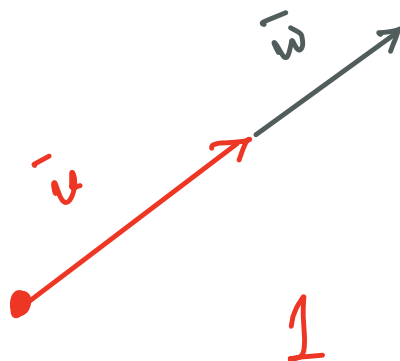
Convolutional.



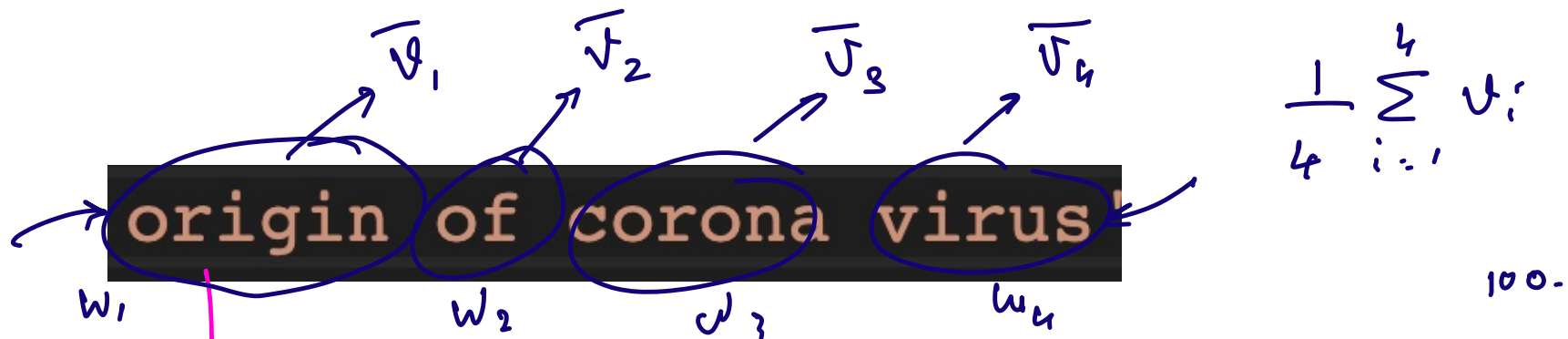
$$|V| \approx \underline{\underline{10^5}}$$

$$\text{softmax} = \frac{e^{z_i}}{\sum_{k=1}^{|V|} e^{z_k}} \quad \swarrow$$

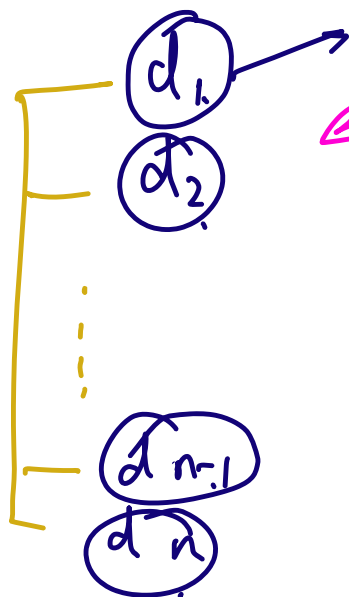
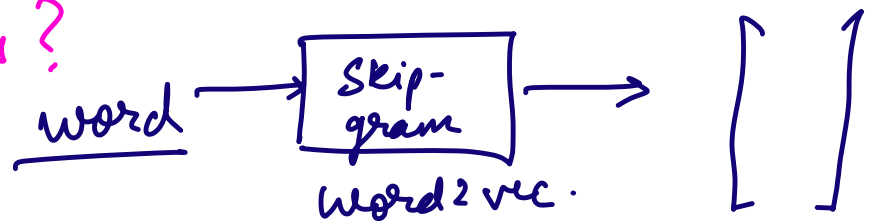
$\rightarrow O(|V|) \rightarrow$  This will take too long!!



$$\arg \max_{\theta} \sum_{(w,c) \in D} \log \sigma(v_c \cdot v_w) +$$



① How to encode this?



② How to encode these?