

- ↳ Text Cleaning & Text Preprocessing
- ↳ Text encoding
- ↳ Word Embedding =

Word2Vec

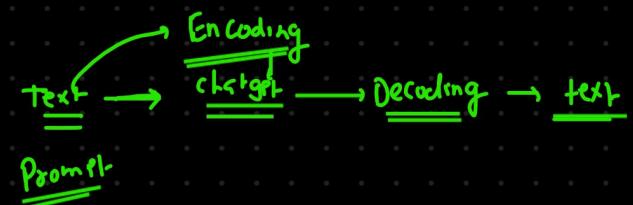
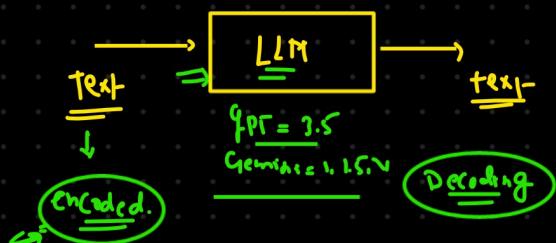
BERT embedding

Transformer

Attention all you need ⇒ Python



- ① Vector
- ② Vector Similarity
- ③ Embedding ⇒ Word2Vec
- ④ Python Practical



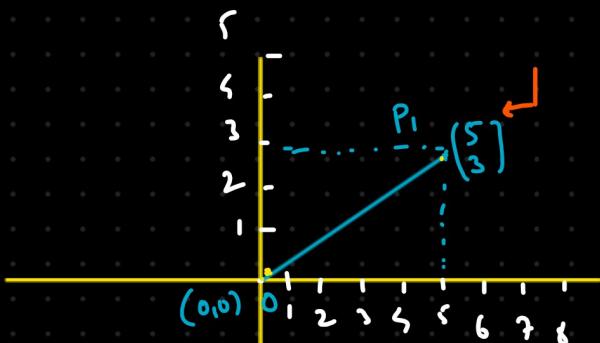
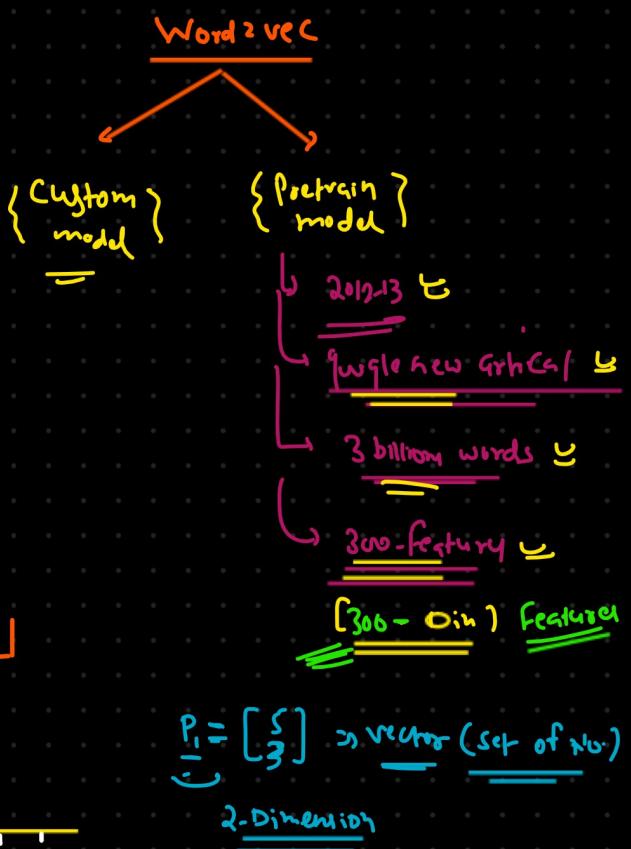
Word2vec

Neural network

CBOW
Skipgram

Word2vec \Rightarrow Word \rightarrow vector

- ① Semantic meaning
- ② Lower Dim
- ③ Sparsity \rightarrow Dense



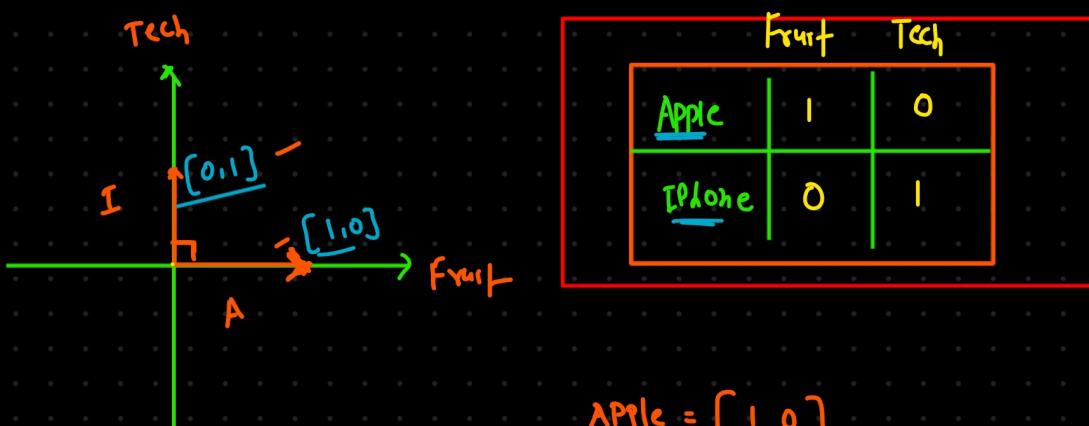
$$\begin{aligned} &= \sqrt{(5-0)^2 + (3-0)^2} \\ &= \sqrt{25 + 9} \\ &= \sqrt{34} \end{aligned}$$

$$\overrightarrow{OP} = \text{Not } \mathbf{0}$$

$$= \sqrt{34} \begin{bmatrix} 5 \\ 3 \end{bmatrix}$$

$$\boxed{\overrightarrow{OP} = x\hat{i} + y\hat{j} + z\hat{k}}$$

Vector \Rightarrow Different



$$\text{Apple} = \begin{bmatrix} 1, 0 \\ \hline \end{bmatrix}$$

$$\text{iPhone} = \begin{bmatrix} 0, 1 \\ \hline \end{bmatrix}$$

① Dot product ③ euclidean Dist }
 ② Cosine Similarity \equiv } Similarity RAG \Rightarrow Similarity Search

$$① \quad \underline{\underline{[1, 0]}}, \underline{\underline{[0, 1]}} = \frac{1 \times 0 + 0 \times 1}{\sqrt{1^2 + 0^2 + 0^2 + 1^2}} = 0$$

$$② \quad \underline{\underline{[1, 0]}}, \underline{\underline{[0, 1]}} = \cos 90^\circ = 0$$

$$③ \quad \sqrt{(1-0)^2 + (0-1)^2} = \sqrt{2}$$

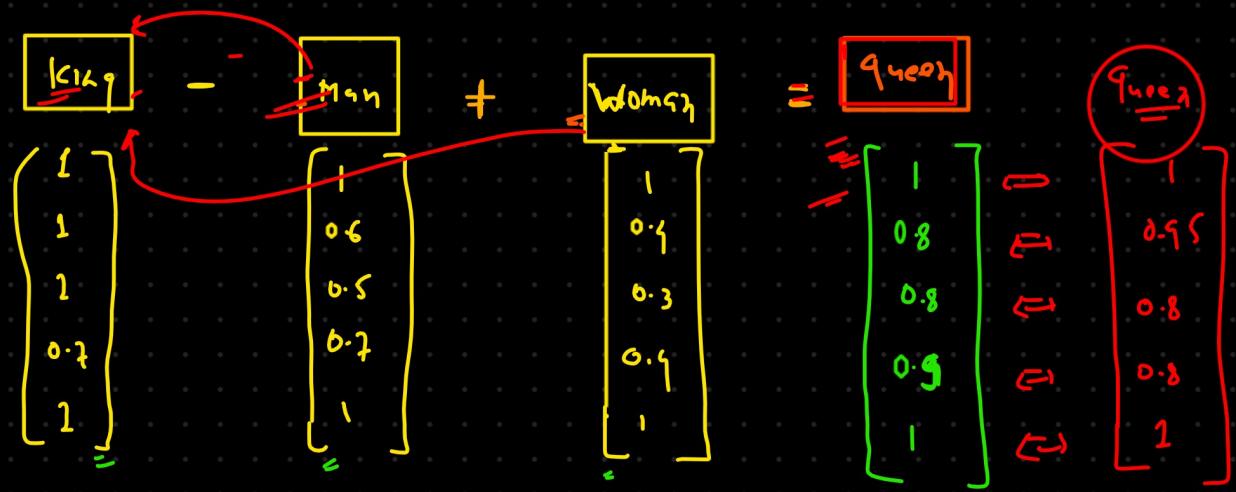
Word2vec

Semantics

	\bar{F}_1 <u>Gender</u>	\bar{F}_2 <u>Wealth</u>	\bar{F}_3 <u>Power</u>	\bar{F}_4 <u>Weight</u>	\bar{F}_5 <u>Specie</u>	
King	1	1	1	0.7	1	(25) $\underline{\underline{(0,1)}}$ (myself)
Queen	1	0.95	0.8	0.8	1	
Man	1	0.6	0.5	0.7	1	
Woman	1	0.4	0.3	0.9	1	
Monkey	1	0	0	0.3	0	

Word

Vector



Dimension \Rightarrow 5D

Athletic operation

mathematically



\downarrow
 $\text{PdF} = 1.0 \text{ feature}$

5 → feature \Rightarrow manually Digt \rightarrow feature

Automated.

log B hours \Rightarrow feature \Rightarrow NH

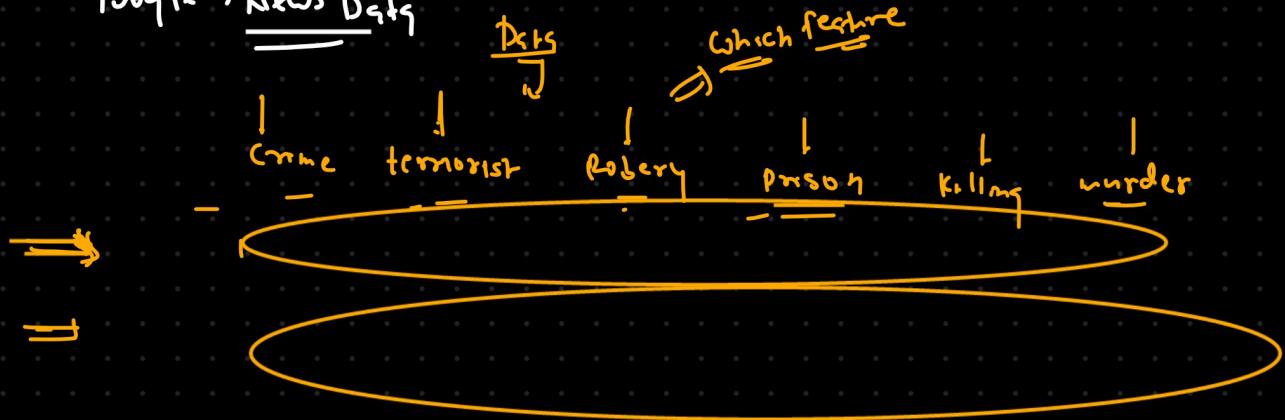
NewsDigest

Representation

As worries about how social media sites might make people addicted grow, the European Union has recently started an investigation into Meta Platforms Inc. and their goods, like Facebook and Instagram. This move is meant to deal with and control

$$\boxed{\text{objec-}} \quad \frac{-F_1}{\underline{s_m}} + \frac{-F_2}{\underline{\text{Add}_{cv}}} - \frac{-F_3}{\underline{g_{now}}} : F_L \quad \underline{F_S} - \underline{F_C}$$

Google \rightarrow News Data



Skipgram, CBOW \Rightarrow NN

CBOW \rightarrow Continuous Bag of words



X
[input text length]

Watch inuron

neuron data

for Science

Y
O/P

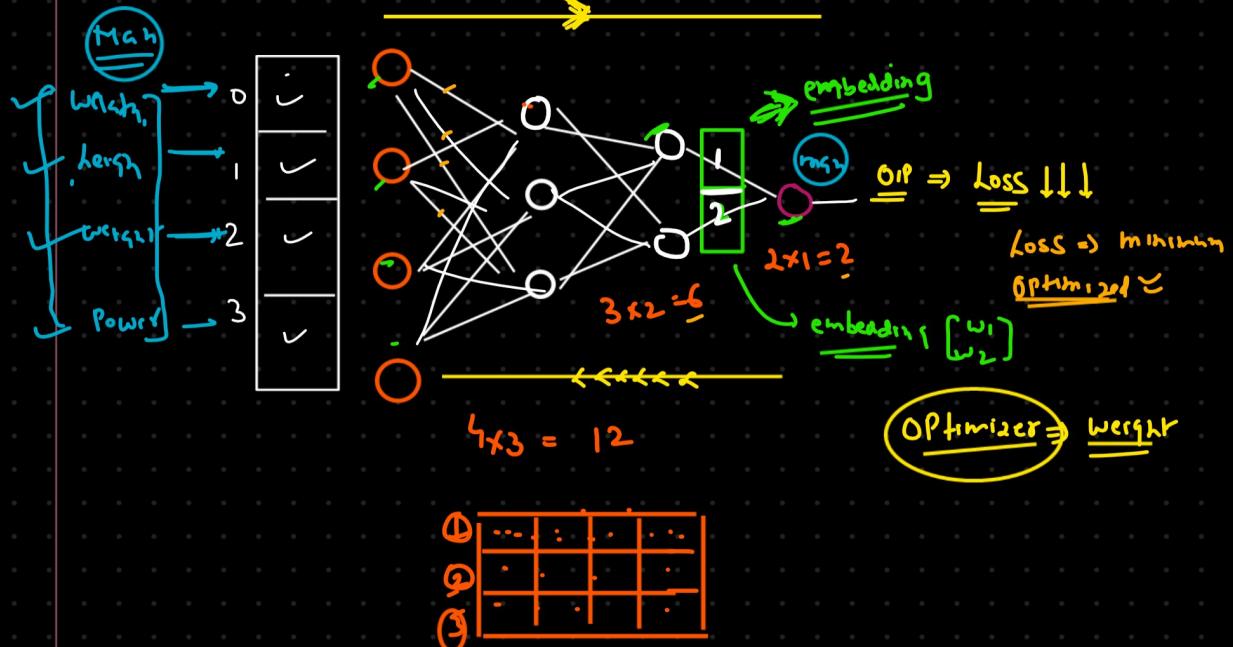
ineuron

For

Data

IIP

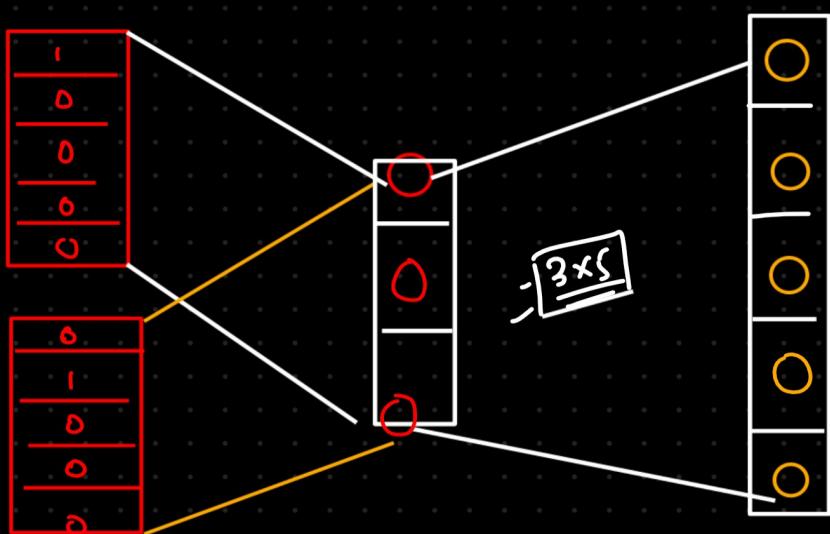
FP \Rightarrow calculation



SIP

Hidden

OIP



Watch **neuron** for Data Science

Context=3

Model Context=3

Howe Ph



- Data | For Scene

VOC \Rightarrow S



5×3



3×5



3×5

