

Word Embedding

- * Word Embeddings are type of word representation that allows words with similar meaning ~~for the~~ to have a similar representation.
- * Word embeddings where individual words are represented as real valued vectors in a predefined vector space.
- * Each word \rightarrow mapped to \rightarrow vector values.
↓
tens of hundreds of dimensions.
- * Contrasted to thousands or million dimensions for sparse word representation.

So Each word is associated with a point in vector space. The no of features is less smaller than the size of the vocabulary.

- \rightarrow Word embedding that is learned jointly with a neural network model on a specific NLP task.
- \rightarrow Requires text to be cleaned and prepared such that each word is one hot encoded.
 - \rightarrow embedded layer is used on the front end of the neural network and is fit in a supervised way.

① Word2vec is a statistical method for efficiently learning a standalone word embedding for text corpus
by Tomas Mikolov [↓] Google 2013

* This representation surprisingly good at capturing syntactic and semantic regularities in language.

king - man + woman = Queen

→ CBOW Continuous Bag of word model

→ Continuous Skip Gram model

CBOW model learns the embedding by predicting the current word based on context.

Continuous Skip Gram model learns by predicting the surrounded words given a current word

* context is defined by a window of neighbouring words

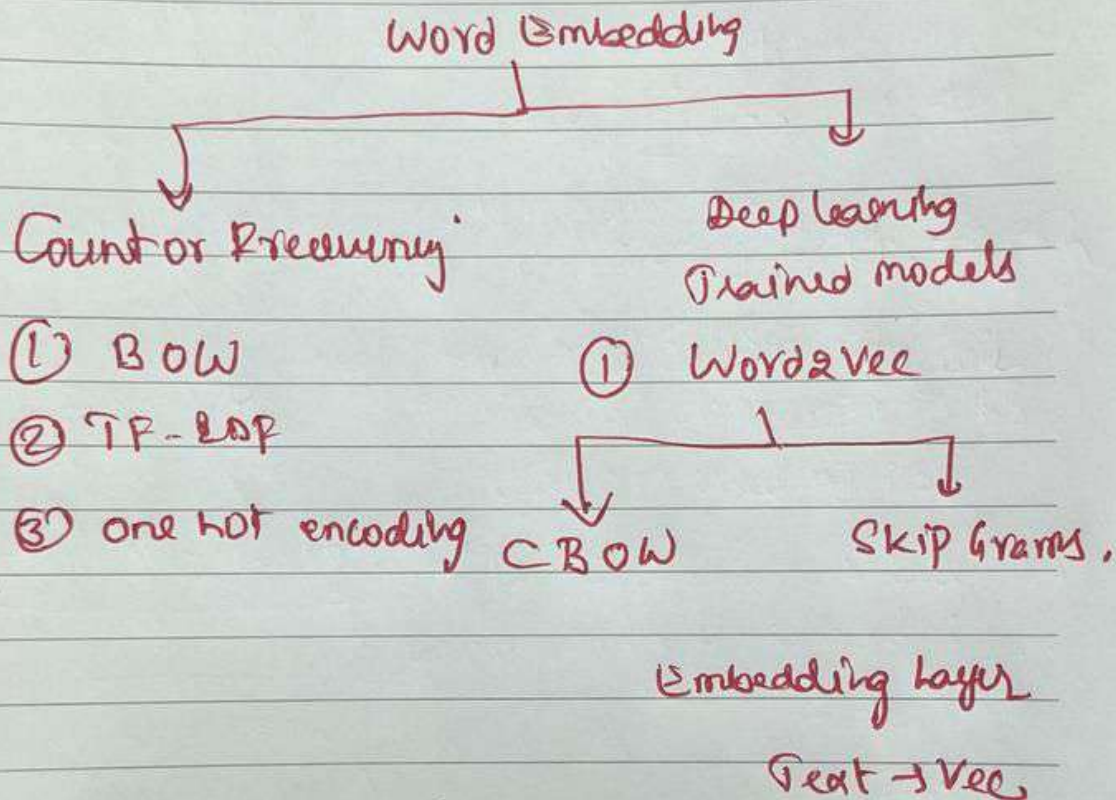
* window is a parameter

② ~~Colover~~

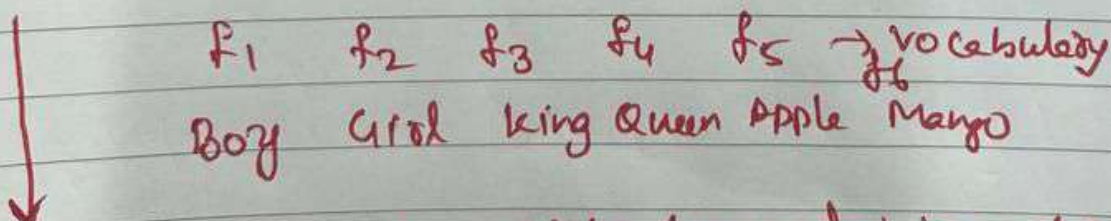
② Glove : Global vector for word representation is an extension of word2vec for efficiently learning word vectors by Pennington & Stanford

* Word embedding :- *

It is a technique which converts word into vectors



Word2Vec :- Feature Representation



- * every word will be converted to vector
- * limited dimensions

- * sparsity will be reduced (not find many zeros)
- * semantic meaning will be maintained

drawbacks of BOW, TF-IDF

(1) { honest }
good } Symmetric relationship

Date _____

(2) Sparse Matrix
→ huge dimensions

	f1	f2	f3	f4	f5	f6
	Boy	Girl	King	Queen	Apple	Mango
Gender	-1	1	-0.92	+0.93	0	0
Royal	0.01		0.95	0.96		
Age	0.03					
Food						
↓						
300 Dimen						

→ trained model

→ semantic relationship

→ related vector

king $[0.96, 0.95]$

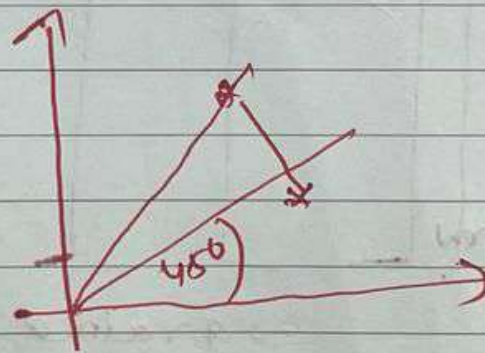
Queen $[-0.96, 0.95]$

Man $[0.95, 0.98]$

Women $[-0.94, -0.96]$

king - Man + Woman = Queen

Cosine Similarity



Distance = $1 - \text{Cosine Similarity}$

$$\cos 45^\circ = 0.7071 \approx \frac{1}{\sqrt{2}}$$

$$\Rightarrow 1 - 0.7071$$

$$\Rightarrow \approx 0.29$$

Ph : +91-40 - 48488275

more towards
very similar

* Word2vec

① CBOW = { Continuous Bag of words }

Corpus : dataset Training

[Krish channel is related to] ^{total 7 words}
data science

Window size = 5

* window must be odd no
 Bigger window
 → 3000 model

Independent feature

O/P

→ Krish, channel, Related, to 25

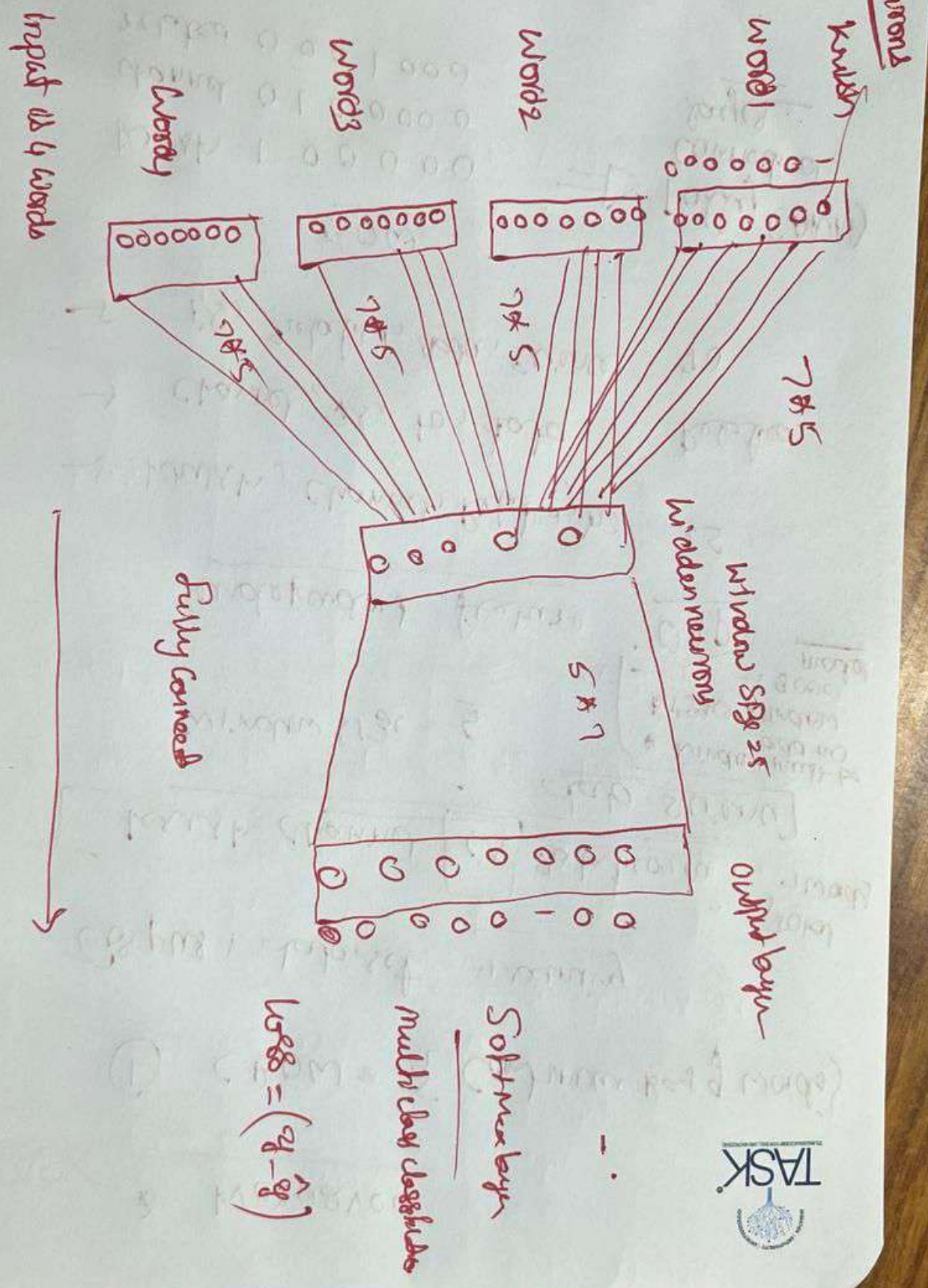
→ channel, is, to, data Related

→ is, related, data, science 30

BOW

Krish	1	0	0	0	0	0
channel	0	1	0	0	0	0
related	0	0	0	1	0	0

⇒ Fully (ANN)
 Connected
layer



Date _____

② Skip gram

Q1P

Q1P

is
Related

crush, channel, related, to
channel, us, to, data

QO

Qs, related, data, Science.