

## Chapter 2

Date .../.../...

### # Vector Space Models for NLP

#### ⇒ 2-D Vector Space

A 2-D vector-space is defined as a set of linearly independent basis vectors with 2 axes. Each axis corresponds to a dimension in the vector space.

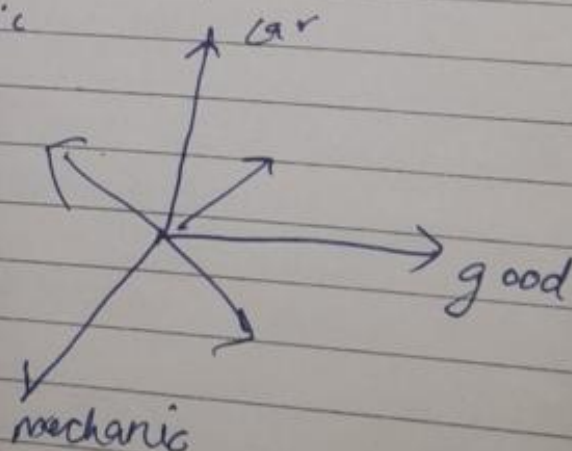
#### ⇒ 3-D Vector Space

A 3-D vector-space is defined as a set of linearly independent basis vectors with 3 axes. Each axis corresponds to a dimension in the vector space.

Linearly independent vectors of size  $N$  will result in  $N$  dimensional axes which are mutually orthogonal to each other.

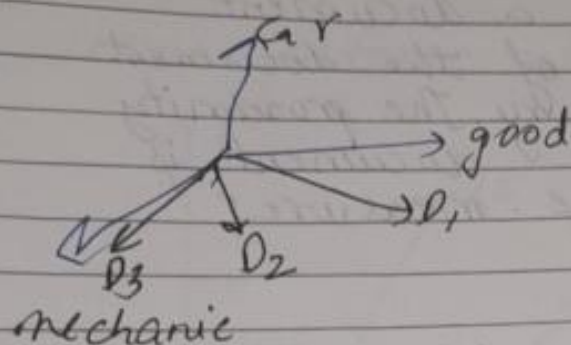
Ex:- Binary Incidence Matrix

	good	car	mechanic
$D_1$	1	1	1
$D_2$	1	0	1
$D_3$	0	1	1



Ex:- TF-IDF Incidence Matrix

	good	car	mechanic
$d_1$	0.91	0	0.0011
$d_2$	0.21	0	0.1
$d_3$	0.15	0	0.921



# Document - Term Matrix

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	...	$d_{12}$
$t_1$	0.1	0.0	-	-	-	-	-
$t_2$	0.1	0.0	-	-	-	-	-
$t_3$	0.0	0.9	-	-	-	-	-
$t_4$	0.0	0.8	-	-	-	-	-
$\vdots$							
$t_{10}$	0.6	0.7	-	-	-	-	-

The columns of the matrix represent the documents as vectors. A document vector is represented by the term present in the document.



## # Query Modelling

Date .... / .... / .....

Each query is modelled as a vector using the same attribute space of documents.

$$q = [q_1, q_2, q_3, \dots, q_n]$$

The relevancy ranking of a document depends on the distance of the document with respect to the query. The proximity of the query with every document is computed using distance measures.

## # Document Similarity

Earlier using the Binary Incidence Matrix, a query returned a set of documents whether the query keywords were found in documents or absent. It did not give any ranking for the retrieved documents.

A similarity measure is a real-value function that qualifies the similarity between two objects. Some of the methods are given below.

Euclidian distance =  $\sqrt{d_1^2 - d_2^2}$

Cosine similarity =  $\frac{d_1 \cdot d_2}{\|d_1\| \|d_2\|} = \frac{\vec{d}_1 \cdot \vec{d}_2}{\|\vec{d}_1\| \|\vec{d}_2\|}$

cosine distance =  $1 - \cos(\vec{d}_1, \vec{d}_2)$   
cluster similarity =  $\frac{\vec{d}_1 \cdot \vec{d}_2}{\|\vec{d}_1\| \|\vec{d}_2\|}$

Jaccard similarity =  $\frac{|\vec{d}_1 \cap \vec{d}_2|}{|\vec{d}_1 \cup \vec{d}_2|}$

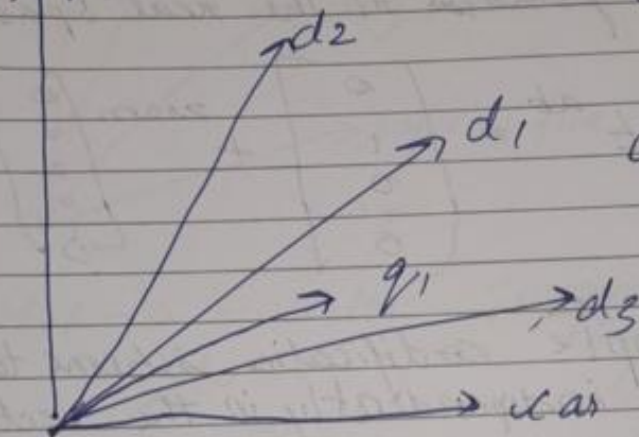
Spiral



Euclidean measure does not work well for unequal sized vectors as the vectors are not normalized. We often use a normalized correlation coefficient, cosine distance for the similarity measure.

$$\text{cosine distance} = 1 - \cos(\vec{d}_1 \cdot \vec{d}_2)$$

plane



$$\cos \alpha_1 = \frac{\vec{q}_1 \cdot \vec{d}_3}{|\vec{q}_1| |\vec{d}_3|}$$

$$\cos \alpha_2 = \frac{\vec{q}_1 \cdot \vec{d}_2}{|\vec{q}_1| |\vec{d}_2|}$$

$$\cos \alpha_3 = \frac{\vec{q}_1 \cdot \vec{d}_1}{|\vec{q}_1| |\vec{d}_1|}$$

### # Proximity Score

A query is considered as a document vector. The proximity of the query with every document is computed using a distance measure. Cosine distance is preferred and it is easy to compute if the document vector distances are normalized. Proximity score (angle) will be considered as relevant and retrieved.

# Vector Representation of words  
 Let  $V$  be the unique terms and  $|V|$  be the size of the vocabulary. Then every vector representing the word  $R^{|V| \times 1}$  would point to a vector in the  $V$  dimensional space.

### # One hot Vector

We can represent each word as an independent vector quantity as follows in the real space  $R^{|V| \times 1}$

$$t^a = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \quad t^{ab} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} + \text{zoom} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

This is a very simple codification scheme to represent words independently in the vector space. This is known as one hot vector.

### # One Hot Vector - 2

In one-hot vector, every word is represented independently. The terms, home, house, flats, apartments are independently coded. With one-hot vector based model, the dot product

$$(t^{\text{house}})^T \cdot t^{\text{apartment}} = 0$$

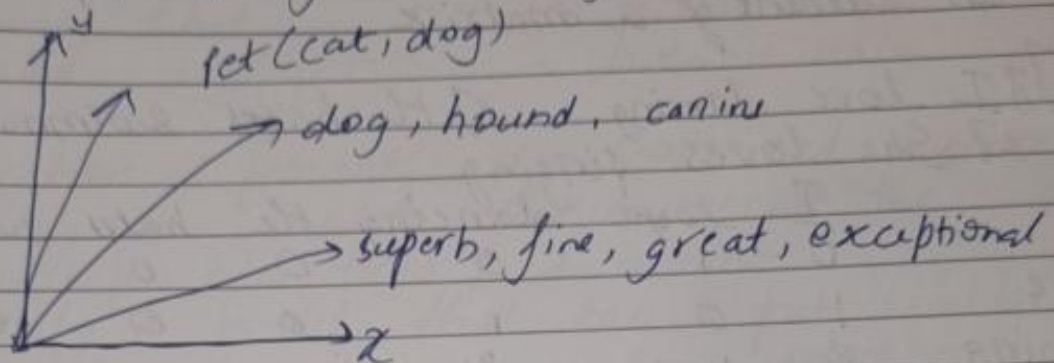
With one-hot vector, there is no notion of similarity or ~~synonyms~~ synonyms.



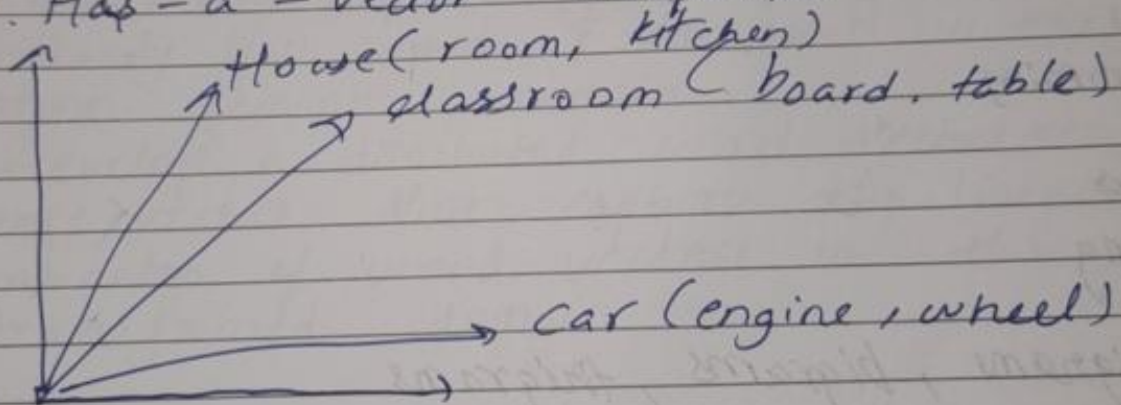
Goal of word vector:-

- Reduce word - vector space into a smaller sub-space
- encode the relationship among words.

→ Relationship among terms - synonyms



→ Has-a - vector - composition



# Contextual Understanding of text

- In order to understand the word and its meaning, it's not enough if we consider only the individual word
- The meaning of context should be central in understanding word/text
- exploit the context dependency of words
- co-occurrence of words

## # Co-occurrence Matrix

A co-occurrence is a combination of terms that are likely to be used in the same context. A co-occurrence matrix stores co-occurrences of words. The count of a pair of words that appears in a context window is represented as an element of a matrix.

[Ex] 1] I love dancing 2] He hates swimming  
3] She loves singing

	I	love	dancing	He	hates	swimming
I	0	1	0	0	0	0
love	1	0	1	0	0	0
dancing	0	1	0	0	0	0
He				1	0	0
hates				0	1	0
swimming				0	0	1
she						
loves						
singing						

## # Unigrams, bigrams, trigrams

- A sequence of two words is called a bigram
- A three-word sequence is called trigram
- N-gram means a sequence of words of length n.



## # Collocations

Collocations is a juxtaposition of two or more words that more often occur together than by chance.

- powerful computer
- brief chat
- broad daylight
- major problem
- pitch dark
- heavy rain

## # Creation of semantically connected vectors

- Identify a model that enumerates the relationship b/w terms & documents
- Identify a model that tries to put similar items closer to each other in some space or structure
- A model that discovers/uncovers the semantic similarity b/w words & documents in the latent semantic domain
- Develop a distributed word vectors or dense vectors that capture the linear combination of word vectors in the transformed domain.

## # Methods to create dense vector

- Latent Semantic Analysis or Latent Semantic Indexing
- Neural networks using skip grams & CBOW
- Skip grams use center of words to predict the surrounding words.
- Brown Clustering - statistical algorithms for assigning words to class based on the frequency of their co-occurrence with other words.



## # Why dense Vectors?

Date .../.../.....

- sparse vectors are too long & ↓ convenient as features in machine learning.
- abstracts more than just frequency count.
- It captures neighborhood words that are connected as synonyms.

## # Singular Value Decomposition

Singular value decomposition is a method to factorize a rectangular / square matrix into three matrices.

$$\begin{array}{ccccc} A & = & U & \Sigma & V^T \\ M \times N & & M \times K & K \times K & K \times N \\ & & \text{left singular} & \text{singular} & \text{right singular} \\ & & \text{vectors} & \text{vector} & \text{matrix} \end{array}$$

Row vectors of  $U$  are called as the left-singular vectors. Rows of  $V^T$  form an orthonormal set.  $\Sigma$  is diagonal matrix and its values are arranged in descending order.

## # Singular Values

- It is a diagonal matrix.
- Singular values are arranged in descending order.
- Singular values reflect the major associative patterns in the data, & ignore the smaller less important influences.
- Highest order dimension captures the most variance in the original dataset or document matrix.



- The next higher dimension captures the next higher variance in the original data set
- arranged in descending order

### # Dimensionality Reduction

- SVD is better suited for measuring the similarity b/w documents by exploiting the similarity patterns that exist in the word co-occurrence.
- The co-occurring terms are mapped in the same dimension thereby reducing the dimensions
- Increases the similarity Matrix  $A$  in the  $m$ -dimensional space & transforms it as  $\hat{A}$  in the reduced dimensional space  
 $k \leq m$

$$\Delta = \|A - \hat{A}\|_2 \text{ should be as less as possible}$$

where  $\|\cdot\|$  is the  $L_2$  norm for the matrices.

Similar documents will be brought to <sup>one</sup> hyper plane.

Similar words will be brought to same axis.

### # Important equation in SVD

Since  $U$  &  $V$  are orthonormal matrices

$$U^T U = V^T V = I$$

$$A^T A = V \Sigma^T U^T U \Sigma V^T = V \Sigma^2 V^T$$

$$A A^T = U \Sigma V^T V \Sigma^T U^T = U \Sigma^2 U^T$$

$$U^T A = U^T U \Sigma V^T = \Sigma V^T$$

$$U \hat{q} = U^T U \Sigma^{-1}$$



$$\Sigma_3^{-1} = \begin{bmatrix} \frac{1}{\sigma_{11}} & 0 & 0 \\ 0 & \frac{1}{\sigma_{22}} & 0 \\ 0 & 0 & \frac{1}{\sigma_{33}} \end{bmatrix}$$

Date .... / .... / .....

$$q = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix} \rightarrow u_q = u_q^T U \Sigma^{-1}$$

$$\cos \theta = \frac{u_q \cdot d}{\|u_q\| \|d\|}$$