

22-12-2020

## Chapter - 5

### Topic modelling

#### Objectives :

- 1) Describe topic modelling and its use cases
- Describe topic modelling algorithms
- Describe the working of LSA and LDA
- Describe topic fingerprinting
- Implement topic modelling using LSA and LDA.

#### Topic modelling :

- A simple way to capture meaning from a collection of documents
- Unsupervised learning algo.

⇒ Why topic modelling :

- Topic discovery :

Find a set of topics that can be used to classify a set of documents

- Discovering themes:-

To get a general idea of what themes or topics are present in document.

- EDA:-

Topic modelling can be used to check whether data is balanced or ~~not~~ skewed, which then can be used to select the ML algo.

- Document clustering:-

Topic modelling allows soft clustering

- Dimensionality reduction

- Historical analysis.

## ⇒ Topic modelling Algorithm +

### Assumptions +

- Topics contain a set of words.
- Documents contain a set of topics.

Analog → LSA (Latent Semantic Analysis)

LDA (Latent Dirichlet Allocation)



→ Better Definition  
(from Analytics Vidhya video)

Topic → A repeated group of statistically significant tokens or words in a corpus.

Topic modelling - Process to find topics from documents in an unsupervised manner.

(Out of Book content)

⇒ Overview of SVD (Singular value decomposition)  
(Since it will be used in LSA ~~SVD~~)

Used for:

- Data reduction
- Data driven generalization of Fourier transform

Let  $X$  be a matrix.

$$X = \begin{bmatrix} | & | & | & \dots & | \\ x_1 & x_2 & x_3 & \dots & x_m \\ | & | & | & \dots & | \end{bmatrix} \text{ where } x_i \in \mathbb{R}^n$$

shape of  $x \rightarrow (n, m)$

$\downarrow$   $(n, 1)$   
row

Using SVD we can represent this  $X$  as a product of 3 other matrices.

$$X = U \cdot \Sigma \cdot V^T$$

$U$ : unitary, orthogonal matrix  
 $\Sigma$ : diagonal matrix  
 $V^T$ : orthogonal matrix

$$U U^T = U^T U = I_{n \times n}$$

$$V V^T = V^T V = I_{m \times m}$$

$\Sigma$  diagonal  $\sigma_1, \sigma_2, \dots, \sigma_m \geq 0$ .

$$\begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ x_1 & x_2 & \dots & x_m \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}_{(n,m)} = \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ U_1 & U_2 & \dots & U_n \\ \downarrow & \downarrow & & \downarrow \end{bmatrix}_{(n,n)} \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_m & \\ & & & 0 \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow & \uparrow \\ \sqrt{v_1} & \sqrt{v_2} & \dots & \sqrt{v_m} \\ \downarrow & \downarrow & \downarrow & \downarrow \end{bmatrix}_{(m,m)}$$

ordered by importance

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

$U$ : eigen vectors  
 $\Sigma$ : importance  
 $V^T$ : mixture of v's that makes  $X$ 's

$U, \Sigma, V^T \rightarrow$  for  $X$ , they are guaranteed to exist and are unique.

⇒ SVD matrix approximation

$$\underline{X} = \begin{bmatrix} \uparrow & \uparrow & \uparrow \\ u_1 & u_2 & \dots & u_m \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} = U \Sigma V^T = X$$

$$= \begin{bmatrix} \uparrow & \uparrow & \uparrow \\ u_1 & u_2 & \dots & u_n \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} \begin{bmatrix} \sigma_1 & \sigma_2 & & \\ & \sigma_3 & & \\ & & \ddots & \\ & & & \sigma_m \\ & & & & 0 \end{bmatrix} \begin{bmatrix} \leftarrow V_1^T \leftarrow \\ \leftarrow V_2^T \leftarrow \\ \leftarrow V_3^T \leftarrow \\ \leftarrow V_m^T \leftarrow \end{bmatrix}$$

$$= \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \dots + \sigma_m u_m v_m^T + 0$$

• Even  $V$  is a large matrix ( $n, n$ ) only the first  $m$  cols are significant for  $u$ . ( $n \gg m$ )

• After the first  $m$  cols, they are scaled by 0, (from diagonal matrix)  $\sigma_i = 0 \forall i > m$ .

∴ We need to use only the first  $m$  rows cols of  $U$  to get  $X$ .

$$X = \hat{U} \hat{\Sigma} V^T$$

↓  
economy  
SVD

We calculate only first  $m$  rows cols of  $u$ .



$$= \sigma_1 U_1 V_1^T + \sigma_2 U_2 V_2^T + \dots + \sigma_m U_m V_m^T$$

Representing them as a sum of rank 1 matrices.

$$= \sigma_1 \underbrace{\begin{matrix} U_1 \\ (n,1) \end{matrix}}_{(n,m)} \begin{matrix} V_1^T \\ (1,m) \end{matrix} + \sigma_2 \begin{matrix} U_2 \\ (n,1) \end{matrix} \begin{matrix} V_2^T \\ (1,m) \end{matrix} + \dots + \sigma_m \begin{matrix} U_m \\ (n,1) \end{matrix} \begin{matrix} V_m^T \\ (1,m) \end{matrix}$$

The above expression contains  $m$  terms.

But we will truncate it to  $n$  terms which will give us the best estimate of  $X$ .

$$= \sigma_1 \begin{matrix} U_1 \\ (n,1) \end{matrix} \begin{matrix} V_1^T \\ (1,m) \end{matrix} + \sigma_2 \begin{matrix} U_2 \\ (n,1) \end{matrix} \begin{matrix} V_2^T \\ (1,m) \end{matrix} + \sigma_3 \begin{matrix} U_3 \\ (n,1) \end{matrix} \begin{matrix} V_3^T \\ (1,m) \end{matrix} + \dots + \sigma_n \begin{matrix} U_n \\ (n,1) \end{matrix} \begin{matrix} V_n^T \\ (1,m) \end{matrix}$$

truncate it to  $n$  terms.

Now,

$$X \approx \tilde{U} \tilde{\Sigma} \tilde{V}^T = \tilde{X}$$

How to select  $n$ ,

According to Eckart - Young Theorem  
[1936],

$$\underset{\substack{\text{argmin} \\ \tilde{X}, \text{ s.t.} \\ (\tilde{X}) = n}}{\|X - \tilde{X}\|_F} = \tilde{U} \tilde{Z} \tilde{V}$$

After truncating,  
 $U$  is not unitary.

⇒ Interpretation of SVD in terms  
of correlation

[Continuing the book]

⇒ LSA [Latent Semantic Analysis]

• Our corpus is represented as a  
term-to-document matrix.

Rows - terms (words)

columns - documents.

	doc 1	doc 2	doc 3
we			
a			
an			



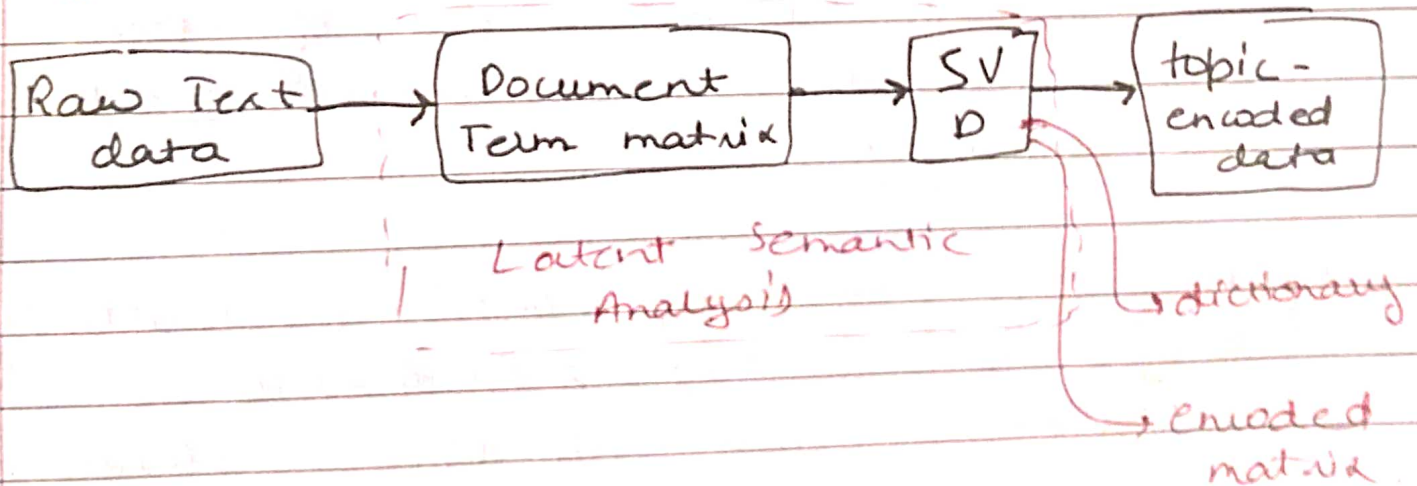
Using SVD, this matrix can be broken down into product of 3 matrices.

$$\begin{array}{ccccc} \text{Term-to} & = & \text{Term-to} & \times & \text{topic} & \times & \text{topic-to} \\ \text{document} & & \text{topics} & & \text{importance} & & \text{documents} \\ m \times m & & m \times n & & n \times n & & n \times m \end{array}$$

⇒ From Youtube

### Latent features

- Features that are hidden in the data which cannot be directly measured.
- These features are essential to the data, but are not original features of the dataset.



18<sup>th</sup>-12-2020

- In LSA, we have to choose topic number of topics beforehand.

To select the optimal number of topics, we create LSA model for different topics in range and check their coherence values.

## ⇒ Latent Dirichlet Allocation

- LDA is more often used for topic modelling.
- LDA is a generative statistical model that allows a set of items to be sorted into unobserved groups by similarity.
- LDA is reasonably good particularly useful for finding reasonably accurate mixture of topics within a given document.

⇒ From Youtube:  
Luis Serrano

The problem:

- Let's suppose that we have some documents, each document can belong to a topic or a combination of other topics.

LDA approach:

- We get two Dirichlet distributions one associated documents to topics and other associated topics to words.
- With the help of these Dirichlet dist we generate a multinomial dist.



- First we generate topics and then we generate words.

⇒ ~~From~~ From ppt  
(mphil in advance computer science)

⇒ Introduction to probabilistic topic models ÷

- We want to find themes (or topics) in document.
  - useful for search or browsing
- We don't want to do supervised topic classification - rather not fix topics in advance, nor do manual annotation.
- Need an approach which automatically traces out the topics.
- This is essentially a clustering problem - can think of both words and documents are being clustered.

## ⇒ Key assumptions behind LDA

- Documents exhibit multiple topics
- LDA is a probabilistic approach model with a corresponding generative process.
  - each doc. is assumed to be generated by this (simple process)
- A topic is a distribution over a fixed vocabulary
  - these topics are assumed to be generated first, before the documents.
- Only the number of topics is specified in advance.

## ⇒ The Generative process

To generate a document:

1. Randomly choose a distribution over topics
2. For each word in the document:
  - a. Randomly choose a topic from the dist over topics
  - b. Randomly choose a word from the corresponding topic (dist over the vocabulary).



## ⇒ Topic Fingerprinting

### Document fingerprinting

A set of numbers that summarizes a document's content and allows you to perform simple math functions.

- We can use topic modelling to figure out what topics are present in a document and its relevance.

- Suppose we use topic modelling with  $\text{num of topics} = 50$ , then we can represent and visualize each document as a vector of  $\text{len} = 50$ .

We can use these vectors to find similar documents and use them for any other math operations.