

1a) LDA :-

latent Dirichlet

In natural language processing

allocation is a generative statistical model that allows a set of observations to be explained by unobserved latent groups that explain why some parts of the data are similar to observations are wordly collected into documents. It posits that each document is a mixture of small number of topics and that each word's collection is attribute to one of the document topics.

How to create the topics from the corpus

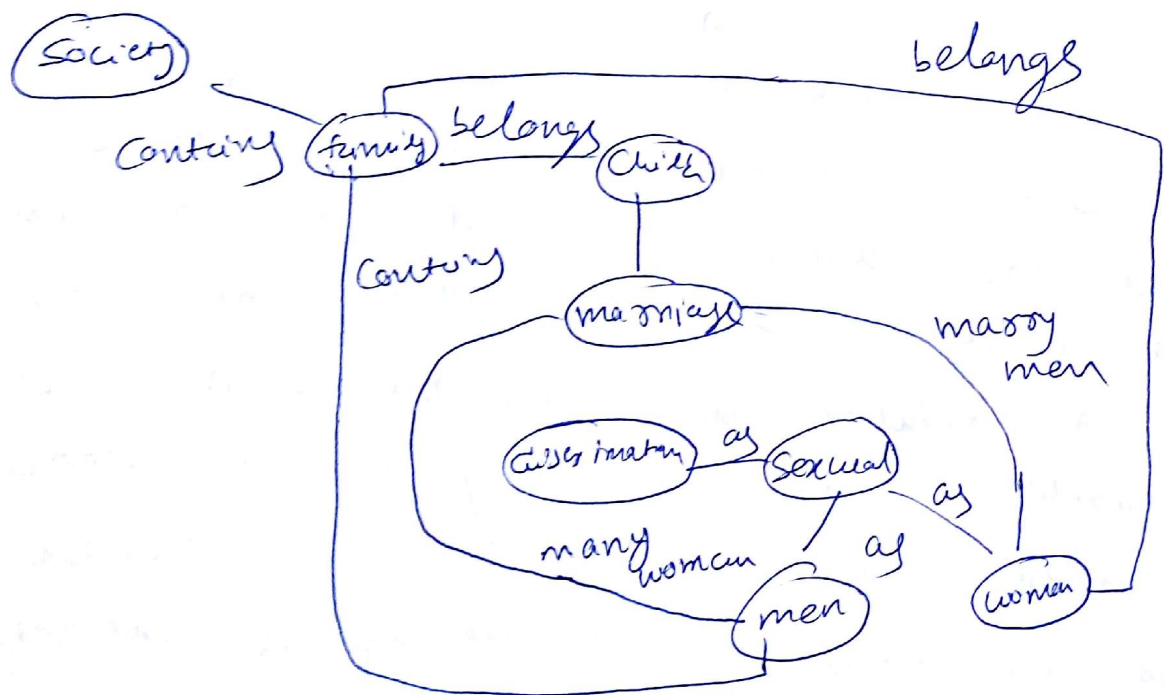
In LDA, each document may be viewed as a mixture of various topics where each document is considered to have a set of topics that are assigned to it. For example, an LDA model might have topics that can be classified as cat-related and dog-related. A topic has probabilities of generating various words. Such as milk, meat and kiltan which can be classified and implemented by the viewer as "cat-related" the dog. The dog-related topic likewise has the probability of generating each word.

1b) Knowledge graph for topic 3 in Yale Law Journal

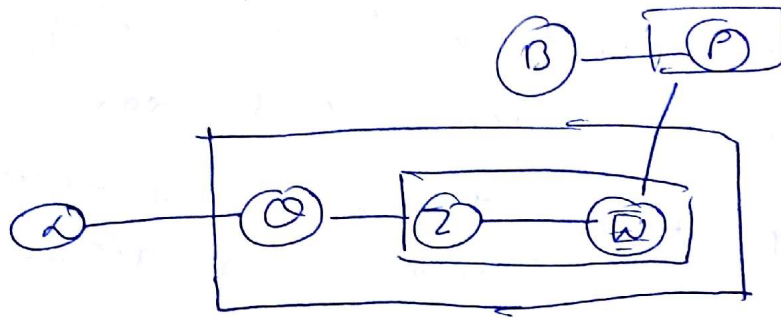
In the figure given in the problem set there are top eight topics were displayed each topic will be illustrated with its top-most frequent words. Each word's position along the x-axis denotes its specificity to the document.

Topic 3 in the Yale Law Journal - has the following words: women, sexual, men, sex-child, family-children, gender woman, marriage, discrimination, male, social, female, parently.

The most important words which were spread among the x-axis is the topic 3 as the basis for the construction of the knowledge graph



1c) Determining generativity or specificity of term in a topic



The dependencies among the many variables can be captured concisely. The boxes are plates representing display - the outer plate represents documents, while the inner plate represents the repeated choice of topics and words within a document.

### Generative process

Documents are represented as random mixtures over latent topics where each topic is characterised by a distribution over words. LDA assumes the following generative process for a corpus  $D$  consisting of  $m$  documents each of length  $n_i$ .

(1) Choose  $\theta_i \sim \text{Dir}(\alpha)$  where  $i \in \{1, \dots, m\}$  and  $\text{Dir}(\alpha)$  is a  $d$ -dimensional distribution.

(2) Choose  $\phi_k \sim \text{Dir}(\beta)$  where  $k \in \{1, \dots, K\}$

(3) for each word position  $i^j, j$  where  $j = \{1, \dots, n_i\}$  and  $i \in \{1, \dots, m\}$



## 1d) Inference algorithm in LDA :

The goal of topic modeling is to automatically discover the topics from a collection of documents. The documents and words are observed, the topic structure is hidden. The topics, per document topic distribution, per-document per-word topic assignment. We use observed variables to infer the hidden structure.

We can measure the content spread of each sentence by a word count.

Step 1 :- You tell the algorithm how many topics we think there are.

Step 2 :- The algorithm will assign every word to a temporary topic.

Step 3 :- The algorithm will check and update the topic assignments.

The posterior computation over hidden variables given a document.

$$P(z, \phi, \theta | w, \alpha, \beta) = P(z, \phi, \theta, w, \alpha, \beta) / P(w | \alpha, \beta)$$

The document is represented as continuous mixture.

$$P(w | \alpha, \beta) = \sum_k P(\phi | \alpha) \left( \prod_{n=1}^N P(w_n | \phi, \beta) \right)$$

for topic  $k$ , term  $v$ .

$$\alpha_{kv} = \beta_{kv} + \sum_n \mathbb{I}(w_n = v) \phi_{kn}$$

2a) Given the term/document matrix

Document	online	festival	books	figure	delhi
D <sub>1</sub>	1	0	1	0	1
D <sub>2</sub>	2	1	2	1	1
D <sub>3</sub>	0	0	1	1	1
D <sub>4</sub>	1	2	0	2	0
D <sub>5</sub>	3	1	0	0	0
D <sub>6</sub>	0	1	1	1	2
D <sub>7</sub>	2	0	1	2	1
D <sub>8</sub>	1	1	0	1	0
D <sub>9</sub>	1	0	2	0	0
D <sub>10</sub>	0	1	1	1	1

Step 1:

Given also the distance matrix, there are 3 clusters D<sub>2</sub>, D<sub>5</sub>, D<sub>7</sub> as per the diagram as we get distance as 0.0 for above 3 which means that D<sub>2</sub>, D<sub>5</sub>, D<sub>7</sub> are the centroids. The remaining documents have moved into those 3 different clusters using K-means  $K=3$

D<sub>2</sub> :- D<sub>1</sub>, D<sub>6</sub>, D<sub>9</sub>, D<sub>10</sub>    D<sub>1</sub> :- D<sub>3</sub>, D<sub>4</sub>,    D<sub>5</sub> :- D<sub>8</sub>

The first row of the distance matrix corresponds to the distance of each object to the first centroid and the second row is the distance of each object to the second centroid and so on. Minimum distance grouping is done.



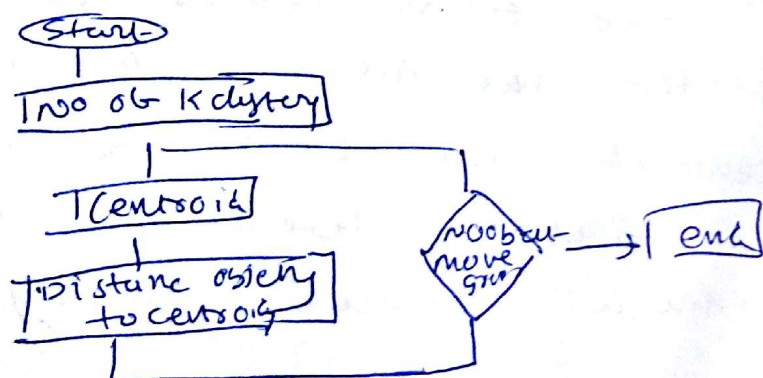
## 2) Clustering:

Clustering/segmentation is one of the most important techniques used in acquisition analysis. It is the process of making a group of abstract objects into classes or the process of making a group of abstract objects into classes or the process of making a group of similar objects. We will partition the observation into a cluster in such a way that they are similar in sense.

Clustering is a method of unsupervised learning, and a common technique for the statistical data analysis used in many fields.

### K-means Clustering:

K-means clustering is an algorithm to classify or to group your objects based on attributes/feature into  $K$  number of groups.  $K$  is positive integer number. The grouping is done by minimizing the sum of squares of distance between data and the corresponding cluster centroid.



There are 3 centroids randomly taken

$$D_2 (2, 1, 2, 1) \quad D_5 (3, 1, 0, 0, 0) \quad D_7 (2, 0, 1, 2, 1)$$

Step 2 :- Now calculate the distance for  $D_1$  from

$D_2, D_5, D_7$

$$D_1 \rightarrow D_2$$

$$\sqrt{(1-2)^2 + (0-1)^2 + (1-2)^2 + (1-0)^2 + (1-1)^2} = \sqrt{1+1+1+1+0} = \sqrt{4} = 2$$

$$D_1 \rightarrow D_5$$

$$\sqrt{(1-3)^2 + (0-1)^2 + (1-0)^2 + (1-0)^2} = \sqrt{4+1+1+1} = \sqrt{7} \approx 2.6$$

$$D_1 \rightarrow D_7$$

$$\sqrt{(1-2)^2 + (0-0)^2 + (1-1)^2 + (0-2)^2 + (1-0)^2} = \sqrt{1+0+0+4+1} = \sqrt{6} \approx 2.2$$

Step 3 :-

group the data into clusters based on their minimum distance.

$$D_2 = \{D_1, D_6, D_9, D_{10}\}$$

$$D_7 = \{D_3, D_4\}$$

$$D_5 = \{D_8\}$$

In the above steps using the K-means algorithm we will cluster the data points based on the centroids and we will reiterate this process by calculating the new mean and new clusters.

2b) The difference between k-means and LDA are as follows

- Both are applied to assign  $k$  topics to set of  $N$  documents. k-means is going to partition the  $N$  documents in  $k$  distinct clusters while LDA assigns a document to a mixture of topics.
- k-means is hard clustering while LDA is soft clustering

LDA Pros:

- LDA is in the exponential family and consequently to the multinomial distribution
- feature set is reduced
- one document can be associated with multiple topics

Cons:

- unable to capture the correlation between the different topics.

k-means Pros:

- simple, easy to implement
- easy to interpret the clustering result.
- It is a great solution for pre-clustering, reducing the space into distinct smaller sub-spaces where other clustering algorithms can be applied

k-means Cons:

- difficult to predict  $k$ -value
- with global clusters, it didn't work well.