

CS5560 Knowledge Discovery and Management

Problem Set 5

July 3 (T), 2017

Name: Revanth Chakilam

Class ID: 02

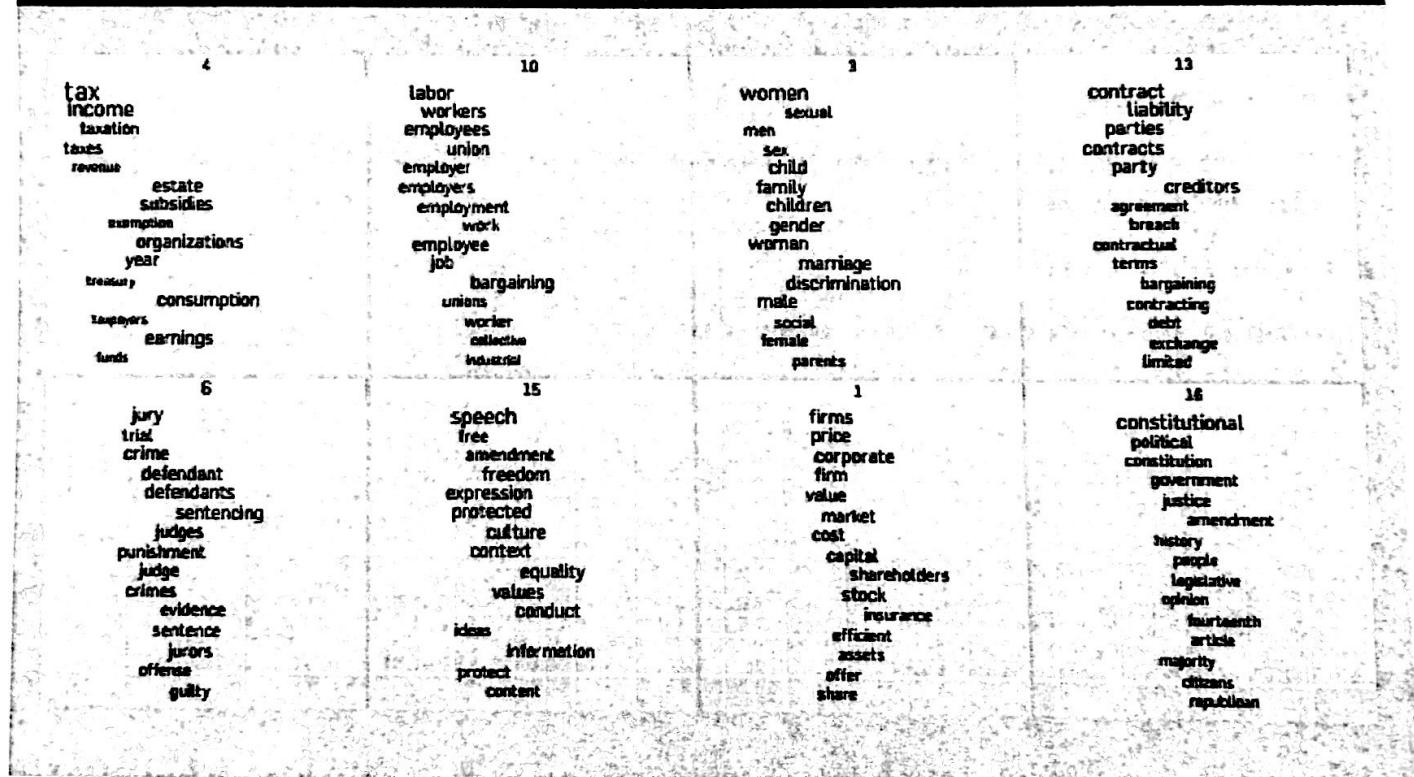
1. LDA

Read the following articles to learn more about LDA

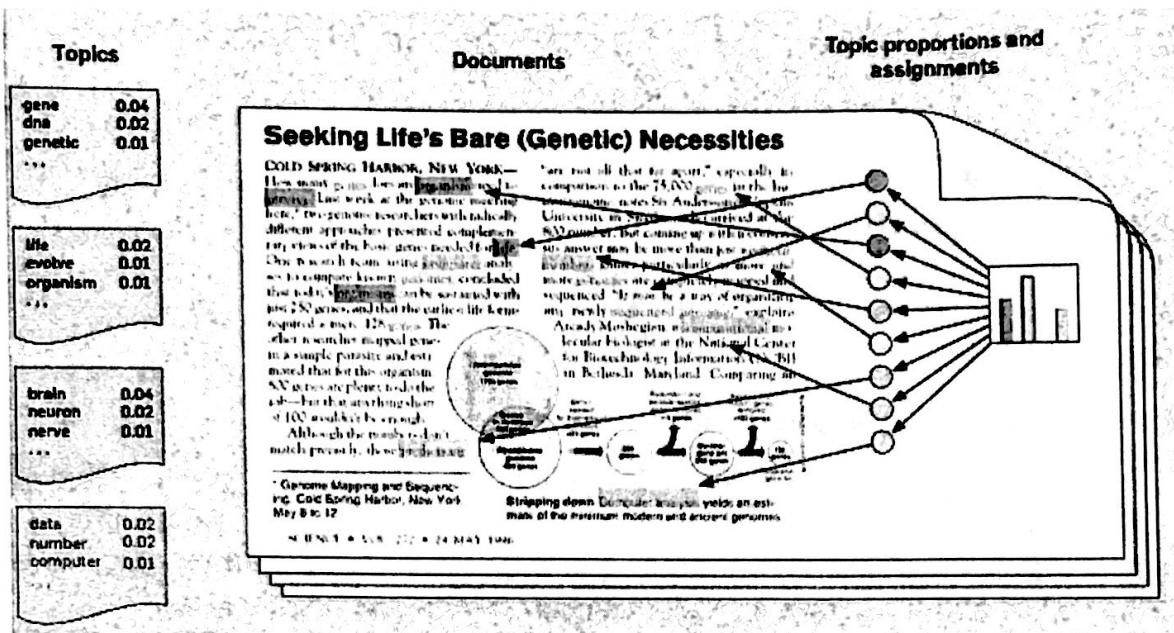
- <https://algobean.com/2015/06/21/laymans-explanation-of-topic-modeling-with-lda-2/>
- <http://engineering.intenthq.com/2015/02/automatic-topic-modelling-with-lda/>

Consider the topics discovered from Yale Law Journal. (Here the number of topics was set to be 20.) Topics about subjects like about discrimination and contract law.

Figure 3. A topic model fit to the *Yale Law Journal*. Here, there are 20 topics (the top eight are plotted). Each topic is illustrated with its topmost frequent words. Each word's position along the x-axis denotes its specificity to the documents. For example "estate" in the first topic is more specific than "tax."



- Describe the overall process to generate such topics from the corpus.
- Draw a knowledge graph for Topic 3 in Yale Law Journal (The First Figure).
- Each topic is illustrated with its topmost frequent words. Each word's position along the x-axis denotes its specificity to the documents. For example "estate" in the first topic is more specific than "tax." (the second figure). Describe how to determine the generality or specificity of the terms in a topic.
- Describe the inference algorithm that was used in LDA.



2. K-means clustering vs. LDA

Read the K-means clustering for text clustering from <https://www.experfy.com/blog/k-means-clustering-in-text-data>

- (a) Describe the steps how the following 10 documents have moved into 3 different clusters using clustered using k-means (K=3).

Document/Term Matrix

| Documents | Online | Festival | Book | Flight | Delhi |
|-----------|--------|----------|------|--------|-------|
| D1 | 1 | 0 | 1 | 0 | 1 |
| D2 | 2 | 1 | 2 | 1 | 1 |
| D3 | 0 | 0 | 1 | 1 | 1 |
| D4 | 1 | 2 | 0 | 2 | 0 |
| D5 | 3 | 1 | 0 | 0 | 0 |
| D6 | 0 | 1 | 1 | 1 | 2 |
| D7 | 2 | 0 | 1 | 2 | 1 |
| D8 | 1 | 1 | 0 | 1 | 0 |
| D9 | 1 | 0 | 2 | 0 | 0 |
| D10 | 0 | 1 | 1 | 1 | 1 |

Distance Matrix

| Documents | Distance from 3 clusters | | | | |
|-----------|--------------------------|-----|-----|---------------|----------|
| | D2 | D5 | D7 | Min. Distance | Movement |
| D1 | 2.0 | 2.6 | 2.2 | 2.0 | D2 |
| D2 | 0.0 | 2.6 | 1.7 | 0.0 | |
| D3 | 2.4 | 3.6 | 2.2 | 2.2 | D7 |
| D4 | 2.8 | 3.0 | 2.6 | 2.6 | D7 |
| D5 | 2.6 | 0.0 | 2.8 | 0.0 | |
| D6 | 2.4 | 3.9 | 2.6 | 2.4 | D2 |
| D7 | 1.7 | 2.8 | 0.0 | 0.0 | |
| D8 | 2.6 | 2.0 | 2.8 | 2.0 | D5 |
| D9 | 2.0 | 3.0 | 2.6 | 2.0 | D2 |
| D10 | 2.2 | 3.5 | 2.4 | 2.2 | D2 |

(b) Describe the difference (pro and con) of k-means clustering and the LDA topic discovery model.

1. ② LDA :- (Latent Dirichlet Allocation)

When we consider NLP - Natural Language Processing, LDA is a generative statistical model that allows sets of observations to be explained by unobserved groups that explain why some parts of the data are similar. If observations are words collected into documents. It postulates that each document is a mixture of small number of topics and that each words collection is attributable to one of the document topics.

→ Creation of topics from Corpus

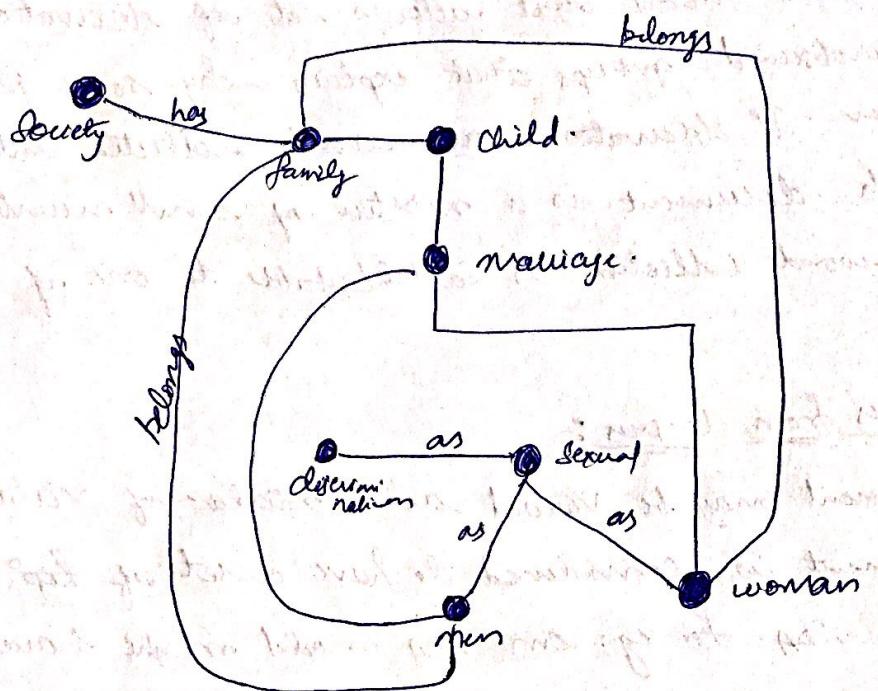
In LDA, each document may be viewed as a mixture of various topics where each document is considered to have a set of topics that are assigned to it via LDA. For eg: our LDA model might have topics that can be classified as CAT related or DOG-Related. A topic has probabilities of generating various words which are classified and interpreted by the viewer as "CAT-related". And the other related topic like wise has the probabilities of generating each word. Puppy, bark, etc might have high probability.

③ KG for the topic-3 in Yale Law Journal:

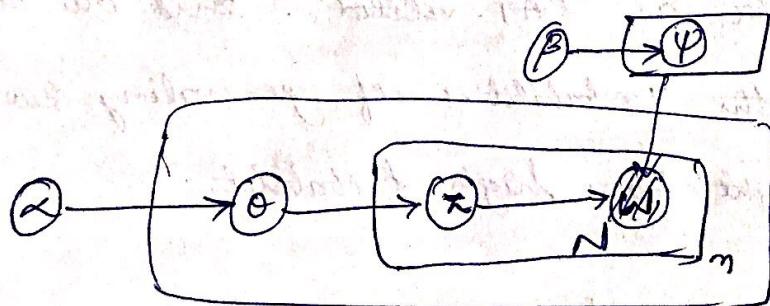
If we consider the diagram in the question, there were around 8 topics were displayed. Each topic will be illustrated with its top-most frequent words. Each words position along the X-axis denotes its specificity to the documents.

(2)

The most important words which were spread among the X-axis is the topic 3 on the basis for the construction of Knowledge Graph.



② To determine generality (or) specificity of terms in a topic 8.



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

Generative process:

Documents are represented as random mixtures over latent topics where each topic is characterized by a distribution over words.

LDA assumes the following generation process for a Corpus D .

- Choose $\theta_i \sim Dir(\alpha)$ where $i \in \{1, \dots, M\}$ and $Dir(\alpha)$ is a dirichlet distribution.
- Choose $\psi_k \sim Dir(\beta)$ where $k \in \{1, \dots, K\}$
- For each word position i, j where $j \in \{1, \dots, N_i\}$ and $i \in \{1, \dots, M\}$

D) Inference Algorithms in LDA:

The goal of the topic modelling 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, are obscured variables to infer the hidden structure.

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

- Step - ① : The algorithm tells how many topics we think are there.
 - ② : The algorithm will assign every word to a temporary topic.
 - ③ : The algorithm will check and update the topic assignments.
- The posterior computation over hidden variables given a document.

$$P(z_1, \phi, \theta | w_1, \alpha, \beta) = P(z_1, \phi, \theta, w_1 | \alpha, \beta) / P(w_1 | \alpha, \beta)$$

$$P(w_1 | \alpha, \beta) = \int P(\theta | \alpha) \left(\prod_{n=1}^N P(w_n | \theta, \beta) \right) d\theta.$$

$$\eta_{kv} = \beta_{kv} + \sum_{d=1}^D \sum_{n=1}^{N_d} \mathbb{I}(w_{dn}=v) \psi_{dkv}$$

$$\psi_{dkv} \propto \exp \left\{ \text{Eq} \left[\log(\theta_{dk}) + \log(\phi_{kvn}) \right] \right\}$$

2Q] K-means clustering VS LDA :-

Clustering: Clustering / Segmentation is one of the important techniques used in acquisition analytics. It is the process of making a group of abstract objects into classes of the similar objects. We will partition the observations into a cluster in such a way that they are similar in sense. A method of unsupervised learning, and a common technique for the statistical data analysis and in many fields.

K-means clustering: K-means clustering is an algorithm to classify or to group your objects based on attribute/features into k-number of groups. K is positive integer number. The grouping is done by minimizing the sum of squares of distance b/w data and the corresponding cluster centroid.

(a) Document Text Matrix:

In the given figures, there are total 10 documents.

→ Distance matrix is also provided. There are 3 clusters D_2, D_5, D_7 as per the diagram as we get distance as 0.0 for above 3 which indicates that D_2, D_5, D_7 are the centroids. The remaining documents have moved onto other 3 different clusters using k-means $k=3$.

$$D_2 : - D_1, D_6, D_9, D_{10}$$

$$D_7 : - D_3, D_4, \quad D_5, D_8$$

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 based on minimum distance grouping is done.

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

→ Let's calculate the distance from D1 from D2, D5, D7

$$D_1 \rightarrow D_2$$

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

$$D_1 \rightarrow D_5$$

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

$$D_1 \rightarrow D_7$$

$$\sqrt{(1-2)^2 + (0-0)^2 + (1-1)^2 + (0-2)^2 + (1-1)^2} = \sqrt{1+4} = \sqrt{5} = 2.2$$

Same way calculating the distance from each point to centroid.

→ Let now group the data into clusters based on these minimum distance.

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

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

$$D_5 : \{D_8\}$$

Using k-means algorithm we will cluster the data points based on the centroid and we will re-iterate this process by calculating the new mean and new clusters.

(b) Differences b/w k-means and the LDA are follows:-

→ If both can applied to assign k-topics to a set of N-documents k-mean is going to partition the N-documents in k-disjoint clusters while LDA assigns a document to a mixture of topics.

→ k-means is hard clustering while LDA is soft clustering

Advantages: (LDA)

→ LDA is in the exponential family and conjugate to the multinomial distributions.

⑥

- feature set is reduced
- One document can be annotated with multiple topics.

Disadvantages: LDA

- unable to capture the correlation b/w the different topics.

Advantages: K-means

- simple, easy to implement
- easy to interpret the clustering result.
- It is a great solution for pre-clustering, reducing the space into disjoint smaller sub-spaces where other clustering algorithms can be applied.
- The clusters are non-hierarchical and they don't overlap.
- It is computationally faster.

Disadvantages: K-means

- Difficult to predict k-value.
- With global cluster, it didn't work well
- Doesn't work well with non-circular cluster shape - number of cluster and initial seed value need to be specified beforehand.
- Applicable only when mean is specified.
- Sensitive to the outliers.