

CS5560 Knowledge Discovery and Management

Problem Set 5

July 3 (T), 2017

Name: Yalamanchili Sowmya
Class ID: 16246716 (30)

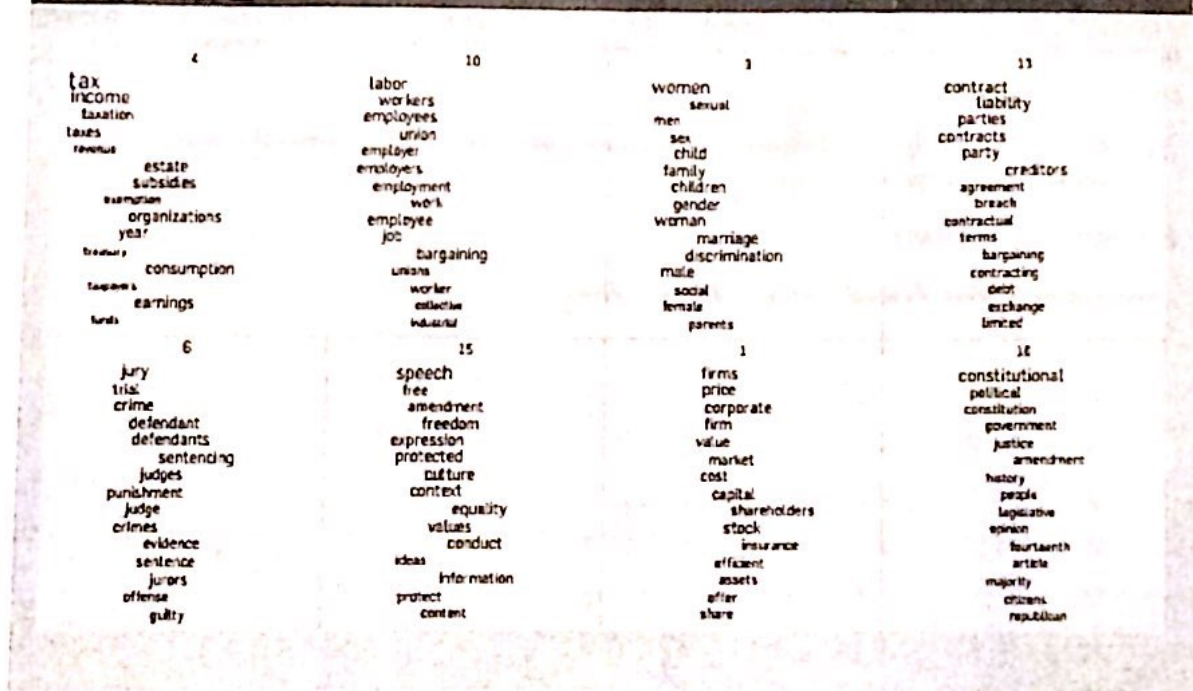
1. LDA

Read the following articles to learn more about LDA

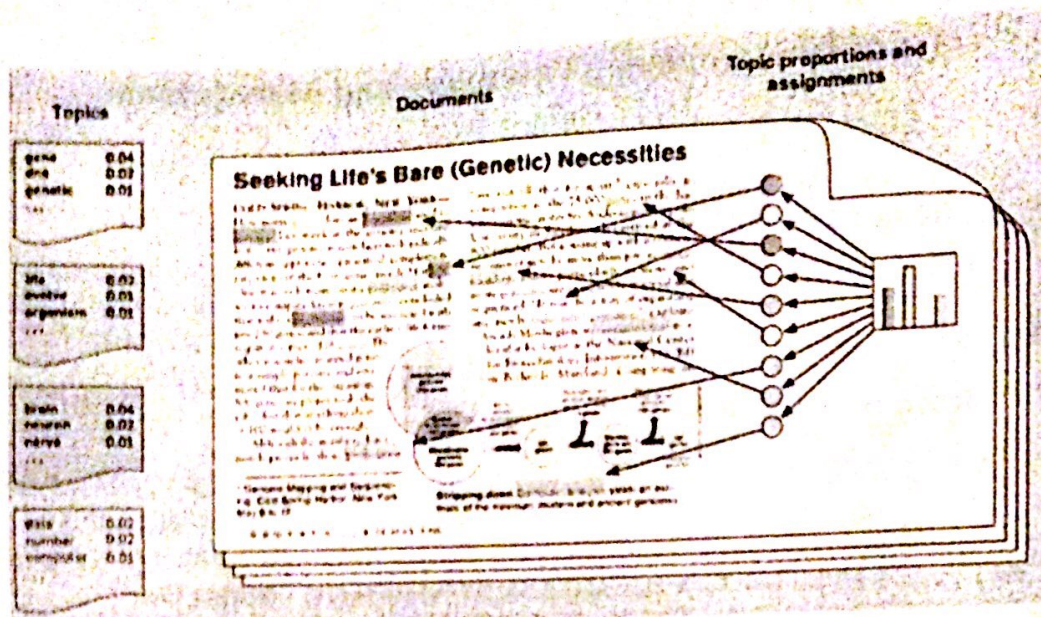
- <https://alqobbeans.com/2015/06/21/laymans-explanation-of-topic-modeling-with-lda-2/>
- <http://engineering.intenthq.com/2015/02/automatic-topic-modelling-with-lda/>

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 top-most 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				Movement
	D2	D5	D7	Min. Distance	
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 (a) LDA (Latent Dirichlet allocation)

In natural language processing, latent Dirichlet allocation 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 posits that each document is a mixture of small number of topics and that each word's collection is attributable to one of the document's topics.

How to create the topics from the corpus?

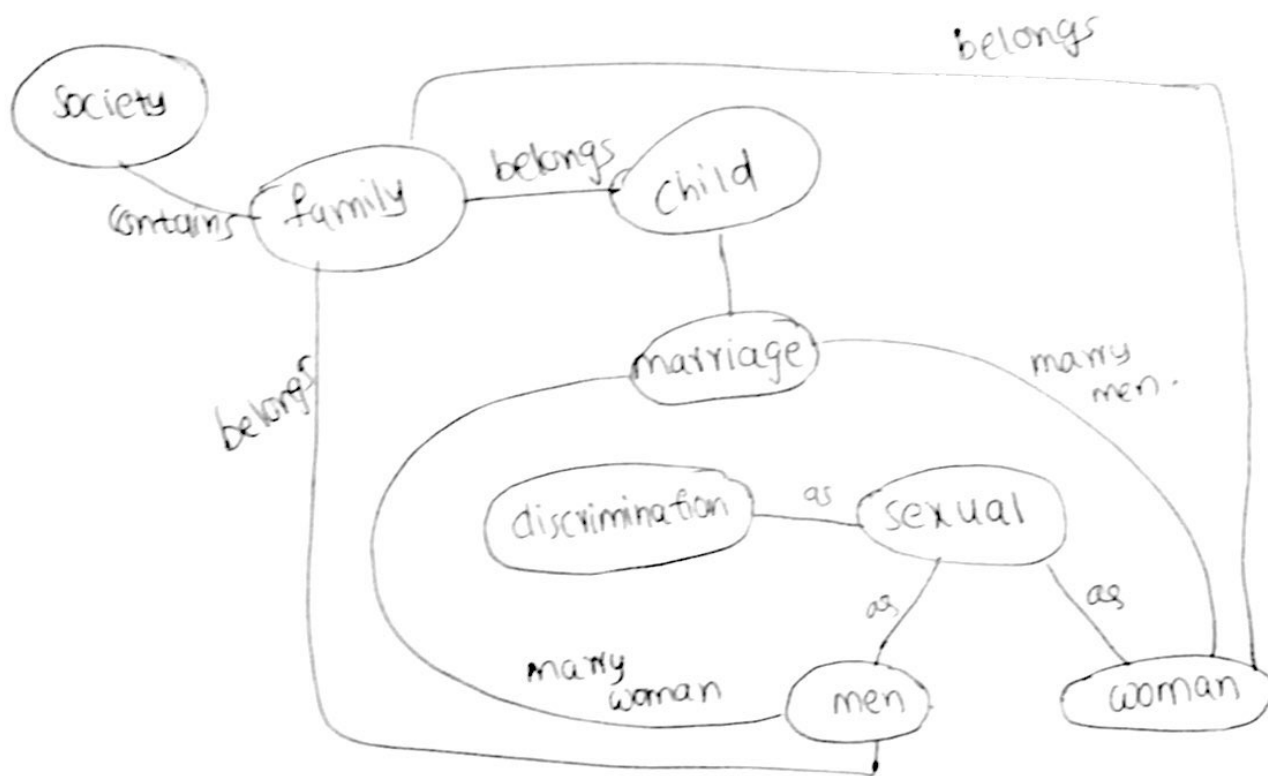
In LDA, each document may be viewed as a mixture of various topics where each document is considered to have a set of topics that are assigned to it via LDA. For example, an LDA model might have topics that can be classified as CAT related and DOG related. A topic has probabilities of generating various words, such as milk, meow and kitten which can be classified and interpreted by the viewer as "CAT related". The DOG related topic likewise has the probabilities of generating each word: puppy, bark and bone might have high probability.

① (b) Knowledge graph for Topic 3 in Yale Law Journal

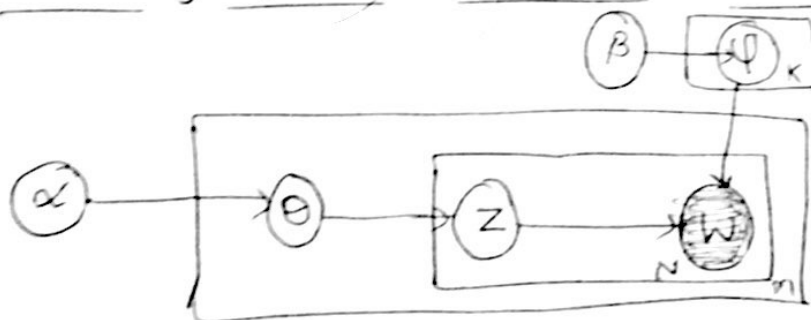
In the figure given in the problem set there are top eight topics were displayed. Each topic will be illustrated with its top-most frequent words. Each word's position along the x-axis denotes its specificity to the documents.

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

The most important words which were spread among the 1-axis is the topic 3 are the basis for the construction of the knowledge graph.



1c) Determining generality or specificity of terms in a topic.



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

Generative process

Name: yalamanchili sumy

class ID: 30

Documents are represented as random mixtures over latent topics, where each topic is characterized 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 Dirichlet distribution
- ② choose $\varphi_k \sim \text{Dir}(\beta)$ where $k \in \{1, \dots, K\}$
- ③ for each word positions i, j where $j \in \{1, \dots, N_i\}$ and $i \in \{1, \dots, M\}$

The generality and specificity of the terms was determined by their document frequency (DF) the more documents a term occurred in, the more general it was assumed to be.

① 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 infer 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

$$p(z, \phi, \theta | w, \alpha, \beta) = \frac{p(z, \phi, \theta, w | \alpha, \beta)}{p(w | \alpha, \beta)}$$

The document represented as continuous mixture:

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

for topic k , term v

$$\lambda_{kv} = \beta_{kv} + \sum_d \sum_n \mathbb{I}[w_{dn} = v] \varphi_{dnk}$$

for each document d $\gamma_{dk} = \alpha_k + \sum_n \varphi_{dnk}$

for each word n

$$\varphi_{dnk} \propto \exp \left\{ E_q \left[\log(\theta_{dk}) + \log(\phi_{kwn}) \right] \right\}$$

② clustering

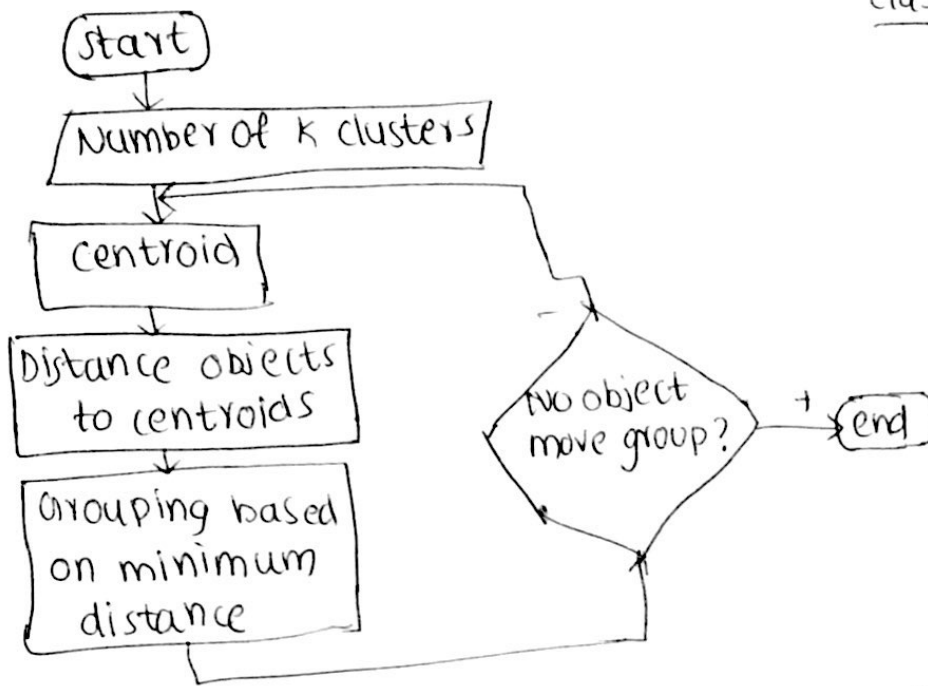
(clustering) segmentation is one of the most important techniques used in Acquisition Analytics. It is the process of making a group of abstract objects into classes of the similar objects. We will partition the observations into a cluster in such a way that they are similar in sense.

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/features into k number of groups. k is positive integer number.

The grouping is done by minimizing the sum of squares of distances between data and the corresponding cluster centroid.



The figure describes how the K-means clustering Algorithm works.

(2a) Given the term/document matrix

Documents	online	Festival	Book	flight	Delhi
D ₁	1	0	1	0	1
D ₂	2	1	2	1	1
D ₃	0	0	1	1	1
D ₄	1	2	0	2	0
D ₅	3	1	0	0	0
D ₆	0	1	1	1	2
D ₇	2	0	1	2	1
D ₈	1	1	0	1	0
D ₉	1	0	2	0	0
D ₁₀	0	1	1	1	1

Total
10
documents

step 1:-

Given also the distance matrix. There are 3 clusters D₂, D₅, D₇ as per the diagram as we get distance as 0.0 for above 3 which indicates that D₂, D₅, D₇ are the centroids. The remaining documents have moved into those 3 different clusters using K-means K=3.

p₂:- D₁, D₆, D₉, D₁₀ D₇:- D₃, D₄, D₅:- D₈

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.

There are 3 centroids randomly taken

Step 2: $D_2 (2, 1, 2, 1, 1)$ $D_5 (3, 1, 0, 0, 0)$ $D_7 (2, 0, 1, 2, 1)$.

Now calculate the distance for D_1 from D_2, D_5, D_7

$D_1 \rightarrow D_2$

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

$D_1 \rightarrow D_5$

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

$D_1 \rightarrow D_7$

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

likewise we will calculate the sum of squares of distance from each data point to the centroid.

Step 3:

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 :- $\{D_8\}$

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

(2b) The differences between k-means and the LDA are as follows.

- If both are applied to assign k topics to a set of N documents, k-means 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

LDA pros

- LDA is in the exponential family and conjugate 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 disjoint smaller sub-spaces where other clustering algorithms can be applied
- The clusters are non-hierarchical and they do not overlap
- It is computationally faster.
- The clusters are globular.

k-means cons

- Difficult to predict k-value
- With global clusters, 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.