

Name: Sadanand Kallakuri

Class ID: 9

PROBLEM SET 5

1.a. Latent Dirichlet Allocation (LDA)

It is an probabilistic is used to generate the topics. LDA is the iterative model which requires 3 parameters, which are number of topics and deep knowledge of dataset.

The performance of LDA using perplexity is evaluated. To evaluate the LDA model, one document is taken and split in two. The first half is fed into LDA to compute the topics composition, from that composition, then, the word distribution is estimated. This is then composed with the word distribution of the 2nd half of the document. Then a measure of distance is extracted. Perplexity is often used to select the best no. of topics of the LDA model.

LDA Algorithm

Input: Words $w \in$ documents d

Output: Topic assignments z and counts n_{dk} , n_{kw} and n_k

begin

randomly initialize z and increment counters

for each iteration do

for $i=0 \rightarrow N-1$ do

word $\leftarrow w[i]$

topic $\leftarrow z[i]$

$n_{d, \text{topic}} = 1$; $n_{\text{word}, \text{topic}} = 1$; $n_{\text{topic}} = 1$

for $k=0 \rightarrow k-1$ do

$$p(z=k) = (n_{d,k} + \alpha_k) \frac{n_{k,w} + \beta_w}{n_k + \beta \times W}$$

end

topic \leftarrow sample from $p(z)$

$z[i] \leftarrow$ topic

$n_{d,topic} += 1$; $n_{word,topic} += 1$; $n_{topic} += 1$

end

end

return $z, n_{d,k}, n_{k,w}, n_k$

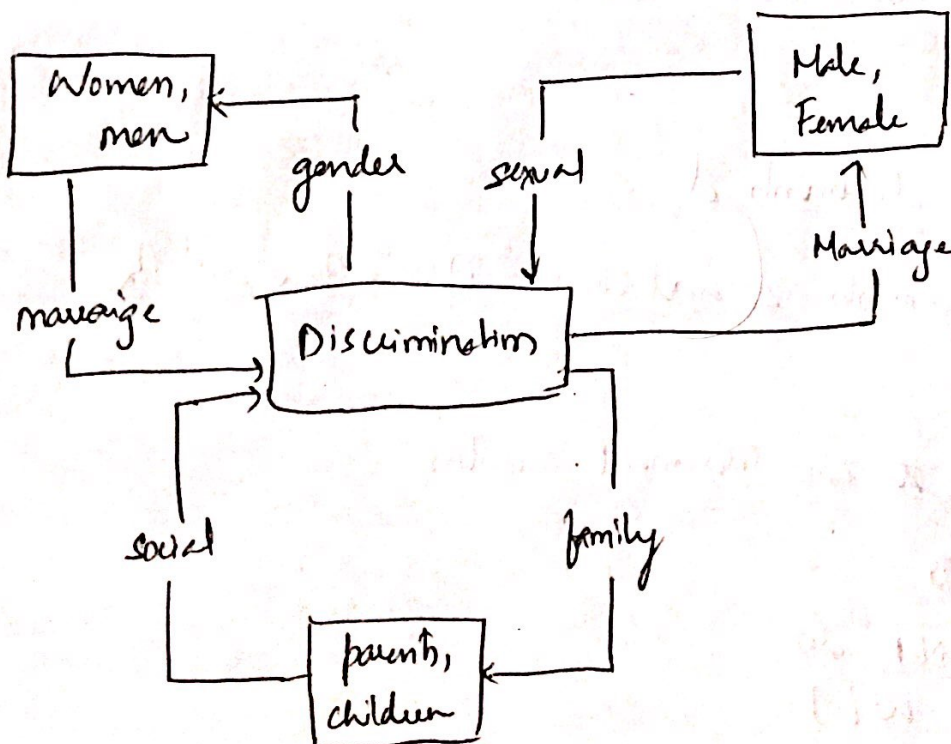
end

Step 1: Decide how many topics we need.

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

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

1b) Knowledge Graph



1c) How much prevalent are topics in the document?

Since the words in Doc Y are assigned to Topic F and Topic P in a 50-50 ratio, the remaining "fish" words seems equally likely to be about either topic.

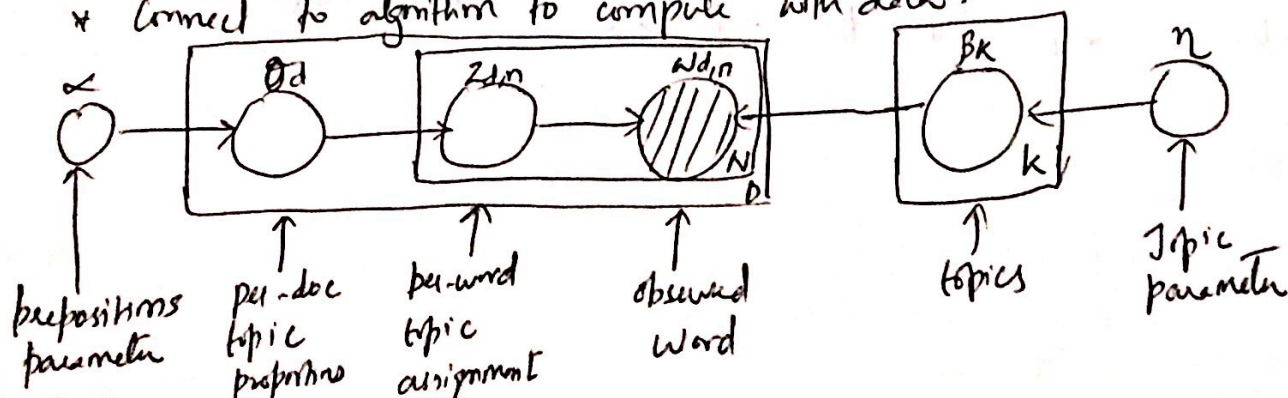
	Doc X		Doc Y
F	Fish	?	Fish
F	Fish	F	Fish
F	Eat	F	Milk
F	Eat	P	kitten
F	Vegetables	P	kitten

1d) Inference algorithm used in LDA

- * Each topic is a distribution over words.
- * Each document is a mixture of corpus-wide topics.
- * Each word is drawn from one of the topics.
- * Observe the documents.
- * Goal is to infer the hidden variables, compute distribution conditional on the documents.

$P(\text{topics, proportions, assignments} / \text{documents})$

- * Encode assumption.
- * Define a factorization of the joint distribution.
- * Connect to algorithm to compute with data.



$$P(\beta, \theta, z, w) = \left(\prod_{i=1}^K p(\beta; \ln) \right) \left(\prod_{d=1}^D P(z_{d,n} | \theta_d) P(w_{d,n} | \beta, z_{d,n}) \right)$$

2. a) k-means clustering vs LDA

We have to create $K=3$ clusters.

Lets choose D_2, D_5 and D_7 as initial three seeds.

Now we have to calculate euclidean distance from other documents to D_2, D_5 and D_7 .

$O \rightarrow$ Online $F \rightarrow$ Festival $B \rightarrow$ Book $P \rightarrow$ Flight $D \rightarrow$ Delhi

$$D_1 \text{ to } D_2 = \sqrt{(O_1 - O_2)^2 + (F_1 - F_2)^2 + (B_1 - B_2)^2 + (P_1 - P_2)^2 + (D_1 - D_2)^2}$$

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

$$D_1 \text{ to } D_5 = \sqrt{(1-3)^2 + (0-1)^2 + (1-0)^2 + (0-0)^2 + (1-0)^2} = \sqrt{7} = 2.6$$

$$D_1 \text{ to } D_7 = \sqrt{(1-2)^2 + (0-0)^2 + (1-1)^2 + (0-2)^2 + (1-1)^2} = \sqrt{5} = 2.2$$

$$D_2 \text{ to } D_2 = 0$$

$$D_2 \text{ to } D_5 = \sqrt{(2-3)^2 + (1-1)^2 + (2-0)^2 + (1-0)^2 + (1-0)^2} = \sqrt{7} = 2.6$$

$$D_2 \text{ to } D_7 = \sqrt{(2-2)^2 + (1-0)^2 + (2-1)^2 + (1-2)^2 + (1-1)^2} = \sqrt{3} = 1.7$$

$$D_3 \text{ to } D_2 = \sqrt{6} = 2.4$$

$$D_4 \text{ to } D_2 = \sqrt{8} = 2.8$$

$$D_7 \text{ to } D_7 = 0$$

$$D_3 \text{ to } D_5 = \sqrt{13} = 3.6$$

$$D_4 \text{ to } D_5 = \sqrt{9} = 3$$

$$D_7 \text{ to } D_2 = \sqrt{3} = 1.7$$

$$D_3 \text{ to } D_7 = \sqrt{5} = 2.2$$

$$D_4 \text{ to } D_7 = \sqrt{7} = 2.6$$

$$D_7 \text{ to } D_5 = \sqrt{8} = 2.8$$

$$D_5 \text{ to } D_2 = \sqrt{7} = 2.6$$

$$D_6 \text{ to } D_2 = \sqrt{6} = 2.4$$

$$D_8 \text{ to } D_2 = \sqrt{6} = 2.4$$

$$D_5 \text{ to } D_5 = 0$$

$$D_6 \text{ to } D_5 = \sqrt{15} = 3.8$$

$$D_8 \text{ to } D_5 = \sqrt{5} = 2.2$$

$$D_5 \text{ to } D_7 = \sqrt{8} = 2.8$$

$$D_6 \text{ to } D_7 = \sqrt{7} = 2.6$$

$$D_8 \text{ to } D_7 = \sqrt{5} = 2.2$$

$$D_9 \text{ to } D_2 = \sqrt{4} = 2$$

$$D_9 \text{ to } D_5 = \sqrt{9} = 3$$

$$D_9 \text{ to } D_7 = \sqrt{7} = 2.6$$

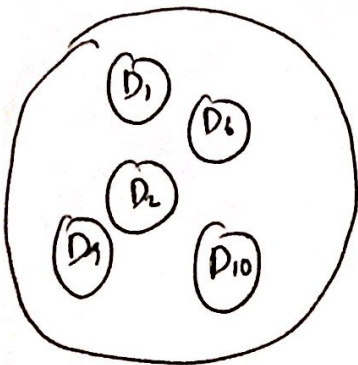
$$D_{10} \text{ to } D_2 = \sqrt{5} = 2.2$$

$$D_{10} \text{ to } D_5 = \sqrt{12} = 3.4$$

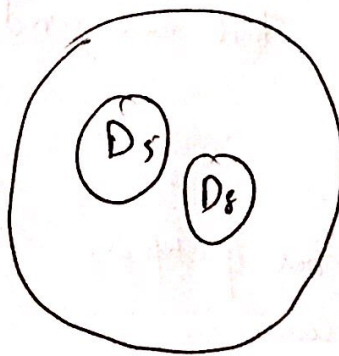
$$D_{10} \text{ to } D_7 = \sqrt{6} = 2.4$$

<u>Documents</u>	<u>D_2</u>	<u>D_5</u>	<u>D_7</u>	<u>Min distance</u>	<u>Cluster</u>
D_1	2.0	2.6	2.2	2.0	D_2
D_2	0.0	2.6	1.7	0.0	D_2
D_3	2.4	3.6	2.2	2.2	D_7
D_4	2.8	3.0	2.6	2.6	D_7
D_5	2.6	0.0	2.8	0.0	D_5
D_6	2.4	3.9	2.6	2.4	D_2
D_7	1.7	2.8	0.0	0.0	D_7
D_8	2.6	2.0	2.8	2.0	D_5
D_9	2.0	3.0	3.6	2.0	D_2
D_{10}	2.2	3.5	2.4	2.2	D_2

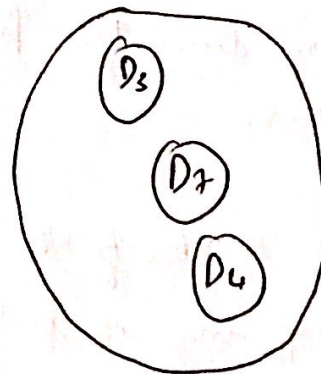
D_2 Cluster



D_5 Cluster



D_7 Cluster



2. b) Pros & Cons of K-Means Clustering:

Pros:

- ① Fast \rightarrow Computational cost $\rightarrow O(k \times n \times d)$
- ② Robust and easier to understand.
- ③ Gives best result when data set are distinct or well separated from each other.
- ④ It is a great solution for pre-clustering
- ⑤ Works great for spherical clusters.

Cons:

- ① K-value is not known and is difficult to predict.
- ② Does not work well with clusters of different size and different density.

LDA topic Discovery Model

Pros:

- ① We can infer the content spread of each sentence by a word count.
- ② We can derive the proportions that each word constitutes in given topics.

Cons:

- ① We have to specify the number of topics.
- ② LDA's efficiency is pretty low.
- ③ LDA cannot capture correlations.
- ④ Unsupervised
- ⑤ Uses BOV.