

①② LDA
Latent Dirichlet

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

To evaluate a LDA model, one document is taken and split in two. The first half is fed to LDA to compute the topics composition from that composition then word distribution is estimated. This distribution is then computed with word distribution of 2nd half of document. A measure of distance is extracted.

LDA Algorithm

Input : words $w \in$ documents d
Output : topic assignments z and counts $n_d, k, n_{k,w}$ and n_k
begin.

```

Randomly initialize  $z$  and increment counter
for each iteration  $n$  do
  for  $i \leftarrow 0 \rightarrow N-1$  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
  
```

$$p(z=k | i) = (n_{d,k} + \alpha_k) \frac{n_{k,w} + \beta_w}{n_k + \beta_{\text{w}}}$$

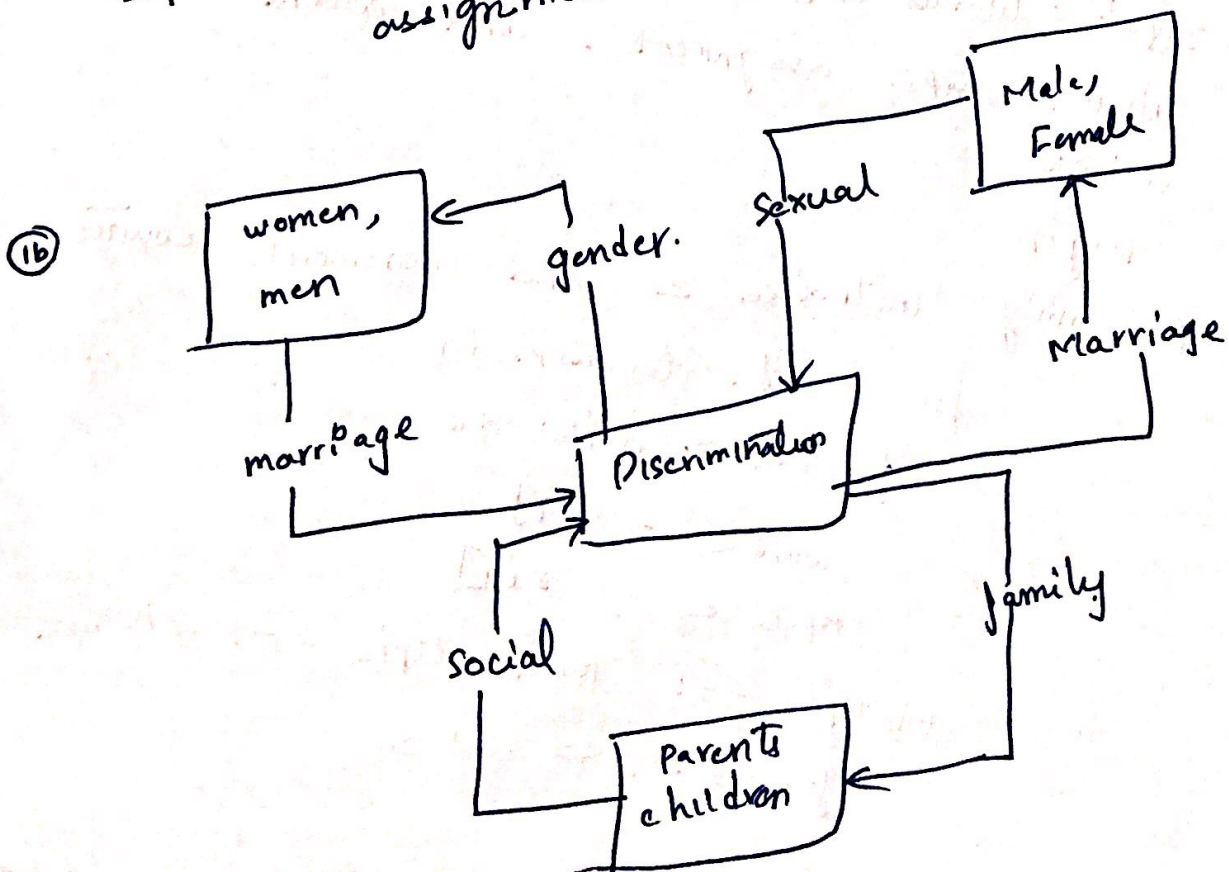
end
 topic \leftarrow sample from $p(z | i)$
 $z[i] \leftarrow$ topic
 $n_{d, \text{topic}} += 1$, $n_{\text{word}, \text{topic}+1} = 1$, $n_{\text{topic}} += 1$
 end
 end

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

end

step 1 : Decide how many topics we need
 The algorithm will assign every word to a temporary topic

step 2 : The algorithm will check or update topic assignment

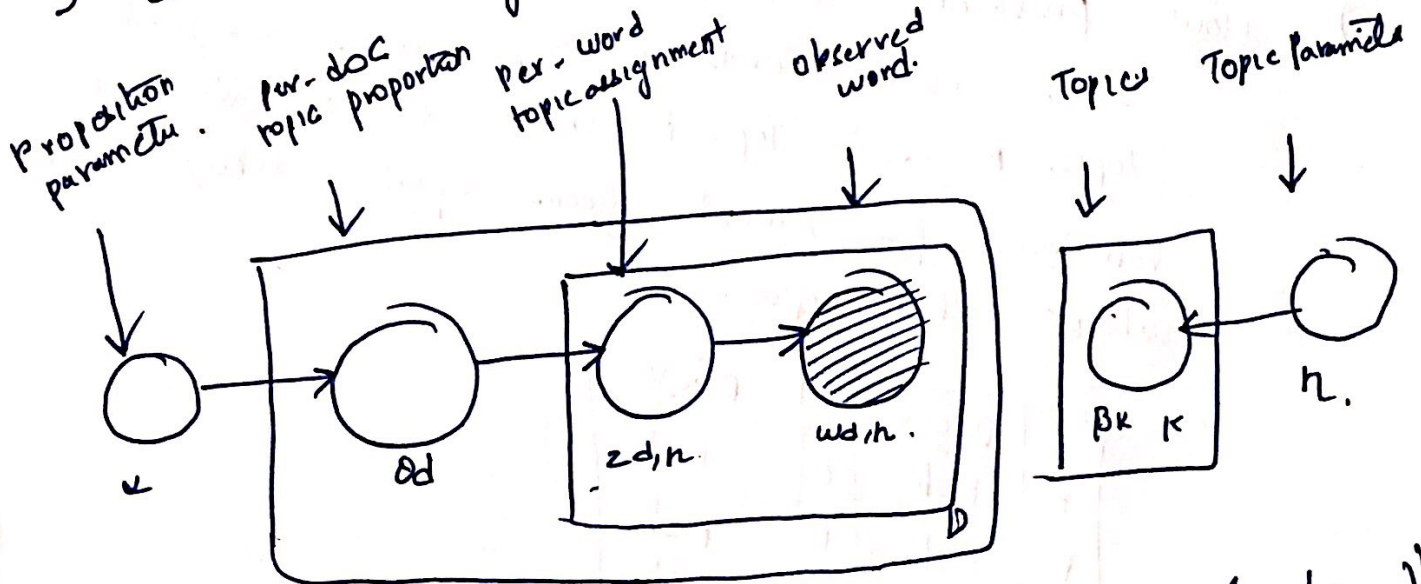


- ⑩ How prevalent are topics in the document?
- Since the words in DocY are assigned to Topic F and Topic P in a 50-50 ratio, the remaining "fish", words seem equally likely to be about either topic.

	Doc X		DocY
F	Fish	2	Fish.
F	Fish	F	fish.
F	cat	F	Milk.
F	cat	P	Kitten
F	Vegetable	P	Kitten

- ⑪
- 1) each topic is a distribution over words
 - 2) each document is a mixture of corpus wide topics
 - 3) each word is drawn from one of these topics
 - 4) we only observe the documents
 - 5) The other structure are hidden variable
 - 6) our goal is to infer the hidden variable i.e. compute their distribution conditioned on documents.
 - 7) Encode assumption (topics, proportions, assignments | documents)
 - 8) Define a factorization of the joint distribution

a) connect to algorithm to compute with data



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

@@ we have to create \$K=3\$ clusters
 Let's choose \$D_2, D_5\$ & \$D_7\$ as initial three seeds

Now we have to calculate euclidean distance from other documents \$D_1, D_3\$ & \$D_7\$
 \$O \rightarrow\$ Online, \$F \rightarrow\$ Festival, \$B \rightarrow\$ Book, \$F \rightarrow\$ Flight, \$D \rightarrow\$ Delhi

$$\begin{aligned} D_1 \text{ to } D_2 &= \sqrt{(O_1 - O_2)^2 + (F_1 - F_2)^2 + (B_1 - B_2)^2 + (T_1 - T_2)^2 + (P_1 - P_2)^2} \\ &= \sqrt{(1-2)^2 + (0-1)^2 + (1-2)^2 + (1-1)^2} \\ &= \sqrt{4} = 2 \end{aligned}$$

$$\begin{aligned} D_1 \text{ to } D_5 &= \sqrt{(1-3)^2 + (0-0)^2 + (1-0)^2 + (0-0)^2 + (1-0)^2} \\ &= \sqrt{4} = 2.6 \end{aligned}$$

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

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

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

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

$$D_3 \text{ to } D_2 = \sqrt{6} = 2.4 \quad D_4 \text{ to } D_2 = \sqrt{1} = 2.8 \quad D_4 \text{ to } D_7 = 0$$

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

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

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

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

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

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

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

$$D_5 \text{ to } D_2 = \sqrt{6} = 2.4$$

$$D_7 \text{ to } D_2 = \sqrt{6} = 2.4$$

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

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

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

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

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

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

$$D_6 \text{ to } D_2 = \sqrt{4} = 2$$

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

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

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

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

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

Documents

D₁

D₂

D₃

D₄

D₂

2.0

0.0

2.4

2.8

D₅

2.6

2.6

3.6

3.0

D₇

2.2

1.7

2.2

2.6

Minds

2.0

0.0

0.2

2.6

clusters

D₂

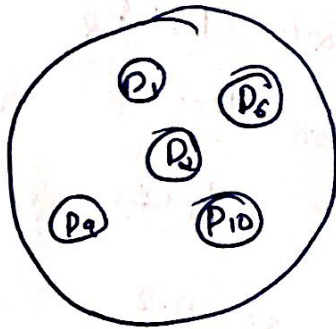
D₂

D₇

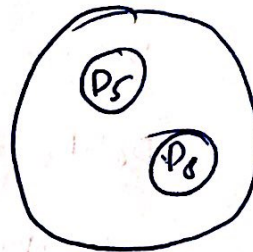
D₂

P ₆	2.6	0.0	2.8	0.0	P ₅
P ₆	2.4	3.9	2.6	2.4	P ₂
P ₇	1.7	2.8	0.0	0.0	P ₇
P ₈	2.6	2.0	2.8	2.0	P ₅
P ₉	2.0	3.0	3.6	2.0	P ₂
P ₁₀	2.2	3.5	2.4	2.2	P ₂

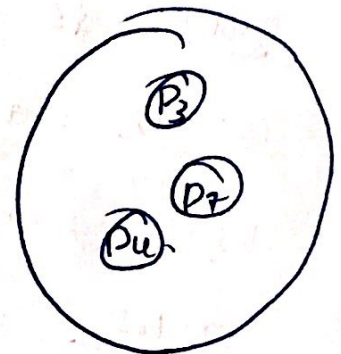
P₂ cluster



P₅ cluster



P₇ cluster



Q6 K Means clustering

Pros

- ① Fast, robust and easier to understand
- ② gives Best result when data set are distinct or well separated from each other.
- ③ It is a great solution for pre-clustering
- ④ works great for spherical clusters

Cons

- ① K-Value is not known and is difficult to predict
- ② There is no unique solution for a certain value since initial partitions can be different
- ③ 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 contribute in given topic

Cons

- ① we have to specify number of topics
- ② LDA's efficiency is pretty low when compared to machine learning algorithms
- ③ LDA cannot capture correlations
- ④ unsupervised (sometimes we need supervision. sentiment analysis)
- ⑤ uses Bow (assumes words are exchangeable)