1) Latent Drichlet Allocation is a probabilistic used to generated the topics. LDA is the iterative model which requires 3 parameters, which are number of topics and dupa prior, knowledge of the dataset

We evaluate the LDA performance using perplexity to Evaluate the LDA model. One document is taken and split in two. the first half is led in LDA to compute the topics composition, from that composition, then, the word distribution is Estimated.

This distribution is then composed with the word distribution is Estimated. This distribution is then composed with the measure of distance is Entrouted. perplexity is often used to select the best number of topics of the LDA model

LDA

Input : words WE Document d

Output : topic assignment Z and counts $n_{d,k}$, $n_{k,w}$ and $n_k$

begin
    Randomly initialized Z and increment Countles for Each interation do

    for i = 0 → N-1 do

        word ← w[i]
        topic ← z[i]

        $n_d$, topic -1 ; $n_{word, topic}$ -1 $n_{topic}$ -1

        for k = 0 → K-1 do

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

        End
        topic ← Sample from $p(z)$
        $z[i]$ ← topic

        $n_d$ → topics 1 ; $n_{word, topic}$ +=1 ; $n_{topic}$ +=1
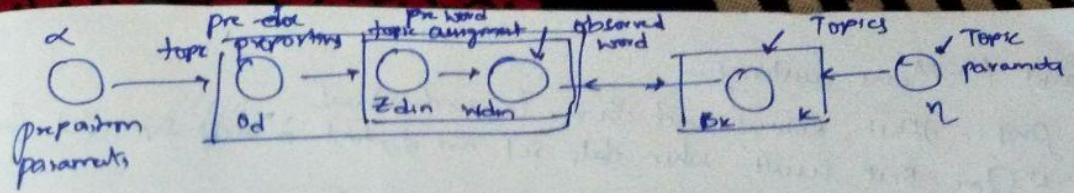
    End
End
return $n_{d,k}$, $n_{k,w}$, $n_k$

1. decide How many topics you need

2. The algorithm will assign Every word to a tempora word
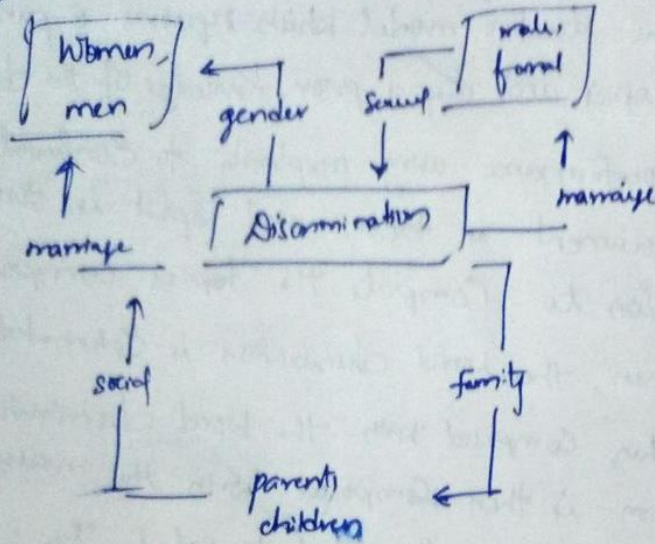
2. will check and update top

$$P(\beta, \theta, z, W) = \left(\prod_{t=1}^{k} P(\beta, I_\eta)\right)\left(\prod_{d=1}^{D} P(z_{d,n}/\theta_d)\right) P(W_{d,n}/\beta: k\, z_{d,n})$$

2. a) K=3   centers $D_2\ b_5\ D_7$

| | |
|---|---|
| $D_1$ to $D_2 = \sqrt{4} = 2$ | $D_3$ to $D_2 = \sqrt{6} = 2.4$ |
| $D_1$ to $D_5 = \sqrt{7} = 2.6$ | $D_3$ to $D_2 = \sqrt{13} = 3.6$ |
| $D_1$ to $D_7 = \sqrt{5} = 2.2$ | $D_3$ to $D_5 = \sqrt{5} = 2.2$ |
| $D_2$ to $D_5 = \sqrt{7} = 2.6$ | $D_3$ to $D_7 = \sqrt{7} = 2.6$ |
| $D_4$ to $D_2 = \sqrt{8} = 2.8$ | $D_5$ to $D_2 = \sqrt{7} = 2.6$ |
| $D_4$ to $D_5 = \sqrt{9} = 3$ | $D_5$ to $D_5 = 0 = 0$ |
| $D_7$ to $D_7 = \sqrt{3} = 1.7$ | $D_5$ to $D_7 = \sqrt{8} = 2.8$ |
| $D_7$ to $D_5 = \sqrt{8} = 2.8$ | $D_7$ to $D_2 = 2.2$ |
| $D_7$ to $D_2 = \sqrt{7} = 2.4$ | $D_9$ to $D_5 = 3.4$ |
| $D_5 , D_5 = \sqrt{6} = 2.2$ | $D_9$ to $D_7 = 3.2$ |
| $D_8$ to $D_7 = \sqrt{5} = 2.2$ | $D_{10}$ to $D_2 = 2.4$ |
| | $D_{10}$ to $D_5 = 2.2$ |

| Doc | $D_2$ | $D_5$ | $D_7$ | Min | clusts |
|---|---|---|---|---|---|
| $D_1$ | 2.0 | 2.6 | 2.2 | 2.0 | $D_2$ |
| $D_2$ | 0 | 2.6 | 2.7 | 0 | $D_2$ |
| $D_3$ | 2.4 | 3.6 | 2.2 | 2.2 | $D_7$ |
| $D_4$ | 2.8 | 5.0 | 2.6 | 2.6 | $D_7$ |
| $D_5$ | 2.6 | 0 | 2.8 | 0 | $D_5$ |
| $D_6$ | 2.4 | 3.9 | 2.6 | 2.4 | $D_2$ |
| $D_7$ | 1.7 | 2.8 | 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 | $D_2$ |

(1c) Since the words in the document Y are assigned to Topic C and Topic P in a 50-50 ratio, the remaining "fish" word Seems Equally likely to be about either topic

| DocX | | DocY | |
|------|------|------|------|
| F | Fish | ? | Fish |
| F | Fish | F | Fish |
| F | Eat | F | Milk |
| F | Eat | P | Kitten |
| F | Vegetably | P | Kitten |

(1d)
1) Each topic is distribution over words.
2) Each document is a mixture of Corpus-wide topics.
3) Each word is drawn from one of those topics.
4) We only observe the documents.
5) The other structure are hidden variables.
6) Our goal is to infer the hidden values i.e Compute their distribution
   Conditional on the document.

7) Encode assumption
8) Define a factorization of the joint distribution

2.b) K-Means clustering :-

Pros :- 1) Fast, Robust and Easure to understand
2) Give Best result when data set are distint or well seperated from
   Each other.
3) It is a great solution for pre-clustery.
4) works great from spherical cluster.

LDA Topic Discovery Model.

Pros: We Can infer Context spread of 8th sentence by a word Count.
① We Can derive the properties that Each word consthtury is given
   topics.

Cons:
1) We have to specify the number of topics.
2) LDA's efficiency is pretty low when compart to m.L. Algo
3) LDA Cannot Compare Co-relation.
4) unsupervised
5) uses BoW (words are Exchangble)