

Latent Dirichlet Allocation is a generative model.

In ML it is used for the task of Topic Modeling. Topic Modeling is clustering of documents into topics which are latent.

Some Preliminaries

Multinomial Distribution

- Generalization of binomial distribution

Suppose we have vector $x \in \mathbb{R}^K$ such that only one coordinate is 1 and rest are 0.

If we denote $P(x_k = 1) = \mu_k$ distribution of x is given as

$$P(x|\mu) = \prod_{k=1}^K \mu_k^{x_k}$$

where $\mu = (\mu_1, \dots, \mu_K)$
 $\mu_k \geq 0 \quad \sum_k \mu_k = 1$

Now consider dataset D of N independent observations x_1, \dots, x_N

Likelihood

$$P(D|\mu) = \prod_{i=1}^N \prod_{k=1}^K \mu_k^{x_{ik}}$$

$$= \prod_{i=1}^N \mu_k^{x_{ik}} \quad (\sum x_{ik})$$

$$= \prod_{k=1}^K \mu_k^{m_k} \quad (m_k = \sum_{i=1}^N x_{ik}) \quad (\text{no. of observations of } x_k = 1)$$

Now the Maximum Likelihood Solution for μ is

$$\mu_k^{ML} = \frac{m_k}{N}$$

We consider the joint distribution of quantities m_1, \dots, m_K conditioned on μ and N

$$\text{Mult}(m_1, m_2, \dots, m_K | \mu, N)$$

$$= \binom{N}{m_1, m_2, \dots, m_K} \prod_{k=1}^K \mu_k^{m_k}$$

where $\binom{N}{m_1, m_2, \dots, m_K} = \frac{N!}{m_1! m_2! \dots m_K!}$

$$\sum_{k=1}^K m_k = N$$

This is the multinomial distribution.

Dirichlet Distribution :-

Conjugate Prior for Multinomial distribution

$$p(\mu|\alpha) \text{ varies as } \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

$$\text{where } 0 \leq \mu_k \leq 1$$

$$\sum_k \mu_k = 1$$

$\alpha_1, \dots, \alpha_K$ are the parameters of the distribution

$$\alpha = (\alpha_1, \dots, \alpha_K)$$

The Dirichlet distribution is given as

$$\text{Dir}(\mu|\alpha) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

$$\alpha_0 = \sum_{k=1}^K \alpha_k$$

$$\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du \quad (\text{extension of factorial function for real})$$

$$\Gamma(1) = 1$$

$$\text{Now } p(\mu|D, \alpha)$$

$$= p(D|\mu) p(\mu|\alpha)$$

$$\propto \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1}$$

$$\text{Thus } p(\mu|D, \alpha) = \text{Dir}(\mu|\alpha + m)$$

$$= \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_0 + m_1) \dots \Gamma(\alpha_0 + m_K)} \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1}$$

$$m = (m_1, \dots, m_K)$$

Topic Modeling & LDA

Given a corpus of documents, we consider that each word in each document has a latent assignment into one of K topics. LDA is standard algo. for the task.

Each document is a mixture of topics.

Each topic is a probability distribution over words.

High Level Idea as Generative Model :-

Each document is generated by choosing a topic

Each word is generated using particular topic

Notations & Algorithm :- (From a blog by Aaron Courville)

D - Documents
V - Vocabulary / Dictionary

We represent a document d with a vector $w_d \in \mathbb{N}^n$ where n_d is number of terms in document d .

$$w_d = (w_{d,1}, w_{d,2}, \dots, w_{d,n_d})$$

where $w_{d,i}$ stands for index of word i in position i of document d .

V.

Let K - number of latent topics in corpus

in corpus

θ_d a vector in Δ^{K-1} i.e. $\theta_d \in [0,1]^K$

Share of topic k in doc. d .

$\beta_k \in \Delta^{V-1}$ vector of term probabilities in topic k

β_{kv} is the probability of observing word v in topic k

Consider the latent variable

$z_{d,n} \in \{1, \dots, K\}$ is topic assigned to n th word in document d .

The Algorithm

1. Draw K vectors $\theta_d \in \Delta^{K-1}$ from $\text{Dir}(\theta_d|\eta)$

2. Draw D vectors $\theta_d \in \Delta^{K-1}$ from $\text{Dir}(\theta_d|\alpha)$

3. Each word $w_{d,n}$ in doc d is generated in two steps:

i. Draw $z_{d,n} \in \{1, \dots, K\}$ according to topic probabilities θ_d .

ii. Draw $w_{d,n}$ using term probabilities $\beta_{z_{d,n}}$.

Now draw θ_d from $\text{Dir}(\theta_d|\alpha)$

Draw θ_d from $\text{Dir}(\theta_d|\alpha)$