

Latent Dirichlet Allocation

Saki

August 21, 2018

Abstract

1 Reference

Blei, D.M., Ng, A.Y. and Jordan, M.I., 2003. Latent dirichlet allocation. Journal of machine Learning research, 3(Jan), pp.993-1022.

2 Model

Generative Process for each word

$$\begin{aligned}\theta &\sim \text{Dir}(\alpha) \\ z_n &\sim \text{Multinomial}(\theta) \\ w_n &\sim p(w_n|z_n, \beta)\end{aligned}$$

Joint Probabilities of θ, z_n, w_n given α, β

$$p(\theta, Z, W|\alpha, \beta) = p(\theta|\alpha) \prod_{n=1}^N p(z_n|\theta)p(w_n|z_n, \beta)$$

3 Intuition

Imagine a big box, which corresponds to a document. We put small boxes, which corresponds to topics, into the big box.

The size of the small boxes are determined by θ (e.g. Sports box takes 50% of the big box, Politics box 20%, and Economics box 30%). ($p(\theta|\alpha)$)

Once the small boxes are fit, we chose one small box($p(z_n|\theta)$), and throw a ball (i.e. a word) into the chosen small box. ($p(w_n|z_n, \beta)$)

This means that we pick a word which is likely to occur when we write about the chosen topic (e.g. 'baseball' in sport topic), and write down the word on the document.

Repeat this process for the number of words in the document.