

Topic Modeling Tutorial for Communications Researchers

Contents

1 Overview	1
2 History of Topic Modeling	1
2.1 Latent Dirichlet Allocation	1
3 Computational implementation of LDA with R	2
References	2

1 Overview

We present a concise summary of a collection of machine learning techniques that, together, are called “topic modeling”. We discuss how these methods may be used in communications research, and we apply topic models to illustrative examples.

With the tremendous rise in computing speed and memory capacity over the last quarter century, researchers working at the interface of quantitative methods and social sciences began to treat written texts as data. While text analysis is still in its infancy, scientists nevertheless have made great progress towards computational dissection and interpretation of written documents. Among the most foundational contributions is the development of probabilistic topic models. We detail below, with limited use of statistical terminology, how these methods work and why they may be useful in communications research. We also provide computer code (in the R statistical language) that implements these methods.

2 History of Topic Modeling

2.1 Latent Dirichlet Allocation

Blei, Ng, & Jordan (2003)

3 Computational implementation of LDA with R

We present below instructions and code for using LDA in the R statistical environment (R Core Team, 2015).

References

Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *The Journal of Machine Learning Research*, 3, 993–1022.

R Core Team. (2015). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. Retrieved from <https://www.R-project.org/>