# "Topic modeling" is WE1S's main method for studying collections of documents at big data scale.

Topic modeling is a leading method of machine learning that discovers "topics" in texts by analyzing the statistical co-occurrence of words. It finds out which words tend to come up together in a document set (and in individual texts) when people discuss something or, as in newspapers, many things. An article about a political election, for example, might frequently use words naming a nation's capital city (e.g., "London") with other words like "Parliament." Co-occurring words suggest "topics."

Topic modeling not only identifies topics but also indicates their weights. (For example, politics might be more important than economics in one newspaper.) It also shows which specific documents participate strongly in a topic (or several topics at once, since an article that mentions "London" might talk partly about politics but also partly about economics because London is also a financial capital). Topic models thus not only help with "distant reading" large collections but also guide researchers to specific documents to "close read" because of their statistical association with a topic.

Topic modeling is an "unsupervised" machine-learning method, meaning that it does not require pre-training of the computer on a gold standard of texts that humans previously analyzed. The specific kind of topic modeling that WE1S uses is Latent Dirichlet Allocation (LDA) as implemented in the MALLET Machine Learning for Language Toolkit.



**Topic model of newspaper articles visualized in Dfr-browser (a topic model interface by Andrew Goldstone)**

---

**Audience for key finding**: Scholarly, Digital Humanists

**Introductions toTopic Modeling**: WE1S bibliography

**Recommended articles**:

* Blei, David M. "Probabilistic Topic Models." (2012)

* Underwood, Ted. "Topic Modeling Made Just Simple Enough" (2012)

* Mohr, John, and Petko Bogdanov, "Topic Models: What They Are and Why They Matter" (2013)