

Beyond LDA

How to integrate Topic Modeling into the Historian's Toolbox

Melanie Althage, M.A., Professur für Digital History, Humboldt-Universität zu Berlin

Ausgangslage

Die Quellen der Gegenwart sind digital



Digitale Quellen legen als maschinell prozessierbare Daten den Einsatz digitaler Methoden nahe.

Dissertation

Titel: „Mining the Historian's Web – Methodenkritische Reflexion quantitativer Verfahren zur Analyse genuin digitaler Quellen am Beispiel der historischen Fachkommunikation“

Ziel:

- Untersuchung der Adaptierbarkeit etablierter Textanalysemethoden der Digital Humanities und Computerwissenschaften für historische Quellen und Forschungsfragen
- Ausarbeitung eines Kriterienkatalogs zu ihrer Integration in den Werkzeugkasten der Historiker:innen

Bedarf

Methodenkritik

- Prüfung: Welches Verfahren passt zu meiner Fragestellung?
- Entwicklungskontext: Ursprüngliche Erkenntnisinteressen und Ziele
- Formalisierung welcher theoretisch-methodologischen Annahmen?
- Auf Basis welcher Art von Daten entwickelt und optimiert?
- Welche Parameter gibt es? Wie wirken sie sich auf den Output aus?

Beispiel aus dem digitalen Werkzeugkasten



Zur (diachronen) Untersuchung von Diskursen, Motiven, thematischen Trends, u.v.m. in den digitalen Geistes- und Geschichtswissenschaften etabliert:

Topic Modeling

... ist eine auf **Wahrscheinlichkeitsrechnung** basierende Methode aus dem Bereich des unüberwachten maschinellen Lernens zur **Gruppierung von Dokumenten** einer Textsammlung anhand ihrer **gemeinsamen Sprachgebrauchsmuster**.

Entwicklungskontext: Information Retrieval (IR) und Natural Language Processing (NLP)

Ursprüngliches Ziel: Beschreibung und Exploration umfangreicher Datenkollektionen zur Optimierung von Suchheuristiken

Intuition gemäß distributioneller Semantik: Über die gemeinsame Vorkommenshäufigkeit lexikalischer Einheiten (z.B. Wörter) in einem bestimmten Kontext (z.B. Dokumenten) lassen sich Aussagen über den Inhalt der Dokumente formulieren.

Der Quasi-Standard für Topic Modeling:



Latent Dirichlet Allocation (Blei et al. 2003)
Leicht zu implementieren dank gebrauchsfertiger Tools und niedrigschwelliger Tutorials



Herausforderungen für geschichtswissenschaftliche Anwendungsfälle

Die **Kontextabhängigkeit** und insb. **Historizität** der Daten wird bei der Modellierung nicht berücksichtigt:

- Topics sind über das gesamte Korpus statisch
- keine Berücksichtigung von Inter-Topic-Beziehungen
- Indifferenz ggü. Ordnung des Korpus
- keine Berücksichtigung externer Informationen



Populäre Verfahren wie LDA stammen aus Forschungsbereichen mit je eigenen theoretisch-methodologischen Annahmen und Erkenntnisinteressen, die im Kontrast zu geschichtswissenschaftlichen Forschungstraditionen stehen können.

Problem

Work in Progress:

Guidelines und Best Practices zur Orientierung: Welche Verfahren zur Topic-Modellierung gibt es über LDA hinaus? Für welche Art von Daten und Fragestellungen sind sie geeignet?

