

Big Data Research as Interaction between Topic Models and Expert Data: A New Approach to Capturing National Online Debates

Cäcilia Zirn, Eike Mark Rinke, Charlotte Löb, & Hartmut Wessler



TASK

Dataset

- 1 year of scraping
- 101 sources
- Leading national blogs & news websites
- 6 countries
 - Australia, Germany, Lebanon, Turkey, Switzerland, USA
- 4 languages
 - English, German, Arabic, Turkish
- 1.6 million articles

Find the needle in the haystack!

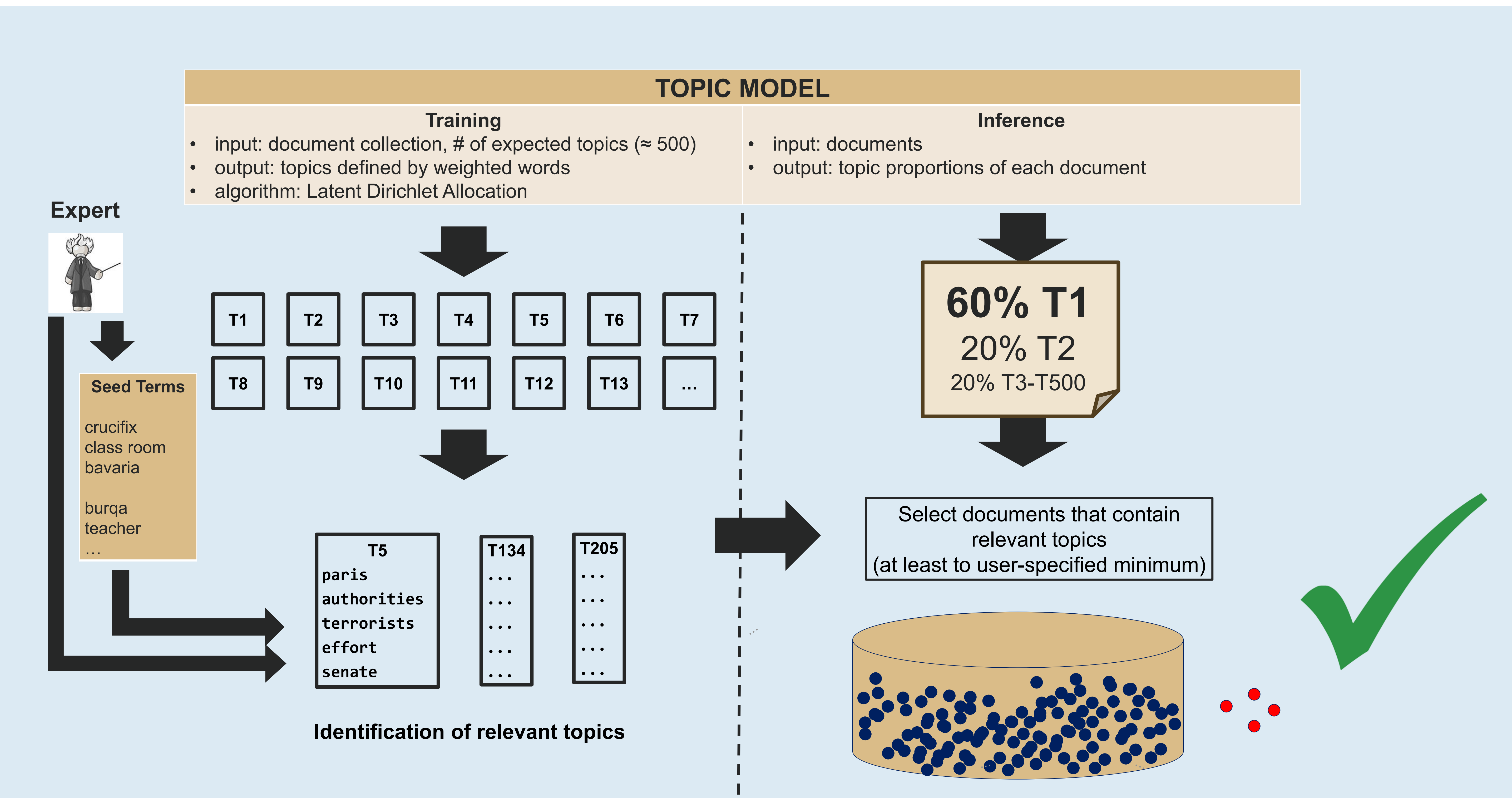
General goal: More insightful topic identification

- Using expert domain knowledge for better automated identification of documents from prespecified broad topic domain ("meta-topic")
- For classification of documents by meta-topic
- For cross-media and cross-national comparison of topic structures within a specific meta-topic

Specific application: Expert knowledge for improved classification of documents by topic

- Finding articles that belong to at least one topic that are part of meta-topic "role of religion in public life"

TOPIC MODELING APPROACH



WHY DID YOU NOT USE ...

Machine Learning	Information Retrieval	Clustering	Seed-based LDA	Crowd Sourcing
<ul style="list-style-type: none">no training data availableirrelevant classes unknownclasses (relevant vs. irrelevant) unbalanced	<ul style="list-style-type: none">expert-generated list of relevant items probably not completecollecting all possible search terms not feasible (synonyms, paraphrases)	<ul style="list-style-type: none">high cognitive effort required to comprehend clusterswords can be assigned to a single cluster only (though they might be relevant for several topics)only single meaning per word can be used	<ul style="list-style-type: none">presumptions about expected issues bias topic creation and might limitate resultsrequires comprehensive issue (seed) list	<ul style="list-style-type: none">coders must follow complex instructions and decide carefullydifficult to scale to millions of documentsrequires preexisting domain knowledgedifficult to verify coding quality with extremely few positive codes