

LDA and Beyond: Topic Models in the Social Sciences

Brandon Stewart¹

Princeton University

June 21, 2017

¹My sincere thanks to my many collaborators and particularly Justin Grimmer, Molly Roberts and Dustin Tingley from whom many of these slides are derived.

Papers

- Overview of Text Analysis:

- ▶ “Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts” (*Political Analysis*, 2013) with Grimmer

- Structural Topic Model:

- ▶ Structural Topic Models for Open-Ended Survey Responses (*American Journal of Political Science*, 2014) with Roberts, Tingley et al.
- ▶ Computer Assisted Text Analysis for Comparative Politics (*Political Analysis* 2015), with Lucas et al.
- ▶ A Model of Text for Experimentation in the Social Sciences (2016) with Roberts and Airolidi

Copies at BrandonStewart.org

Big Data, Big Analytics

Big Data, Big Analytics

- Massive increase in **unstructured** text due to:

Big Data, Big Analytics

- Massive increase in **unstructured** text due to:
 - ▶ new social structures (the internet, email)

Big Data, Big Analytics

- Massive increase in **unstructured** text due to:
 - ▶ new social structures (the internet, email)
 - ▶ new/improved data collection (wiki surveys, survey experiments)

Big Data, Big Analytics

- Massive increase in **unstructured** text due to:
 - ▶ new social structures (the internet, email)
 - ▶ new/improved data collection (wiki surveys, survey experiments)
 - ▶ digitization efforts (government documents, Google Books)

Big Data, Big Analytics

- Massive increase in **unstructured** text due to:
 - ▶ new social structures (the internet, email)
 - ▶ new/improved data collection (wiki surveys, survey experiments)
 - ▶ digitization efforts (government documents, Google Books)
- Communities leave **digitized** footprints

Big Data, Big Analytics

- Massive increase in **unstructured** text due to:
 - ▶ new social structures (the internet, email)
 - ▶ new/improved data collection (wiki surveys, survey experiments)
 - ▶ digitization efforts (government documents, Google Books)
- Communities leave **digitized** footprints
- Tools to **analyze** text advancing in parallel

Big Data, Big Analytics

- Massive increase in **unstructured** text due to:
 - ▶ new social structures (the internet, email)
 - ▶ new/improved data collection (wiki surveys, survey experiments)
 - ▶ digitization efforts (government documents, Google Books)
- Communities leave **digitized** footprints
- Tools to **analyze** text advancing in parallel
 - ▶ text by itself is useless

Big Data, Big Analytics

- Massive increase in **unstructured** text due to:
 - ▶ new social structures (the internet, email)
 - ▶ new/improved data collection (wiki surveys, survey experiments)
 - ▶ digitization efforts (government documents, Google Books)
- Communities leave **digitized** footprints
- Tools to **analyze** text advancing in parallel
 - ▶ text by itself is useless
 - ▶ importing methods from many different fields

Big Data, Big Analytics

- Massive increase in **unstructured** text due to:
 - ▶ new social structures (the internet, email)
 - ▶ new/improved data collection (wiki surveys, survey experiments)
 - ▶ digitization efforts (government documents, Google Books)
- Communities leave **digitized** footprints
- Tools to **analyze** text advancing in parallel
 - ▶ text by itself is useless
 - ▶ importing methods from many different fields
 - ▶ new analysis techniques can even drive new data availability

Different Methods for Different Goals

Different Methods for Different Goals

- **Supervised**: pursuing a known goal

Different Methods for Different Goals

- **Supervised**: pursuing a known goal
 - ▶ human annotates a subset of documents

Different Methods for Different Goals

- **Supervised**: pursuing a known goal
 - ▶ human annotates a subset of documents
 - ▶ algorithm annotates the rest

Different Methods for Different Goals

- **Supervised**: pursuing a known goal
 - ▶ human annotates a subset of documents
 - ▶ algorithm annotates the rest
 - ▶ usually associated with quantitative research

Different Methods for Different Goals

- **Supervised**: pursuing a known goal
 - ▶ human annotates a subset of documents
 - ▶ algorithm annotates the rest
 - ▶ usually associated with quantitative research
- **Unsupervised**: goal is to learn the goal

Different Methods for Different Goals

- **Supervised**: pursuing a known goal
 - ▶ human annotates a subset of documents
 - ▶ algorithm annotates the rest
 - ▶ usually associated with quantitative research
- **Unsupervised**: goal is to learn the goal
 - ▶ algorithm discovers themes/patterns in the texts

Different Methods for Different Goals

- **Supervised**: pursuing a known goal
 - ▶ human annotates a subset of documents
 - ▶ algorithm annotates the rest
 - ▶ usually associated with quantitative research
- **Unsupervised**: goal is to learn the goal
 - ▶ algorithm discovers themes/patterns in the texts
 - ▶ human interprets the results

Different Methods for Different Goals

- **Supervised**: pursuing a known goal
 - ▶ human annotates a subset of documents
 - ▶ algorithm annotates the rest
 - ▶ usually associated with quantitative research
- **Unsupervised**: goal is to learn the goal
 - ▶ algorithm discovers themes/patterns in the texts
 - ▶ human interprets the results
 - ▶ usually associated with qualitative research

Different Methods for Different Goals

- **Supervised**: pursuing a known goal
 - ▶ human annotates a subset of documents
 - ▶ algorithm annotates the rest
 - ▶ usually associated with quantitative research
- **Unsupervised**: goal is to learn the goal
 - ▶ algorithm discovers themes/patterns in the texts
 - ▶ human interprets the results
 - ▶ usually associated with qualitative research
- Both strategies **amplify** human effort, each in different ways.

Different Methods for Different Goals

- Supervised: pursuing a known goal
 - ▶ human annotates a subset of documents
 - ▶ algorithm annotates the rest
 - ▶ usually associated with quantitative research
- **Unsupervised**: goal is to learn the goal
 - ▶ algorithm discovers themes/patterns in the texts
 - ▶ human interprets the results
 - ▶ usually associated with qualitative research
- Both strategies amplify human effort, each in different ways.

Topic Models in Social Science

Topic Models in Social Science

- Core methods developed in computer science and statistics

Topic Models in Social Science

- Core methods developed in computer science and statistics
 - ▶ used as a way to summarize **unstructured** text

Topic Models in Social Science

- Core methods developed in computer science and statistics
 - ▶ used as a way to summarize **unstructured** text
 - ▶ use words **within** document to infer its subject

Topic Models in Social Science

- Core methods developed in computer science and statistics
 - ▶ used as a way to summarize **unstructured** text
 - ▶ use words **within** document to infer its subject
 - ▶ introduced as a form of **dimension reduction**

Topic Models in Social Science

- Core methods developed in computer science and statistics
 - ▶ used as a way to summarize **unstructured** text
 - ▶ use words **within** document to infer its subject
 - ▶ introduced as a form of **dimension reduction**
- A theory of use in the social sciences

Topic Models in Social Science

- Core methods developed in computer science and statistics
 - ▶ used as a way to summarize **unstructured** text
 - ▶ use words **within** document to infer its subject
 - ▶ introduced as a form of **dimension reduction**
- A theory of use in the social sciences
 - ▶ social scientists want to use topics as a form of **measurement**

Topic Models in Social Science

- Core methods developed in computer science and statistics
 - ▶ used as a way to summarize **unstructured** text
 - ▶ use words **within** document to infer its subject
 - ▶ introduced as a form of **dimension reduction**
- A theory of use in the social sciences
 - ▶ social scientists want to use topics as a form of **measurement**
 - ▶ we are often interest in how observed covariates drive **trends** in language

Topic Models in Social Science

- Core methods developed in computer science and statistics
 - ▶ used as a way to summarize **unstructured** text
 - ▶ use words **within** document to infer its subject
 - ▶ introduced as a form of **dimension reduction**
- A theory of use in the social sciences
 - ▶ social scientists want to use topics as a form of **measurement**
 - ▶ we are often interest in how observed covariates drive **trends** in language
 - ▶ we want to tell a story not just about what, but **how** and **why**

Topic Models in Social Science

- Core methods developed in computer science and statistics
 - ▶ used as a way to summarize **unstructured** text
 - ▶ use words **within** document to infer its subject
 - ▶ introduced as a form of **dimension reduction**
- A theory of use in the social sciences
 - ▶ social scientists want to use topics as a form of **measurement**
 - ▶ we are often interest in how observed covariates drive **trends** in language
 - ▶ we want to tell a story not just about what, but **how** and **why**
- Different focus in social science brings different concerns and an emphasis on **validation**

Topic Models in Social Science

- Core methods developed in computer science and statistics
 - ▶ used as a way to summarize **unstructured** text
 - ▶ use words **within** document to infer its subject
 - ▶ introduced as a form of **dimension reduction**
- A theory of use in the social sciences
 - ▶ social scientists want to use topics as a form of **measurement**
 - ▶ we are often interest in how observed covariates drive **trends** in language
 - ▶ we want to tell a story not just about what, but **how** and **why**
- Different focus in social science brings different concerns and an emphasis on **validation**

Some Example Questions

Some Example Questions

- How do senators present their work to the public? What explains variation in representational style? (Grimmer 2013)

Some Example Questions

- How do senators present their work to the public? What explains variation in representational style? (Grimmer 2013)
- Did electoral reform change the portfolio of issues addressed by politicians in Japan? (Catalinac 2016)

Some Example Questions

- How do senators present their work to the public? What explains variation in representational style? (Grimmer 2013)
- Did electoral reform change the portfolio of issues addressed by politicians in Japan? (Catalinac 2016)
- What are the propaganda strategies of the Chinese government? (Roberts and Stewart)

Some Example Questions

- How do senators present their work to the public? What explains variation in representational style? (Grimmer 2013)
- Did electoral reform change the portfolio of issues addressed by politicians in Japan? (Catalinac 2016)
- What are the propaganda strategies of the Chinese government? (Roberts and Stewart)
- What types of actions do countries take towards each other? (O'Connor, Stewart and Smith 2013)

Some Example Questions

- How do senators present their work to the public? What explains variation in representational style? (Grimmer 2013)
- Did electoral reform change the portfolio of issues addressed by politicians in Japan? (Catalinac 2016)
- What are the propaganda strategies of the Chinese government? (Roberts and Stewart)
- What types of actions do countries take towards each other? (O'Connor, Stewart and Smith 2013)
- How do Muslim clerics supporting violent Jihad differ from those who do not in terms of writing in fatwas? (Nielsen, Forthcoming)

Some Example Questions

- How do senators present their work to the public? What explains variation in representational style? (Grimmer 2013)
- Did electoral reform change the portfolio of issues addressed by politicians in Japan? (Catalinac 2016)
- What are the propaganda strategies of the Chinese government? (Roberts and Stewart)
- What types of actions do countries take towards each other? (O'Connor, Stewart and Smith 2013)
- How do Muslim clerics supporting violent Jihad differ from those who do not in terms of writing in fatwas? (Nielsen, Forthcoming)
- Do presidential candidates move to the center after the convention? (Gross et al 2013)

Some Example Questions

- How do senators present their work to the public? What explains variation in representational style? (Grimmer 2013)
- Did electoral reform change the portfolio of issues addressed by politicians in Japan? (Catalinac 2016)
- What are the propaganda strategies of the Chinese government? (Roberts and Stewart)
- What types of actions do countries take towards each other? (O'Connor, Stewart and Smith 2013)
- How do Muslim clerics supporting violent Jihad differ from those who do not in terms of writing in fatwas? (Nielsen, Forthcoming)
- Do presidential candidates move to the center after the convention? (Gross et al 2013)
- Do the proposals of NPC deputies in China represent the interests of their constituents? (Truex 2016)

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But
Some are Useful

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But Some are Useful

- Data generation process for text \rightsquigarrow unknown

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But Some are Useful

- Data generation process for text \rightsquigarrow unknown
- Complexity of language:

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But Some are Useful

- Data generation process for text \rightsquigarrow unknown
- Complexity of language:
 - Time flies like an arrow

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But Some are Useful

- Data generation process for text \rightsquigarrow unknown
- Complexity of language:
 - Time flies like an arrow, fruit flies like a banana

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But Some are Useful

- Data generation process for text \rightsquigarrow unknown
- Complexity of language:
 - Time flies like an arrow, fruit flies like a banana
 - Make peace, not war

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But Some are Useful

- Data generation process for text \rightsquigarrow unknown
- Complexity of language:
 - Time flies like an arrow, fruit flies like a banana
 - Make peace, not war , Make war not peace

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But Some are Useful

- Data generation process for text \rightsquigarrow unknown
- Complexity of language:
 - Time flies like an arrow, fruit flies like a banana
 - Make peace, not war , Make war not peace
 - “Years from now, you’ll look back and you’ll say that this was the moment, this was the place where America remembered what it means to hope. ”

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But Some are Useful

- Data generation process for text \rightsquigarrow unknown
- Complexity of language:
 - Time flies like an arrow, fruit flies like a banana
 - Make peace, not war , Make war not peace
 - “Years from now, you’ll look back and you’ll say that this was the moment, this was the place where America remembered what it means to hope. ”
- Models **necessarily** fail to capture language \rightsquigarrow useful for specific tasks

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 1: All Quantitative Models of Language are Wrong—But Some are Useful

- Data generation process for text \rightsquigarrow unknown
- Complexity of language:
 - Time flies like an arrow, fruit flies like a banana
 - Make peace, not war , Make war not peace
 - “Years from now, you’ll look back and you’ll say that this was the moment, this was the place where America remembered what it means to hope. ”
- Models **necessarily** fail to capture language \rightsquigarrow useful for specific tasks
- **Validation** \rightsquigarrow demonstrate methods perform task

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 2: Quantitative Methods Augment Humans, Not Replace Them

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 2: Quantitative Methods Augment Humans, Not Replace Them

- Computer-Assisted Reading

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 2: Quantitative Methods Augment Humans, Not Replace Them

- Computer-Assisted Reading
- Quantitative methods organize, direct, and suggest

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 2: Quantitative Methods Augment Humans, Not Replace Them

- **Computer-Assisted** Reading
- Quantitative methods organize, direct, and suggest
- Humans: read and interpret

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 3: There is no Globally Best Method for Automated Text Analysis

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 3: There is no Globally Best Method for Automated Text Analysis

- Supervised methods \rightsquigarrow known categories

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 3: There is no Globally Best Method for Automated Text Analysis

- Supervised methods \rightsquigarrow known categories
- Unsupervised methods \rightsquigarrow discover categories

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 3: There is no Globally Best Method for Automated Text Analysis

- Supervised methods \rightsquigarrow known categories
- Unsupervised methods \rightsquigarrow discover categories
- Debate \rightsquigarrow acknowledge differences, resolved

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 4: Validate, Validate, Validate

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 4: Validate, Validate, Validate

- Quantitative methods \rightsquigarrow variable performance across tasks

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 4: Validate, Validate, Validate

- Quantitative methods \rightsquigarrow variable performance across tasks
- Few theorems to guarantee performance

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 4: Validate, Validate, Validate

- Quantitative methods \rightsquigarrow variable performance across tasks
- Few theorems to guarantee performance
- Apply methods \rightsquigarrow validate

Four Principles of Automated Text Analysis (Grimmer and Stewart)

Principle 4: Validate, Validate, Validate

- Quantitative methods \rightsquigarrow variable performance across tasks
- Few theorems to guarantee performance
- Apply methods \rightsquigarrow validate
- Avoid: blind application of methods

- 1 Introduction
- 2 Four Principles
- 3 Preprocessing
- 4 Latent Dirichlet Allocation
- 5 Structured Topic Models
- 6 Structural Topic Models
- 7 Sample Applications
- 8 Conclusion

- 1 Introduction
- 2 Four Principles
- 3 Preprocessing**
- 4 Latent Dirichlet Allocation
- 5 Structured Topic Models
- 6 Structural Topic Models
- 7 Sample Applications
- 8 Conclusion

Bag of Words

Bag of Words

- Throughout we will use a representation called **bag of words** because we will discard word order.

Bag of Words

- Throughout we will use a representation called **bag of words** because we will discard word order.
- Generally instantiated (at least conceptually) as a **document-term matrix**.

Bag of Words

- Throughout we will use a representation called **bag of words** because we will discard word order.
- Generally instantiated (at least conceptually) as a **document-term matrix**.
- This representation is good at capturing **subject matter** of documents but not nuance.

Bag of Words

- Throughout we will use a representation called **bag of words** because we will discard word order.
- Generally instantiated (at least conceptually) as a **document-term matrix**.
- This representation is good at capturing **subject matter** of documents but not nuance.
- We will breeze through this but these choices are **consequential** (see for example Denny and Spirling 2017, Schofield and Mimno 2016)

Preprocessing

“Political power grows out of the barrel of a gun” - Mao

Preprocessing

“Political power grows out of the barrel of a gun” - Mao

Compound Words: With substantive justification, words can be combined or split to improve inference.

Preprocessing

“Political power grows out of the barrel of a gun” - Mao

Compound Words: An analyst may want to combine words into a single term that can be analyzed.

Preprocessing

“Political power grows out of the **barrel of a gun**” - Mao

Compound Words: An analyst may want to combine words into a single term that can be analyzed.

Preprocessing

[Political], [power], [grows], [out], [of], [the], [barrel of a gun]

Compound Words: An analyst may want to combine words into a single term that can be analyzed.

Preprocessing

[Political], [power], [grows], [out], [of], [the], [barrel of a gun]

Stopword Removal: Removing terms that are not related to what the author is studying from the text.

Preprocessing

[Political], [power], [grows], [out], [of], [the], [barrel of a gun]

Stopword Removal: Removing terms that are not related to what the author is studying from the text.

Preprocessing

[Political], [power], [grows], [out], [barrel of a gun]

Stopword Removal: Removing terms that are not related to what the author is studying from the text.

Preprocessing

[Political], [power], [grows], [out], [barrel of a gun]

Stemming: Takes the ends off conjugated verbs or plural nouns, leaving just the “stem.”

Preprocessing

[Politi**cal**], [power], [grow**s**], [out], [barrel of a gun]

Stemming: Takes the ends off conjugated verbs or plural nouns, leaving just the “stem.”

Preprocessing

[Polit], [power], [grow], [out], [barrel of a gun]

Stemming: Takes the ends off conjugated verbs or plural nouns, leaving just the “stem.”

Preprocessing

Finally, we can turn tokens and documents into a “document-term matrix.”

Imagine we have a second document in addition to the Mao quote, which tokenizes as follows.

Document #1: [polit], [power], [grow], [out], [barrel of a gun]

Document #2: [wessi], [compar], [polit], [wessi]

Output: Term-Document Matrix

	<i>Doc1</i>	<i>Doc2</i>
<i>polit</i>	1	1
<i>power</i>	1	0
<i>grow</i>	1	0
<i>out</i>	1	0
<i>barrel of a gun</i>	1	0
<i>wessi</i>	0	2
<i>compar</i>	0	1

- 1 Introduction
- 2 Four Principles
- 3 Preprocessing
- 4 Latent Dirichlet Allocation
- 5 Structured Topic Models
- 6 Structural Topic Models
- 7 Sample Applications
- 8 Conclusion

- 1 Introduction
- 2 Four Principles
- 3 Preprocessing
- 4 Latent Dirichlet Allocation**
- 5 Structured Topic Models
- 6 Structural Topic Models
- 7 Sample Applications
- 8 Conclusion

Latent Dirichlet Allocation

- Idea: documents exhibit each topic in some proportion. This is an **admixture**.

Latent Dirichlet Allocation

- Idea: documents exhibit each topic in some proportion. This is an **admixture**.
- Each document is a mixture over topics. Each topic is a mixture over words.

Latent Dirichlet Allocation

- Idea: documents exhibit each topic in some proportion. This is an **admixture**.
- Each document is a mixture over topics. Each topic is a mixture over words.
- Latent Dirichlet Allocation estimates:

Latent Dirichlet Allocation

- Idea: documents exhibit each topic in some proportion. This is an **admixture**.
- Each document is a mixture over topics. Each topic is a mixture over words.
- Latent Dirichlet Allocation estimates:
 - ▶ The **distribution over words** for each topic.

Latent Dirichlet Allocation

- Idea: documents exhibit each topic in some proportion. This is an **admixture**.
- Each document is a mixture over topics. Each topic is a mixture over words.
- Latent Dirichlet Allocation estimates:
 - ▶ The **distribution over words** for each topic.
 - ▶ The **proportion of a document in each topic**, for each document.

Latent Dirichlet Allocation

- Idea: documents exhibit each topic in some proportion. This is an **admixture**.
- Each document is a mixture over topics. Each topic is a mixture over words.
- Latent Dirichlet Allocation estimates:
 - ▶ The **distribution over words** for each topic.
 - ▶ The **proportion of a document in each topic**, for each document.

Latent Dirichlet Allocation

- Idea: documents exhibit each topic in some proportion. This is an **admixture**.
- Each document is a mixture over topics. Each topic is a mixture over words.
- Latent Dirichlet Allocation estimates:
 - ▶ The **distribution over words** for each topic.
 - ▶ The **proportion of a document in each topic**, for each document.

Maintained assumptions: Bag of words/fix number of topics ex ante.

What this means in pictures

Say you have
a lot of people.

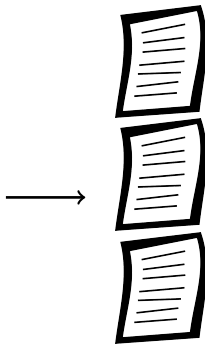


What this means in pictures

Say you have
a lot of people.



Each writes
some texts



What this means in pictures

Say you have
a lot of people.



Each writes
some texts



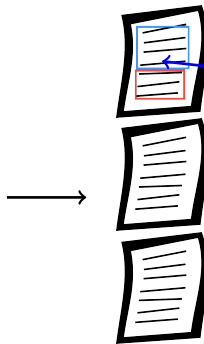
that discuss a few
different topics

What this means in pictures

Say you have
a lot of people.



Each writes
some texts



that discuss a few
different topics

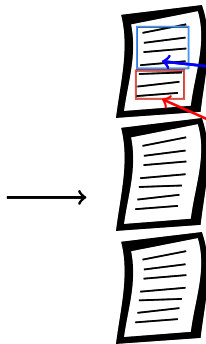
Topic 1

What this means in pictures

Say you have
a lot of people.



Each writes
some texts



that discuss a few
different topics

Topic 1

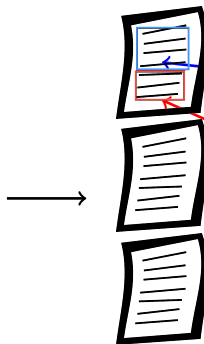
Topic 2

What this means in pictures

Say you have
a lot of people.



Each writes
some texts



that discuss a few
different topics

Topic 1

Topic 2

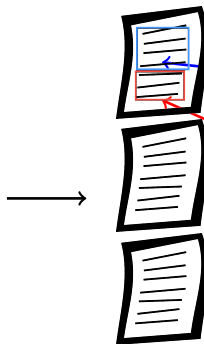
The Latent Dirichlet Allocation estimates:

What this means in pictures

Say you have
a lot of people.



Each writes
some texts



that discuss a few
different topics

Topic 1

Topic 2

The Latent Dirichlet Allocation estimates:

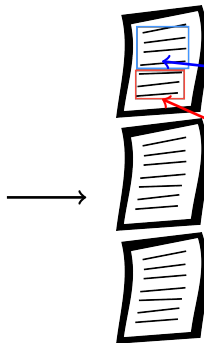
- ① The topics- each is a distribution over words

What this means in pictures

Say you have
a lot of people.



Each writes
some texts



that discuss a few
different topics

Politics

congress, nations,
power, votes, agree-
ment, bargaining

Statistics

estimator, data,
analysis, variance,
model, inference

The Latent Dirichlet Allocation estimates:

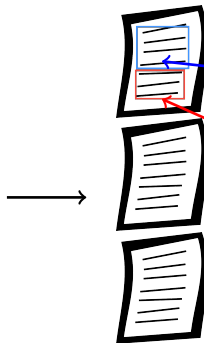
- ① The topics- each is a distribution over words

What this means in pictures

Say you have
a lot of people.



Each writes
some texts



that discuss a few
different topics

Politics

congress, nations,
power, votes, agree-
ment, bargaining

Statistics

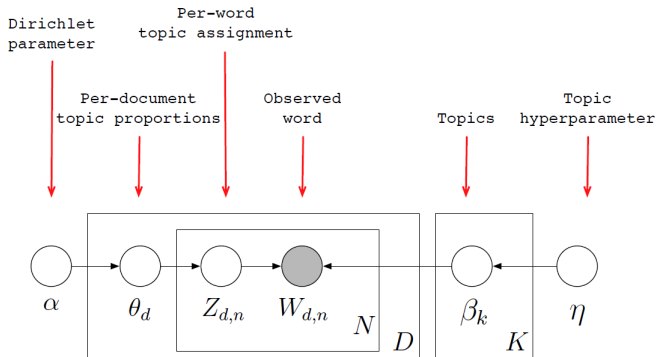
estimator, data,
analysis, variance,
model, inference

The Latent Dirichlet Allocation estimates:

- 1 The topics- each is a distribution over words
- 2 The proportion of each document in each topic

This is a Bayesian Model

Figure: Plate Notation of Latent Dirichlet Allocation



Graphic from David Blei's Website

LDA as a Bayesian Model

$$\beta_k | \eta \sim \text{Dirichlet}(\eta)$$

$$\theta_i | \alpha \sim \text{Dirichlet}(\alpha)$$

$$z_{im} | \theta_i \sim \text{Multinomial}(1, \theta_i)$$

$$w_{im} | \beta_k, z_{imk} = 1 \sim \text{Multinomial}(1, \beta_k)$$

LDA as a Bayesian Model

Unigram Model_{*k*} \sim Dirichlet(η)

Doc. Prop_{*i*} \sim Dirichlet(**Pop. Proportion**)

Word Topic_{*im*} \sim Multinomial(1, **Doc. Prop**_{*i*})

Word_{*im*} \sim Multinomial(1, **Unigram Model**_{*k*})

“Vanilla” Latent Dirichlet Allocation

1) Task:

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

$$f(\mathbf{W}, \boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\alpha})$$

Where:

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

$$f(\mathbf{W}, \beta, \Theta, \alpha)$$

Where:

- $\Theta = N \times K$ matrix with row $\theta_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{iK}) \rightsquigarrow$
proportion of a document allocated to each topic

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

$$f(W, \beta, \Theta, \alpha)$$

Where:

- $\Theta = N \times K$ matrix with row $\theta_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{iK}) \rightsquigarrow$ proportion of a document allocated to each topic
- $\beta = K \times J$ matrix, with row $\beta_k = (\beta_{1k}, \beta_{2k}, \dots, \beta_{kJ}) \rightsquigarrow$ **topics**

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

$$f(W, \beta, \Theta, \alpha)$$

Where:

- $\Theta = N \times K$ matrix with row $\theta_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{iK}) \rightsquigarrow$ proportion of a document allocated to each topic
- $\beta = K \times J$ matrix, with row $\beta_k = (\beta_{1k}, \beta_{2k}, \dots, \beta_{kJ}) \rightsquigarrow$ **topics**
- $\alpha = K$ element long vector, population prior for Θ .

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

$$f(\mathbf{W}, \boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\alpha})$$

Where:

- $\boldsymbol{\Theta} = N \times K$ matrix with row $\boldsymbol{\theta}_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{iK}) \rightsquigarrow$ proportion of a document allocated to each topic
- $\boldsymbol{\beta} = K \times J$ matrix, with row $\boldsymbol{\beta}_k = (\beta_{1k}, \beta_{2k}, \dots, \beta_{kJ}) \rightsquigarrow$ **topics**
- $\boldsymbol{\alpha} = K$ element long vector, population prior for $\boldsymbol{\Theta}$.

3) Optimization

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

$$f(\mathbf{W}, \boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\alpha})$$

Where:

- $\boldsymbol{\Theta} = N \times K$ matrix with row $\boldsymbol{\theta}_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{iK}) \rightsquigarrow$ proportion of a document allocated to each topic
- $\boldsymbol{\beta} = K \times J$ matrix, with row $\boldsymbol{\beta}_k = (\beta_{1k}, \beta_{2k}, \dots, \beta_{kJ}) \rightsquigarrow$ **topics**
- $\boldsymbol{\alpha} = K$ element long vector, population prior for $\boldsymbol{\Theta}$.

3) Optimization

- Variational Inference \rightsquigarrow deterministic approximation

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

$$f(\mathbf{W}, \boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\alpha})$$

Where:

- $\boldsymbol{\Theta} = N \times K$ matrix with row $\boldsymbol{\theta}_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{iK}) \rightsquigarrow$ proportion of a document allocated to each topic
- $\boldsymbol{\beta} = K \times J$ matrix, with row $\boldsymbol{\beta}_k = (\beta_{1k}, \beta_{2k}, \dots, \beta_{kJ}) \rightsquigarrow$ **topics**
- $\boldsymbol{\alpha} = K$ element long vector, population prior for $\boldsymbol{\Theta}$.

3) Optimization

- Variational Inference \rightsquigarrow deterministic approximation
- Collapsed Gibbs Sampling \rightsquigarrow MCMC algorithm

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

$$f(\mathbf{W}, \boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\alpha})$$

Where:

- $\boldsymbol{\Theta} = N \times K$ matrix with row $\boldsymbol{\theta}_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{iK}) \rightsquigarrow$ proportion of a document allocated to each topic
- $\boldsymbol{\beta} = K \times J$ matrix, with row $\boldsymbol{\beta}_k = (\beta_{1k}, \beta_{2k}, \dots, \beta_{kJ}) \rightsquigarrow$ **topics**
- $\boldsymbol{\alpha} = K$ element long vector, population prior for $\boldsymbol{\Theta}$.

3) Optimization

- Variational Inference \rightsquigarrow deterministic approximation
- Collapsed Gibbs Sampling \rightsquigarrow MCMC algorithm
- Spectral/Factorization Methods

“Vanilla” Latent Dirichlet Allocation

1) Task:

- Discover thematic content of documents
- Quickly explore documents

2) Objective Function

$$f(\mathbf{W}, \boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\alpha})$$

Where:

- $\boldsymbol{\Theta} = N \times K$ matrix with row $\boldsymbol{\theta}_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{iK}) \rightsquigarrow$ proportion of a document allocated to each topic
- $\boldsymbol{\beta} = K \times J$ matrix, with row $\boldsymbol{\beta}_k = (\beta_{1k}, \beta_{2k}, \dots, \beta_{kJ}) \rightsquigarrow$ **topics**
- $\boldsymbol{\alpha} = K$ element long vector, population prior for $\boldsymbol{\Theta}$.

3) Optimization

- Variational Inference \rightsquigarrow deterministic approximation
- Collapsed Gibbs Sampling \rightsquigarrow MCMC algorithm
- Spectral/Factorization Methods

4) Validation \rightsquigarrow application-specific

A Statistical Highlighter (With Many Colors)

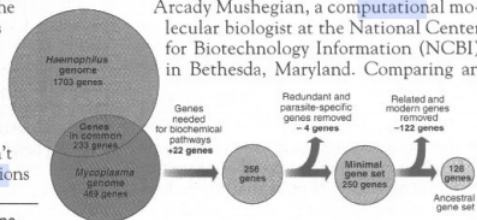
Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

“are not all that far apart,” especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. “It may be a way of organizing any newly sequenced genome,” explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

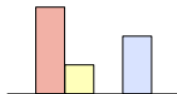


Image from Hanna Wallach

Why does this work \rightsquigarrow Co-occurrence

Where's the information for each word's topic?

Why does this work \rightsquigarrow Co-occurrence

Where's the information for each word's topic?

Reconsider document-term matrix

Why does this work \rightsquigarrow Co-occurrence

Where's the information for each word's topic?

Reconsider document-term matrix

	Word ₁	Word ₂	...	Word _J
Doc ₁	0	1	...	0
Doc ₂	2	0	...	3
⋮	⋮	⋮	⋮	⋮
Doc _N	0	1	...	1

Why does this work \rightsquigarrow Co-occurrence

Where's the information for each word's topic?

Reconsider document-term matrix

	Word ₁	Word ₂	...	Word _J
Doc ₁	0	1	...	0
Doc ₂	2	0	...	3
⋮	⋮	⋮	⋱	⋮
Doc _N	0	1	...	1

We are learning the pattern of what words occur together.

Why does this work \rightsquigarrow Co-occurrence

Where's the information for each word's topic?

Reconsider document-term matrix

	Word ₁	Word ₂	...	Word _J
Doc ₁	0	1	...	0
Doc ₂	2	0	...	3
⋮	⋮	⋮	⋮	⋮
Doc _N	0	1	...	1

We are learning the pattern of what words occur together.

The model wants a topic to contain as few words as possible, but a document to contain as few topics as possible. This **tension** is what makes the model work.

Extensions to LDA

Extensions to LDA

Have there been extensions to LDA proposed?

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

correlated topic models,

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

correlated topic models, dynamic topic models,

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

correlated topic models, dynamic topic models, hierarchical LDA,

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

correlated topic models, dynamic topic models, hierarchical LDA, pachinko allocation, nonparametric pachinko allocation,

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

correlated topic models, dynamic topic models, hierarchical LDA, pachinko allocation, nonparametric pachinko allocation, factorial LDA, gamma-poisson factorization, shared component topic models, dirichlet multinomial regression topic models,

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

correlated topic models, dynamic topic models, hierarchical LDA, pachinko allocation, nonparametric pachinko allocation, factorial LDA, gamma-poisson factorization, shared component topic models, dirichlet multinomial regression topic models, expressed agenda model, structured topic model, nested hierarchical dirichlet process topic model, focused topic model, inverse regression topic model, ideal point topic model, discrete infinite logistic normal topic model

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

correlated topic models, dynamic topic models, hierarchical LDA, pachinko allocation, nonparametric pachinko allocation, factorial LDA, gamma-poisson factorization, shared component topic models, dirichlet multinomial regression topic models, expressed agenda model, structured topic model, nested hierarchical dirichlet process topic model, focused topic model, inverse regression topic model, ideal point topic model, discrete infinite logistic normal topic model, multilingual topic model, markov topic model, relational topic model, syntactic topic model, supervised latent dirichlet allocation

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

correlated topic models, dynamic topic models, hierarchical LDA, pachinko allocation, nonparametric pachinko allocation, factorial LDA, gamma-poisson factorization, shared component topic models, dirichlet multinomial regression topic models, expressed agenda model, structured topic model, nested hierarchical dirichlet process topic model, focused topic model, inverse regression topic model, ideal point topic model, discrete infinite logistic normal topic model, multilingual topic model, markov topic model, relational topic model, syntactic topic model, supervised latent dirichlet allocation

different methods for every problem

Extensions to LDA

Have there been extensions to LDA proposed?

Yes.

correlated topic models, dynamic topic models, hierarchical LDA, pachinko allocation, nonparametric pachinko allocation, factorial LDA, gamma-poisson factorization, shared component topic models, dirichlet multinomial regression topic models, expressed agenda model, structured topic model, nested hierarchical dirichlet process topic model, focused topic model, inverse regression topic model, ideal point topic model, discrete infinite logistic normal topic model, multilingual topic model, markov topic model, relational topic model, syntactic topic model, supervised latent dirichlet allocation

different methods for every problem

What is going on with all of these extensions?

- 1 Introduction
- 2 Four Principles
- 3 Preprocessing
- 4 Latent Dirichlet Allocation
- 5 Structured Topic Models
- 6 Structural Topic Models
- 7 Sample Applications
- 8 Conclusion

- 1 Introduction
- 2 Four Principles
- 3 Preprocessing
- 4 Latent Dirichlet Allocation
- 5 Structured Topic Models**
- 6 Structural Topic Models
- 7 Sample Applications
- 8 Conclusion

Correlated Topic Models (Blei and Lafferty 2007)

Dirichlet distribution \rightsquigarrow Assumes negative covariance between topics

Correlated Topic Models (Blei and Lafferty 2007)

Dirichlet distribution \rightsquigarrow Assumes negative covariance between topics

Logistic Normal Distribution \rightsquigarrow Allows some positive covariance between topics

Correlated Topic Models (Blei and Lafferty 2007)

Dirichlet distribution \rightsquigarrow Assumes negative covariance between topics

Logistic Normal Distribution \rightsquigarrow Allows some positive covariance between topics

$$\beta_k \sim \text{Dirichlet}(\mathbf{1})$$

$$\boldsymbol{\eta}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma} \sim \text{Multivariate Normal}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\theta_i = \frac{\exp(\boldsymbol{\eta}_i)}{\sum_{k=1}^K \exp(\eta_{ik})}$$

$$z_{im} | \theta_i \sim \text{Multinomial}(1, \boldsymbol{\eta}_i)$$

$$w_{im} | \beta_k, z_{imk} = 1 \sim \text{Multinomial}(1, \beta_k)$$

Jihad Example

[illegible]

Fighting	F: Muslim, Jihad, Islam, fight, Jihadi fighters, pathway, almighty, that FREX: jihad, fighting, jihadist fighters, pulpit, approves of us, annotated, to fight, vicinity F: جهاد, قتال, مجاهد, منبر, يوافقنا, مثيل, يقاتل, بجوار FREX: مسلم, جهاد, اسلام, قتل, مجاهد, سبيل, تعالى, دين	•
Social theory	F: person, life, soul/self, knowledge/science, society, work, image, material/physical FREX: imagine, morals, develop, society, product, necessarily, environment, traditions, activity F: تصور, اخلاق, تطور, مجتمع, انتاج, حتم, بين, تقاليد FREX: انفس, حيا, نفس, علم, مجتمع, عمل, صور, ماد	•
Politics	F: Arab, Jews, country, Islam, A.D., year, West, Muslim FREX: capitol, Asia, Iran, South, Washington, A.D., Russia, Turkey F: عاصمت, اسيا, اير, جنوب, اشطنن, م, روسيا, تركيا FREX: عرب, يهود, دول, اسلام, م, سن, غرب, مسلم	•
The Prophet	F: said, prayers (be upon him), peace (be upon him), almighty, messenger, glory, prophet, that FREX: almighty, almighty, glory, bless you, magic, punishment, hypocrisy, sins F: وجل, عز, سيح, تبارك, سحر, عذاب, رياء, ذنوب FREX: قال, صل, سلم, تعالى, رسول, سيح, نب, دين	•
Prayer	F: prayer, pray, son, prophet, sheikh, mosque, fatwas, group FREX: prostration, prostrated, Abd al-Aziz, supplicant, Baz, prayer space, omission, prostration F: ركع, ركعت, عبدالعزيز, ماموم, باز, مصل, سهو, ركوع FREX: صل, صل, سلم, بن, نب, شيخ, مسجد, قنوا	•
Ramadan	F: day, fasting, Ashura, Ramadan, sheikh, group, fatwas, Uthaymeen FREX: wash, one who fasts, fasting, fasting, to break fast, Ramadan, travel, dirty F: غسل, صائم, صيام, صوم, يفطر, رمض, مسافر, نجاس FREX: يوم, صيام, عشر, رمض, شيخ, مجموع, قنوا, عثم	•
Family and Women	F: woman, O, man, girl, one, says, men, people FREX: veil, youth, (sheikh) Tamim, Azzam, tanks, finery, wear, r(typo) F: حجاب, شاب, تميم, عزام, دياب, تيرج, لباس, ر FREX: مرا, يا, رجل, نساء, احد, يقول, رجال, ناس	•
Money, Pilgrimage, and Marriage	F: tithing, money, pilgrimage, permitted, religion, marriage, believe/ratify, divorce FREX: tithing, divorce, banks, divorce, card, banks, to perform pilgrimage, poor F: زكا, مطلق, بنك, مطلق, بطاق, بنوك, يحج, فقراء FREX: زكا, مال, حج, يجوز, دين, زوج, صدق, طلاق	•
Islam and Modernity	F: Islam, land, mankind, people, religion, life, other, God FREX: Europe, civilization, European, mankind, church, goods, generations, their lives F: اوربا, حضار, اورب, بشر, كنيس, متاع, اجيال, حياتهم FREX: اسلام, ارض, بشر, ناس, دين, اخر, ال	•
Hadith	F: Saying, hadith, said, prayers (be upon him), peace (be upon him), Muslim, legally, not FREX: to forbid, analogy, permission, general, evidence, forbid, text, absolutely F: تحريم, قياس, جواز, عموم, ادل, منع, نص, مطلقا FREX: قول, حديث, قال, صل, سلم, مسلم, شرع, ليس	•
Excommunication	F: Apostasy, said, almighty, polytheism, Islam, Apostate, saying, people	

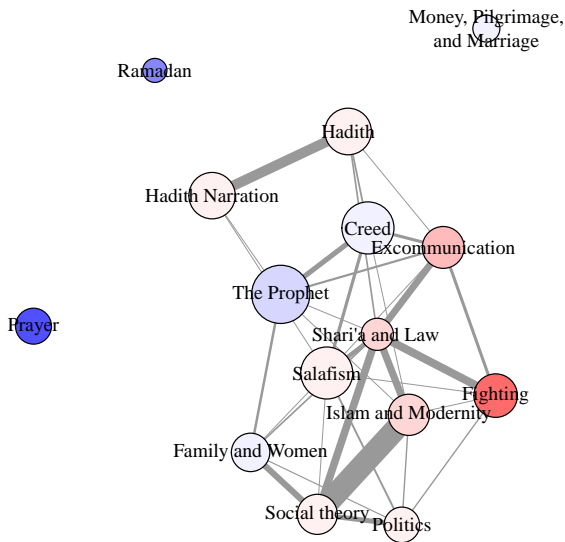
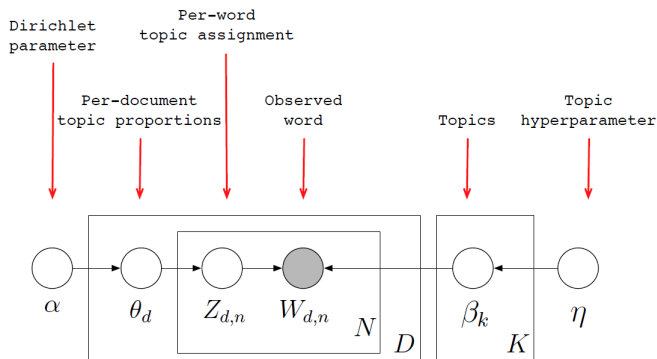


Figure: The network of correlated topics for a 15-topic Structural Topic Model with Jihadi/not-Jihadi as the predictor of topics in Arab Muslim cleric writings.

LDA \rightsquigarrow Dynamic Topic Model

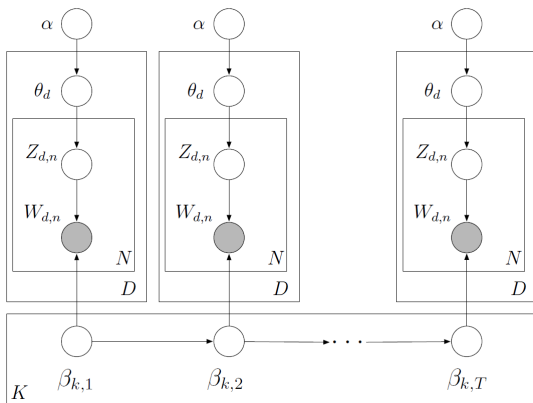
(Blei and Lafferty 2007)

Figure: Plate Notation of Latent Dirichlet Allocation



LDA \rightsquigarrow Dynamic Topic Model

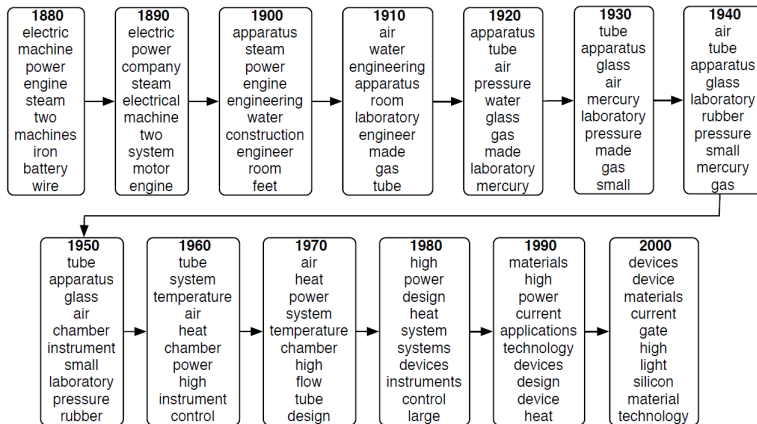
Figure: Dynamic Topic Model



Graphic from David Blei

LDA \rightsquigarrow Dynamic Topic Model

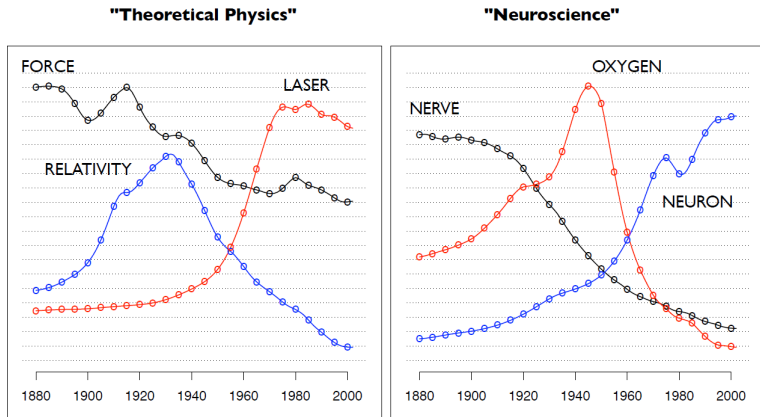
Figure: Topic Evolution over Time



Graphic from David Blei

LDA \rightsquigarrow Dynamic Topic Model

Figure: Word Use in Topics Over Time



Graphic from David Blei

Expressed Agenda Model (Grimmer 2010)

Expressed Agenda Model (Grimmer 2010)

① Assumes:

Expressed Agenda Model (Grimmer 2010)

- ① Assumes:
 - ① Each document is assigned to one topic

Expressed Agenda Model (Grimmer 2010)

① Assumes:

- ① Each document is assigned to one topic
- ② Each author allocates some hidden proportion of time to each topic

Expressed Agenda Model (Grimmer 2010)

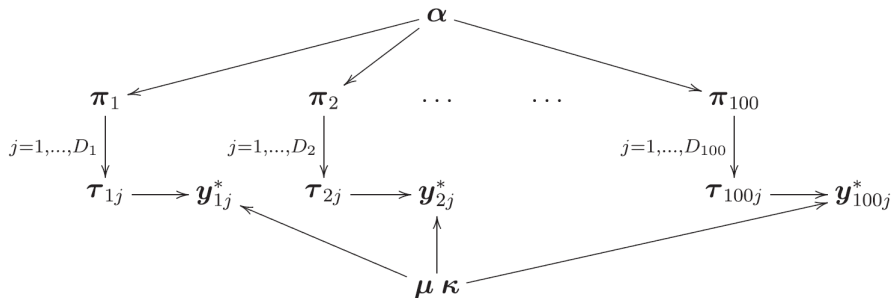
- ① Assumes:
 - ① Each document is assigned to one topic
 - ② Each author allocates some hidden proportion of time to each topic
- ② Grimmer's project seeks to quantitatively represent the content of senators' press releases.

Expressed Agenda Model (Grimmer 2010)

- ① Assumes:
 - ① Each document is assigned to one topic
 - ② Each author allocates some hidden proportion of time to each topic
- ② Grimmer's project seeks to quantitatively represent the content of senators' press releases.
- ③ It is called the **Expressed** Agenda Model because it captures the way they communicate that agenda to constituents.

Expressed Agenda Model

Figure: Expressed Agenda Model



Graphic from Grimmer 2010

- 1 Introduction
- 2 Four Principles
- 3 Preprocessing
- 4 Latent Dirichlet Allocation
- 5 Structured Topic Models
- 6 Structural Topic Models
- 7 Sample Applications
- 8 Conclusion

- 1 Introduction
- 2 Four Principles
- 3 Preprocessing
- 4 Latent Dirichlet Allocation
- 5 Structured Topic Models
- 6 Structural Topic Models**
- 7 Sample Applications
- 8 Conclusion

STM = LDA + Contextual Information

STM = LDA + Contextual Information

- STM provides two ways to include contextual information

STM = LDA + Contextual Information

- STM provides two ways to include contextual information
 - ▶ Topic prevalence can vary by metadata

STM = LDA + Contextual Information

- STM provides two ways to include contextual information
 - ▶ Topic **prevalence** can vary by metadata
 - ★ e.g. city papers cover protests more than provincial papers

STM = LDA + Contextual Information

- STM provides two ways to include contextual information
 - ▶ Topic **prevalence** can vary by metadata
 - ★ e.g. city papers cover protests more than provincial papers
 - ▶ Topic **content** can vary by metadata

STM = LDA + Contextual Information

- STM provides two ways to include contextual information
 - ▶ Topic **prevalence** can vary by metadata
 - ★ e.g. city papers cover protests more than provincial papers
 - ▶ Topic **content** can vary by metadata
 - ★ e.g. city papers talk about protests differently

STM = LDA + Contextual Information

- STM provides two ways to include contextual information
 - ▶ Topic **prevalence** can vary by metadata
 - ★ e.g. city papers cover protests more than provincial papers
 - ▶ Topic **content** can vary by metadata
 - ★ e.g. city papers talk about protests differently
- Including context improves the model:

STM = LDA + Contextual Information

- STM provides two ways to include contextual information
 - ▶ Topic **prevalence** can vary by metadata
 - ★ e.g. city papers cover protests more than provincial papers
 - ▶ Topic **content** can vary by metadata
 - ★ e.g. city papers talk about protests differently
- Including context improves the model:
 - ▶ more accurate estimation

STM = LDA + Contextual Information

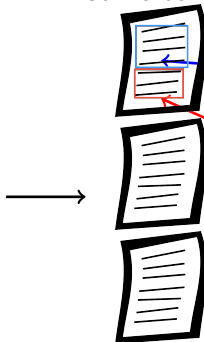
- STM provides two ways to include contextual information
 - ▶ Topic **prevalence** can vary by metadata
 - ★ e.g. city papers cover protests more than provincial papers
 - ▶ Topic **content** can vary by metadata
 - ★ e.g. city papers talk about protests differently
- Including context improves the model:
 - ▶ more accurate estimation
 - ▶ better qualitative interpretability

STM: What this means in pictures

Say you have
a lot of people.



Each writes
some text



that discuss a few
different topics

Politics

congress, nations,
power, votes, agree-
ment, bargaining

Statistics

estimator, data,
analysis, variance,
model, inference

The STM Allows for:

STM: What this means in pictures

Say you have
a lot of people.



Each writes
some text



that discuss a few
different topics

Politics

congress, nations,
power, votes, agree-
ment, bargaining

Statistics

estimator, data,
analysis, variance,
model, inference

The STM Allows for:

- ① The words in each topic to vary by gender

STM: What this means in pictures

Say you have
a lot of people.



Each writes
some text



that discuss a few
different topics

Politics

congress, nations,
power, votes, agree-
ment, bargaining

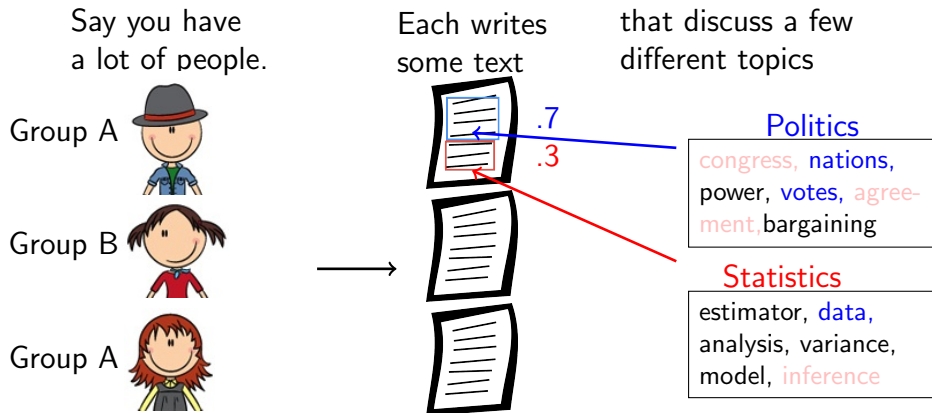
Statistics

estimator, data,
analysis, variance,
model, inference

The STM Allows for:

- ① The words in each topic to vary by gender

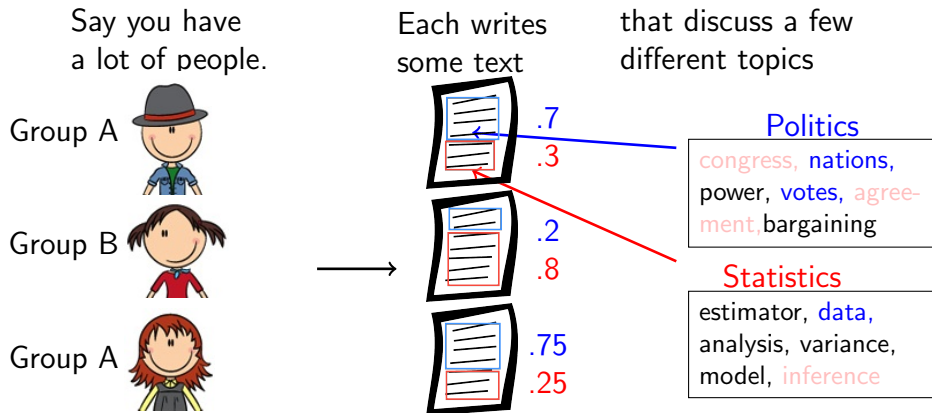
STM: What this means in pictures



The STM Allows for:

- 1 The words in each topic to vary by gender
- 2 The topic proportions to vary by group

STM: What this means in pictures



The STM Allows for:

- 1 The words in each topic to vary by gender
- 2 The topic proportions to vary by group

Mixed-Membership Topic Models

More formal terminology:

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries
 - ▶ Data summarized in a single matrix $D \times V$ matrix \mathbf{W}

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries
 - ▶ Data summarized in a single matrix $D \times V$ matrix \mathbf{W}
- Additional data for STM

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries
 - ▶ Data summarized in a single matrix $D \times V$ matrix \mathbf{W}
- Additional data for STM
 - ▶ Topic prevalence covariates: $D \times P$ matrix \mathbf{X}

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries
 - ▶ Data summarized in a single matrix $D \times V$ matrix \mathbf{W}
- Additional data for STM
 - ▶ Topic prevalence covariates: $D \times P$ matrix \mathbf{X}
 - ▶ Topical content groups: D length vector \mathbf{Y}

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries
 - ▶ Data summarized in a single matrix $D \times V$ matrix \mathbf{W}
- Additional data for STM
 - ▶ Topic prevalence covariates: $D \times P$ matrix \mathbf{X}
 - ▶ Topical content groups: D length vector \mathbf{Y}
- Latent variables

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries
 - ▶ Data summarized in a single matrix $D \times V$ matrix \mathbf{W}
- Additional data for STM
 - ▶ Topic prevalence covariates: $D \times P$ matrix \mathbf{X}
 - ▶ Topical content groups: D length vector \mathbf{Y}
- Latent variables
 - ▶ $D \times K$ matrix θ : proportion of document on each topic.

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries
 - ▶ Data summarized in a single matrix $D \times V$ matrix \mathbf{W}
- Additional data for STM
 - ▶ Topic prevalence covariates: $D \times P$ matrix \mathbf{X}
 - ▶ Topical content groups: D length vector \mathbf{Y}
- Latent variables
 - ▶ $D \times K$ matrix θ : proportion of document on each topic.
 - ▶ $K \times V$ matrix β : probability of drawing a word conditional on topic.

Mixed-Membership Topic Models

More formal terminology:

- User specifies the number of topics: K
- Observed data for standard topic models
 - ▶ Each document ($d \in 1 \dots D$) is a collection of N_d tokens
 - ▶ Each token is a particular word from a dictionary of V entries
 - ▶ Data summarized in a single matrix $D \times V$ matrix \mathbf{W}
- Additional data for STM
 - ▶ Topic prevalence covariates: $D \times P$ matrix \mathbf{X}
 - ▶ Topical content groups: D length vector \mathbf{Y}
- Latent variables
 - ▶ $D \times K$ matrix θ : proportion of document on each topic.
 - ▶ $K \times V$ matrix β : probability of drawing a word conditional on topic.
 - ▶ Low rank approximation to expected counts:
$$\tilde{\mathbf{W}}_{D \times V} \approx \theta_{D \times K} \beta_{K \times V}$$

Technical Details: The Structural Topic Model

- Low rank approximation to expected counts: $\tilde{W}_{D \times V} \approx \theta_{D \times K} \beta_{K \times V}$
- θ , $D \times K$ document-topic matrix
- β , $K \times V$ topic-word matrix
- Each token has a topic drawn from the document mixture
 - ▶ Draw token topic $z_{d,n}$ from $\text{Discrete}(\theta_d)$
 - ▶ Draw observed word $w_{d,n}$ from $\text{Discrete}(\beta_{k=z_{d,n}})$

Technical Details: The Structural Topic Model

- Low rank approximation to expected counts: $\tilde{W}_{D \times V} \approx \theta_{D \times K} \beta_{K \times V}$
- θ , $D \times K$ document-topic matrix \Leftarrow logistic normal glm with covariates
- β , $K \times V$ topic-word matrix
- Each token has a topic drawn from the document mixture
 - ▶ Draw token topic $z_{d,n}$ from $\text{Discrete}(\theta_d)$
 - ▶ Draw observed word $w_{d,n}$ from $\text{Discrete}(\beta_{k=z_{d,n}})$

Technical Details: The Structural Topic Model

- Low rank approximation to expected counts: $\tilde{W}_{D \times V} \approx \theta_{D \times K} \beta_{K \times V}$
 - θ , $D \times K$ document-topic matrix \Leftarrow **logistic normal glm with covariates**
 - ▶ Covariate-specific prior with global topic covariance
 - ▶ $\theta_{d,\cdot} \sim \text{LogisticNormal}(X_d \gamma, \Sigma)$
 - β , $K \times V$ topic-word matrix
-
- Each token has a topic drawn from the document mixture
 - ▶ Draw token topic $z_{d,n}$ from $\text{Discrete}(\theta_d)$
 - ▶ Draw observed word $w_{d,n}$ from $\text{Discrete}(\beta_{k=z_{d,n}})$

Technical Details: The Structural Topic Model

- Low rank approximation to expected counts: $\tilde{W}_{D \times V} \approx \theta_{D \times K} \beta_{K \times V}$
 - θ , $D \times K$ document-topic matrix \Leftarrow **logistic normal glm with covariates**
 - ▶ Covariate-specific prior with global topic covariance
 - ▶ $\theta_{d,\cdot} \sim \text{LogisticNormal}(X_d \gamma, \Sigma)$
 - β , $K \times V$ topic-word matrix \Leftarrow **multinomial logit with covariates**
-
- Each token has a topic drawn from the document mixture
 - ▶ Draw token topic $z_{d,n}$ from $\text{Discrete}(\theta_d)$
 - ▶ Draw observed word $w_{d,n}$ from $\text{Discrete}(\beta_{k=z_{d,n}})$

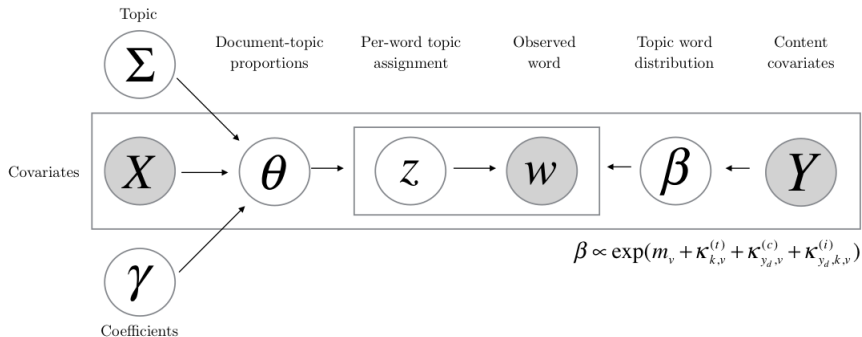
Technical Details: The Structural Topic Model

- Low rank approximation to expected counts: $\tilde{W}_{D \times V} \approx \theta_{D \times K} \beta_{K \times V}$
- θ , $D \times K$ document-topic matrix \Leftarrow **logistic normal glm with covariates**
 - ▶ Covariate-specific prior with global topic covariance
 - ▶ $\theta_{d,\cdot} \sim \text{LogisticNormal}(X_d \gamma, \Sigma)$
- β , $K \times V$ topic-word matrix \Leftarrow **multinomial logit with covariates**
 - ▶ Each topic is now a sparse, covariate-specific deviation from a baseline distribution.
 - ▶ $\vec{\beta}_{k,\cdot} \propto \exp(m + \kappa^{(\text{topic})} + \kappa^{(\text{cov})} + \kappa^{(\text{int})})$
 - ▶ Three parts: topic, covariate, topic-covariate interaction
- Each token has a topic drawn from the document mixture
 - ▶ Draw token topic $z_{d,n}$ from $\text{Discrete}(\theta_d)$
 - ▶ Draw observed word $w_{d,n}$ from $\text{Discrete}(\beta_{k=z_{d,n}})$

Technical Details: The Structural Topic Model

- Low rank approximation to expected counts: $\tilde{W}_{D \times V} \approx \theta_{D \times K} \beta_{K \times V}$
- θ , $D \times K$ document-topic matrix \Leftarrow **logistic normal glm with covariates**
 - ▶ Covariate-specific prior with global topic covariance
 - ▶ $\theta_{d,\cdot} \sim \text{LogisticNormal}(X_d \gamma, \Sigma)$
- β , $K \times V$ topic-word matrix \Leftarrow **multinomial logit with covariates**
 - ▶ Each topic is now a sparse, covariate-specific deviation from a baseline distribution.
 - ▶ $\vec{\beta}_{k,\cdot} \propto \exp(m + \kappa^{(\text{topic})} + \kappa^{(\text{cov})} + \kappa^{(\text{int})})$
 - ▶ Three parts: topic, covariate, topic-covariate interaction
 - ▶ β may instead be **point-estimated**
- Each token has a topic drawn from the document mixture
 - ▶ Draw token topic $z_{d,n}$ from $\text{Discrete}(\theta_d)$
 - ▶ Draw observed word $w_{d,n}$ from $\text{Discrete}(\beta_{k=z_{d,n}})$

Structural Topic Model



Estimation and Implementation of the STM

Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters

Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
 - ▶ bayesian estimation using variational inference

Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
 - ▶ bayesian estimation using variational inference
(initialization from spectral method of moments estimator)

Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
 - ▶ bayesian estimation using variational inference
(initialization from spectral method of moments estimator)
 - ▶ essentially word co-occurences used to discover topics

Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
 - ▶ bayesian estimation using variational inference
(initialization from spectral method of moments estimator)
 - ▶ essentially word co-occurrences used to discover topics
- General to many kinds of corpus structure using covariates

Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
 - ▶ bayesian estimation using variational inference
(initialization from spectral method of moments estimator)
 - ▶ essentially word co-occurrences used to discover topics
- General to many kinds of corpus structure using covariates
- `stm` Package in R

Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
 - ▶ bayesian estimation using variational inference
(initialization from spectral method of moments estimator)
 - ▶ essentially word co-occurrences used to discover topics
- General to many kinds of corpus structure using covariates
- `stm` Package in R
 - ▶ complete workflow: raw texts → figures

Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
 - ▶ bayesian estimation using variational inference
(initialization from spectral method of moments estimator)
 - ▶ essentially word co-occurrences used to discover topics
 - General to many kinds of corpus structure using covariates
 - stm Package in R
 - ▶ complete workflow: raw texts → figures
 - ▶ simple regression style syntax using formulas
- ```
mod.out <- stm(documents,vocab, K=10,
 prevalence= ~paper + s(time),
 data=metadata, init.type="Spectral")
```

# Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
  - ▶ bayesian estimation using variational inference  
(initialization from spectral method of moments estimator)
  - ▶ essentially word co-occurrences used to discover topics
- General to many kinds of corpus structure using covariates
- stm Package in R
  - ▶ complete workflow: raw texts → figures
  - ▶ simple regression style syntax using formulas

```
mod.out <- stm(documents,vocab, K=10,
 prevalence= ~paper + s(time),
 data=metadata, init.type="Spectral")
```
  - ▶ many functions for summarization, visualization and checking

# Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
  - ▶ bayesian estimation using variational inference  
(initialization from spectral method of moments estimator)
  - ▶ essentially word co-occurrences used to discover topics
- General to many kinds of corpus structure using covariates
- stm Package in R
  - ▶ complete workflow: raw texts → figures
  - ▶ simple regression style syntax using formulas

```
mod.out <- stm(documents,vocab, K=10,
 prevalence= ~paper + s(time),
 data=metadata, init.type="Spectral")
```
  - ▶ many functions for summarization, visualization and checking
- Complete vignette online with examples

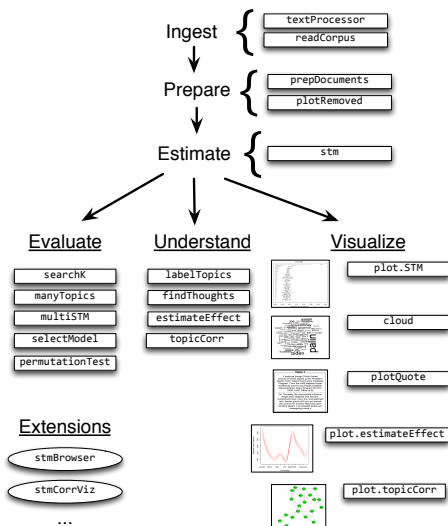
# Estimation and Implementation of the STM

- Define a probabilistic model and estimate parameters
  - ▶ bayesian estimation using variational inference  
(initialization from spectral method of moments estimator)
  - ▶ essentially word co-occurrences used to discover topics
- General to many kinds of corpus structure using covariates
- stm Package in R
  - ▶ complete workflow: raw texts → figures
  - ▶ simple regression style syntax using formulas

```
mod.out <- stm(documents,vocab, K=10,
 prevalence= ~paper + s(time),
 data=metadata, init.type="Spectral")
```
  - ▶ many functions for summarization, visualization and checking
- Complete vignette online with examples

You can do this with your data!

# stm is Full of Functions to Help You!



- 1 Introduction
- 2 Four Principles
- 3 Preprocessing
- 4 Latent Dirichlet Allocation
- 5 Structured Topic Models
- 6 Structural Topic Models
- 7 Sample Applications
- 8 Conclusion

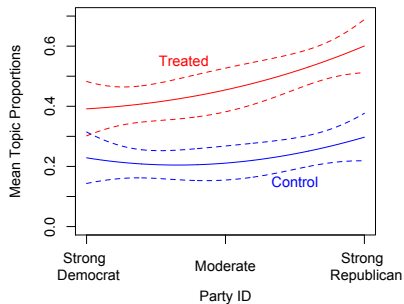
- 1 Introduction
- 2 Four Principles
- 3 Preprocessing
- 4 Latent Dirichlet Allocation
- 5 Structured Topic Models
- 6 Structural Topic Models
- 7 Sample Applications**
- 8 Conclusion

# Applications



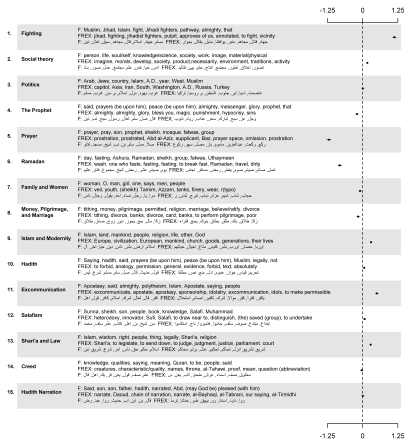
# Applications

- Open-Ended Survey Response  
(Roberts et al 2014)



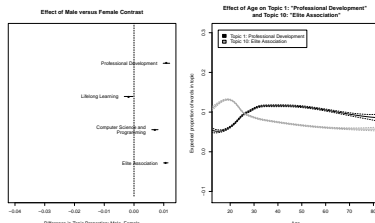
# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)



# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)
- Student Text in MOOCs (Reich et al 2015)





# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)
- Student Text in MOOCs (Reich et al 2015)
- Constitutional Moments (Stewart and Young)
- Chinese/Western Newswires (Roberts et al 2016)

# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)
- Student Text in MOOCs (Reich et al 2015)
- Constitutional Moments (Stewart and Young)
- Chinese/Western Newswires (Roberts et al 2016)
- **How Teams Make Forecasts** (Horowitz et al 2017)

# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)
- Student Text in MOOCs (Reich et al 2015)
- Constitutional Moments (Stewart and Young)
- Chinese/Western Newswires (Roberts et al 2016)
- How Teams Make Forecasts (Horowitz et al 2017)
- Self-Censorship of Bloggers (Roberts et al 2017)

# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)
- Student Text in MOOCs (Reich et al 2015)
- Constitutional Moments (Stewart and Young)
- Chinese/Western Newswires (Roberts et al 2016)
- How Teams Make Forecasts (Horowitz et al 2017)
- Self-Censorship of Bloggers (Roberts et al 2017)
- **Snowden Leak Response** (Lucas et al 2015)



# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)
- Student Text in MOOCs (Reich et al 2015)
- Constitutional Moments (Stewart and Young)
- Chinese/Western Newswires (Roberts et al 2016)
- How Teams Make Forecasts (Horowitz et al 2017)
- Self-Censorship of Bloggers (Roberts et al 2017)
- Snowden Leak Response (Lucas et al 2015)

Many applications **outside** our group:

# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)
- Student Text in MOOCs (Reich et al 2015)
- Constitutional Moments (Stewart and Young)
- Chinese/Western Newswires (Roberts et al 2016)
- How Teams Make Forecasts (Horowitz et al 2017)
- Self-Censorship of Bloggers (Roberts et al 2017)
- Snowden Leak Response (Lucas et al 2015)

Many applications **outside** our group:

- Over 25 published articles in over 20 journals

# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)
- Student Text in MOOCs (Reich et al 2015)
- Constitutional Moments (Stewart and Young)
- Chinese/Western Newswires (Roberts et al 2016)
- How Teams Make Forecasts (Horowitz et al 2017)
- Self-Censorship of Bloggers (Roberts et al 2017)
- Snowden Leak Response (Lucas et al 2015)

Many applications **outside** our group:

- Over 25 published articles in over 20 journals
- Used in education, politics, climate science, sociology, law, statistics, public policy, transportation systems

# Applications

- Open-Ended Survey Response (Roberts et al 2014)
- Fatwas of Jihadi Clerics (Lucas et al 2015)
- Student Text in MOOCs (Reich et al 2015)
- Constitutional Moments (Stewart and Young)
- Chinese/Western Newswires (Roberts et al 2016)
- How Teams Make Forecasts (Horowitz et al 2017)
- Self-Censorship of Bloggers (Roberts et al 2017)
- Snowden Leak Response (Lucas et al 2015)

Many applications **outside** our group:

- Over 25 published articles in over 20 journals
- Used in education, politics, climate science, sociology, law, statistics, public policy, transportation systems
- More on the way...

# Conclusion

# Conclusion

- Emerging opportunities for text analysis in the social sciences

# Conclusion

- Emerging opportunities for text analysis in the social sciences
  - ▶ new models for text analysis using document **context** (STM)

# Conclusion

- Emerging opportunities for text analysis in the social sciences
  - ▶ new models for text analysis using document **context** (STM)
  - ▶ open-source R package **stm**



# Conclusion

- Emerging opportunities for text analysis in the social sciences
  - ▶ new models for text analysis using document **context** (STM)
  - ▶ open-source R package **stm**  
(auxiliary packages `stmCorrViz` and `stmBrowser`)

# Conclusion

- Emerging opportunities for text analysis in the social sciences
  - ▶ new models for text analysis using document **context** (STM)
  - ▶ open-source R package **stm**  
(auxiliary packages `stmCorrViz` and `stmBrowser`)
  - ▶ other great packages for LDA in R (`mallet`, `topicmodels`, `lda`)

# Conclusion

- Emerging opportunities for text analysis in the social sciences
  - ▶ new models for text analysis using document **context** (STM)
  - ▶ open-source R package **stm**  
(auxiliary packages `stmCorrViz` and `stmBrowser`)
  - ▶ other great packages for LDA in R (`mallet`, `topicmodels`, `lda`)
- Talk has necessarily skipped over many, many important details—be sure to read more!

Go try out the software today!

# Suggested Reading

- Blei (2012) “Probabilistic Topic Models” *Transactions of the ACM*.
- Wallach, Mimno and McCallum (2009) “Rethinking LDA: Why Priors Matter” *NIPS*
- Grimmer and Stewart (2013) “Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts” *Political Analysis*
- Roberts et al. (2014) “Structural topic models for open-ended survey responses” *American Journal of Political Science*
- Boyd-Graber, Mimno and Newman (2016) “Care and Feeding of Topic Models: Problems, Diagnostics and Improvements” in *Handbook of Mixed Membership Models*

For more information

BrandonStewart.org

structuraltopicmodel.com