

Text as Data

Justin Grimmer

Professor
Department of Political Science
Stanford University

May 15th, 2019

Measurement via repurposed discovery methods

- 1) Discovery categories, measure prevalence of categories
- 2) Once we fix **interpretation**, accuracy/precision/recall well defined

LDA Revisited

$$\boldsymbol{\theta}_k \sim \text{Dirichlet}(\mathbf{1})$$

$$\boldsymbol{\pi}_i | \boldsymbol{\alpha} \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$\tau_{im} | \boldsymbol{\pi}_i \sim \text{Multinomial}(1, \boldsymbol{\pi}_i)$$

$$x_{im} | \boldsymbol{\theta}_k, \tau_{imk} = 1 \sim \text{Multinomial}(1, \boldsymbol{\theta}_k)$$

LDA Revisited

Unigram Model_{*k*} \sim Dirichlet(**1**)

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

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

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

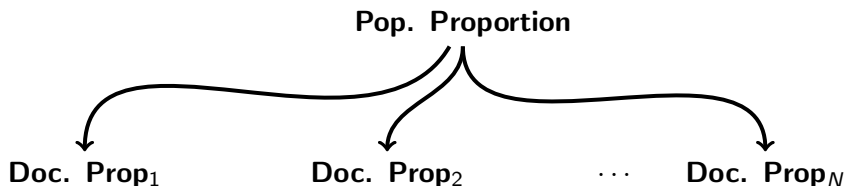
A General Hierarchical Structure

LDA:

Pop. Proportion

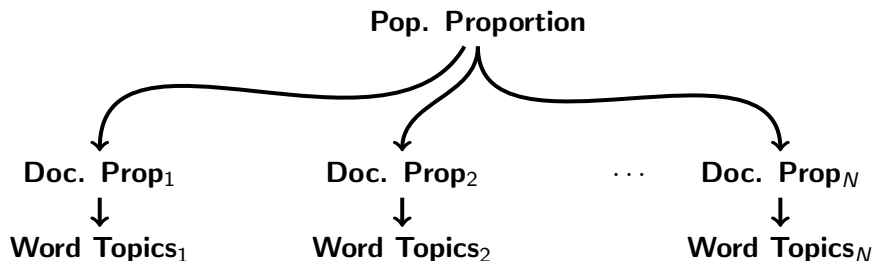
A General Hierarchical Structure

LDA:



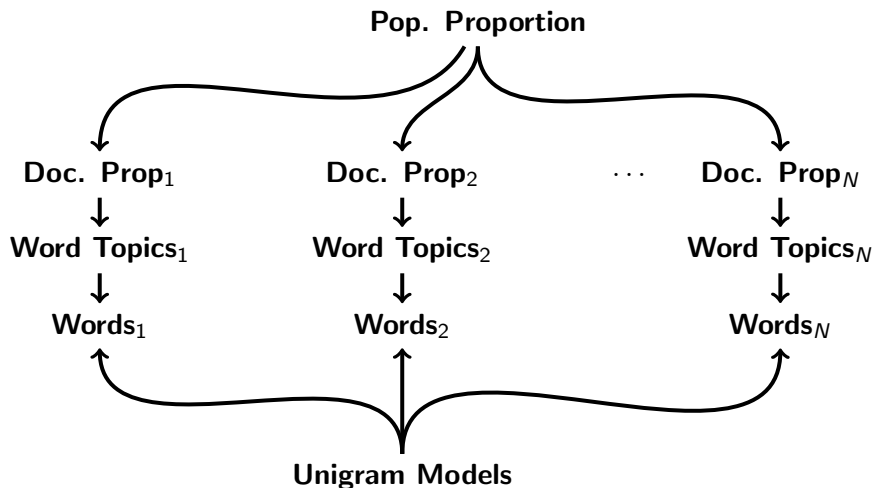
A General Hierarchical Structure

LDA:



A General Hierarchical Structure

LDA:



A General Hierarchical Structure

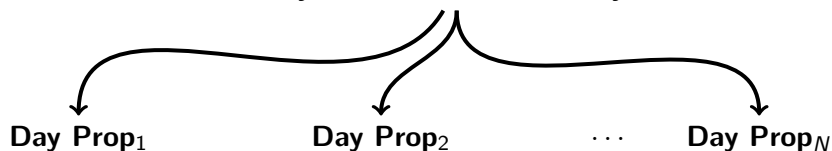
Dynamic Topic Model (Quinn et al 2010)

Dynamic Prior Across Days

A General Hierarchical Structure

Dynamic Topic Model (Quinn et al 2010)

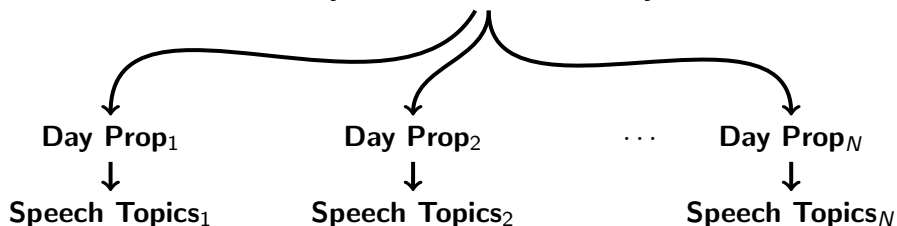
Dynamic Prior Across Days



A General Hierarchical Structure

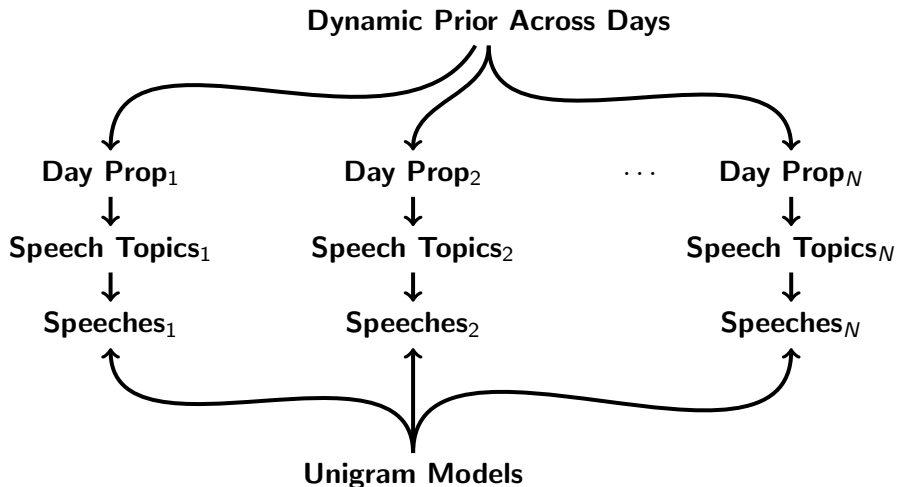
Dynamic Topic Model (Quinn et al 2010)

Dynamic Prior Across Days



A General Hierarchical Structure

Dynamic Topic Model (Quinn et al 2010)



A General Hierarchical Structure

Expressed Agenda Model (Grimmer 2010)

Average Attention Across Authors

A General Hierarchical Structure

Expressed Agenda Model (Grimmer 2010)

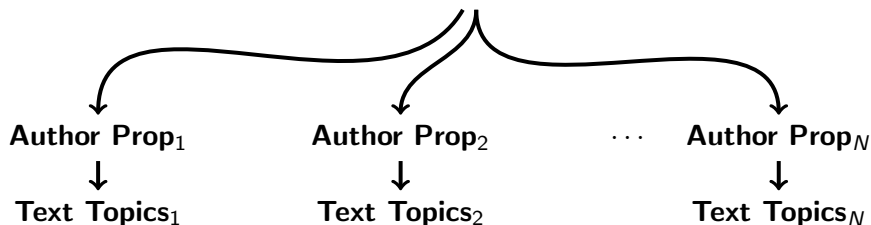
Average Attention Across Authors



A General Hierarchical Structure

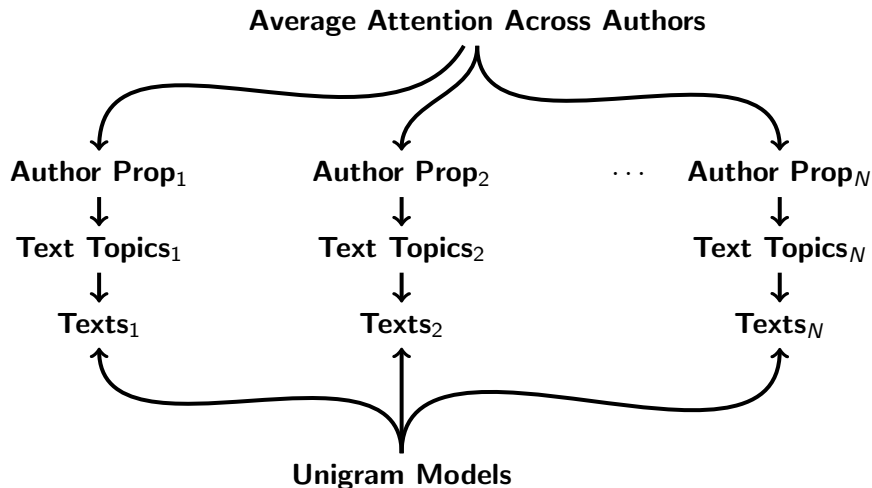
Expressed Agenda Model (Grimmer 2010)

Average Attention Across Authors



A General Hierarchical Structure

Expressed Agenda Model (Grimmer 2010)



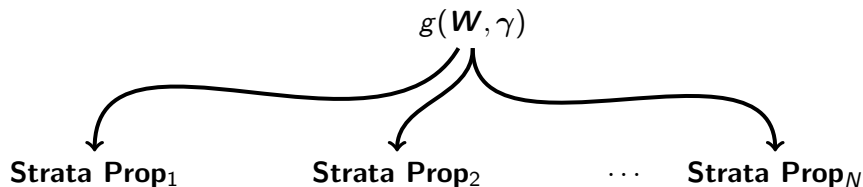
A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airoldi 2014)

$$g(\mathbf{W}, \gamma)$$

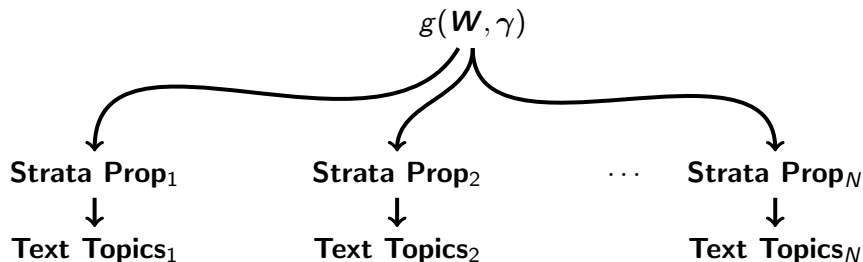
A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airolidi 2014)



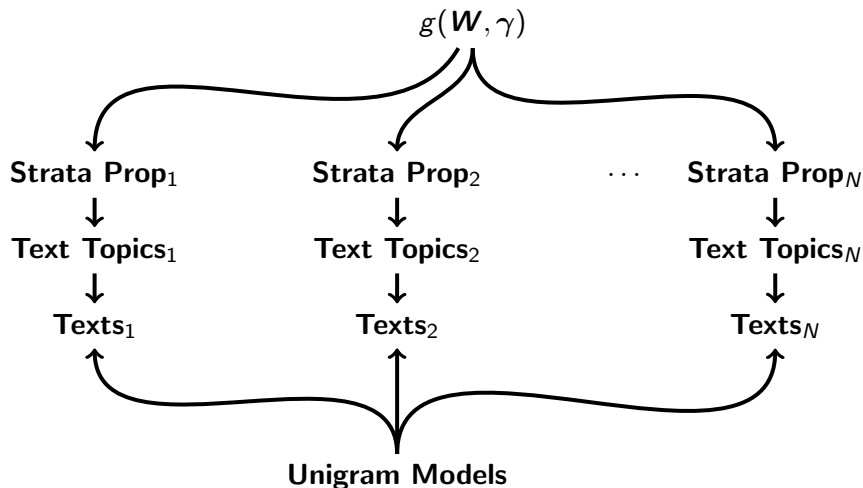
A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airolidi 2014)



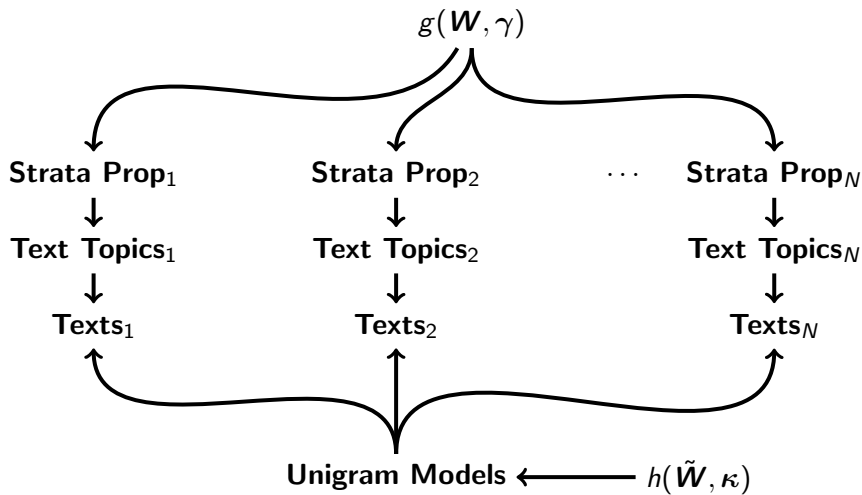
A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airolidi 2014)



A General Hierarchical Structure

Structural Topic Model (Roberts, Stewart, Airolidi 2014)



R Code

A General Hierarchical Structure

Conditioning on Unknown Covariates \rightsquigarrow levels of mixtures at proportions
(Grimmer 2013; Wallach 2008)

Mixture of Top. Attn. Models

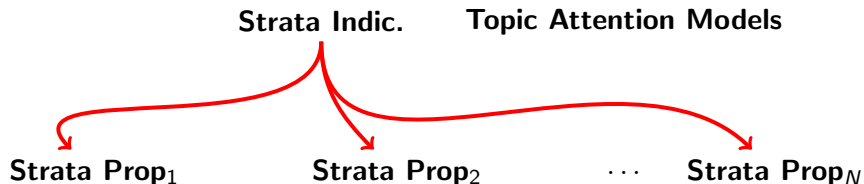
A General Hierarchical Structure

Conditioning on Unknown Covariates \rightsquigarrow levels of mixtures at proportions
(Grimmer 2013; Wallach 2008)

Strata Indic. **Topic Attention Models**

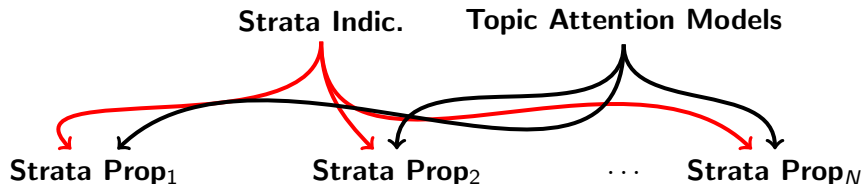
A General Hierarchical Structure

Conditioning on Unknown Covariates \rightsquigarrow levels of mixtures at proportions
(Grimmer 2013; Wallach 2008)



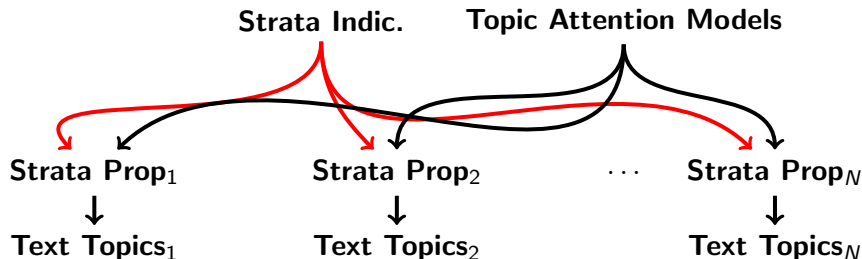
A General Hierarchical Structure

Conditioning on Unknown Covariates \rightsquigarrow levels of mixtures at proportions
(Grimmer 2013; Wallach 2008)



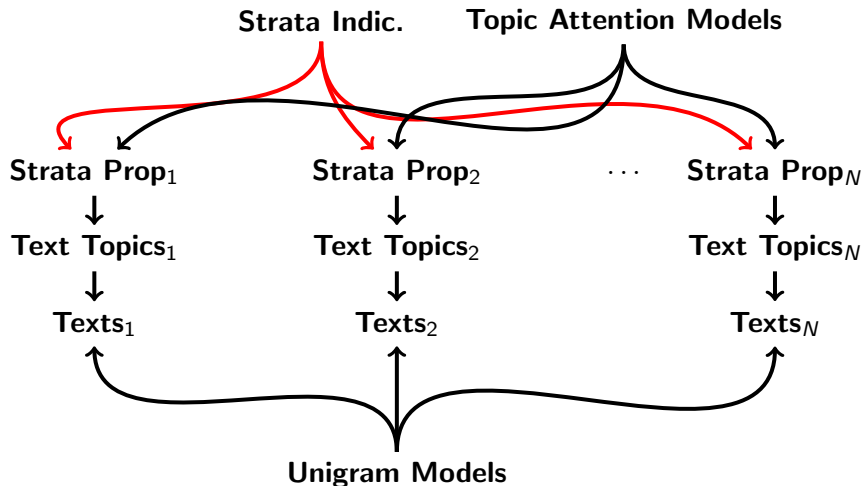
A General Hierarchical Structure

Conditioning on Unknown Covariates \rightsquigarrow levels of mixtures at proportions
(Grimmer 2013; Wallach 2008)



A General Hierarchical Structure

Conditioning on Unknown Covariates \rightsquigarrow levels of mixtures at proportions
(Grimmer 2013; Wallach 2008)



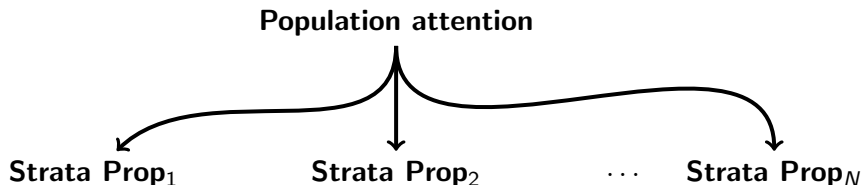
A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics \rightsquigarrow hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2017)

Population attention

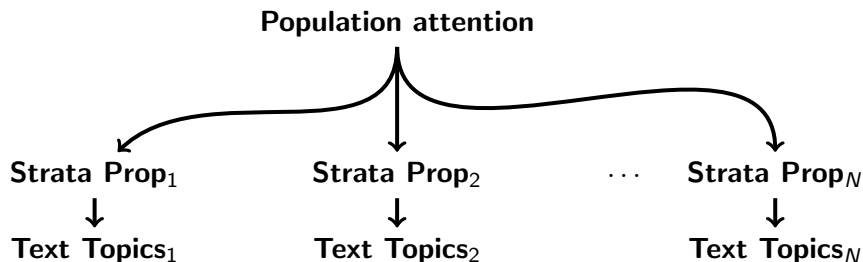
A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics \rightsquigarrow hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2017)



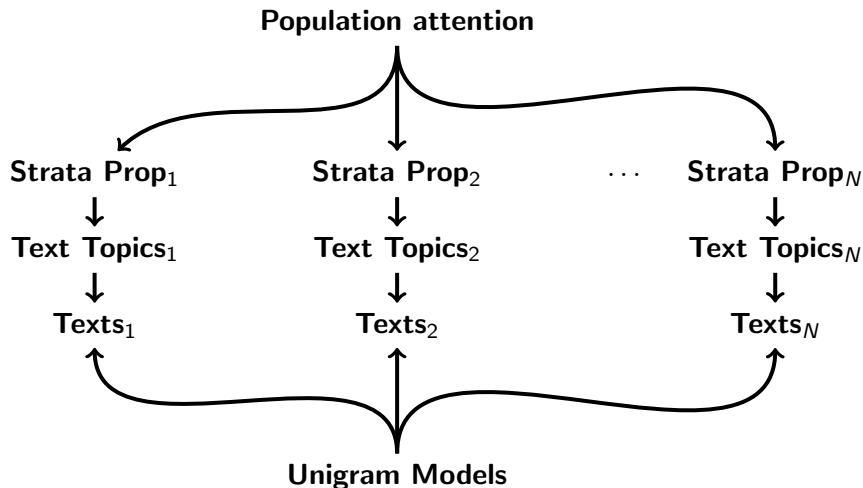
A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics \rightsquigarrow hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2017)



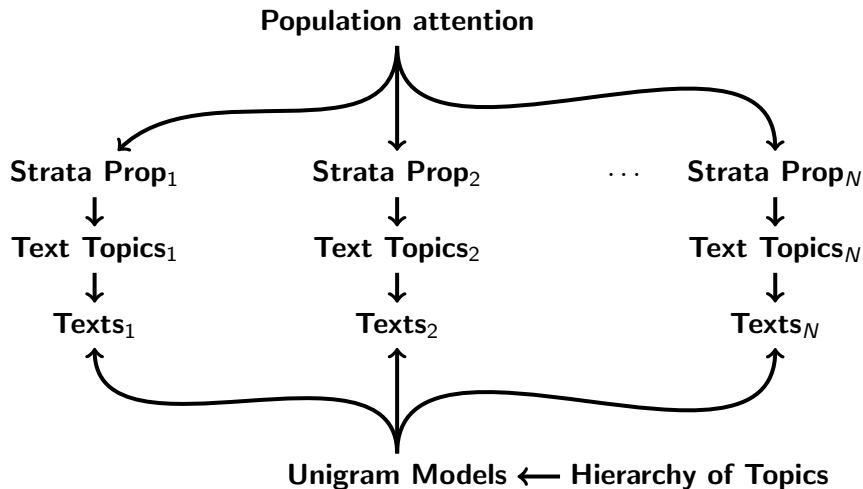
A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics \rightsquigarrow hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2017)



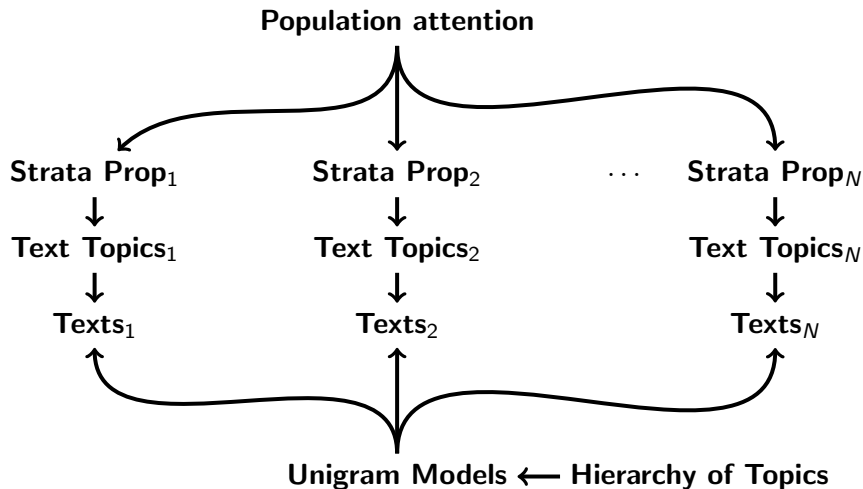
A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics \rightsquigarrow hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2017)



A General Hierarchical Structure

Conditioning on Unknown Covariates for Topics \rightsquigarrow hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2017)



Why Encode Structure in Extensions of LDA?

Why Encode Structure in Extensions of LDA?

- Substantive reasons

Why Encode Structure in Extensions of LDA?

- Substantive reasons
 - Additional structure corresponds to substantively interesting content

Why Encode Structure in Extensions of LDA?

- Substantive reasons
 - Additional structure corresponds to substantively interesting content
 - Avoids potential ad-hoc secondary analysis

Why Encode Structure in Extensions of LDA?

- Substantive reasons
 - Additional structure corresponds to substantively interesting content
 - Avoids potential ad-hoc secondary analysis
 - Clear data generating process

Why Encode Structure in Extensions of LDA?

- Substantive reasons
 - Additional structure corresponds to substantively interesting content
 - Avoids potential ad-hoc secondary analysis
 - Clear data generating process
- Statistical reasons

Why Encode Structure in Extensions of LDA?

- Substantive reasons
 - Additional structure corresponds to substantively interesting content
 - Avoids potential ad-hoc secondary analysis
 - Clear data generating process
- Statistical reasons
 - **Smoothing** \rightsquigarrow borrow information across groups intelligently

Why Encode Structure in Extensions of LDA?

- Substantive reasons
 - Additional structure corresponds to substantively interesting content
 - Avoids potential ad-hoc secondary analysis
 - Clear data generating process
- Statistical reasons
 - **Smoothing** \rightsquigarrow borrow information across groups intelligently
 - **Uncertainty** \rightsquigarrow potential for better uncertainty estimates

Why Encode Structure in Extensions of LDA?

- Substantive reasons
 - Additional structure corresponds to substantively interesting content
 - Avoids potential ad-hoc secondary analysis
 - Clear data generating process
- Statistical reasons
 - **Smoothing** \rightsquigarrow borrow information across groups intelligently
 - **Uncertainty** \rightsquigarrow potential for better uncertainty estimates
 - **Improved topics** \rightsquigarrow small word conditions, structure could help

Plan for the Class

- 1) Discuss model with unknown covariates for strata proportions \rightsquigarrow presentational style
- 2) Discuss model with hierarchy of topics \rightsquigarrow mirrors genre

Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

Describe model that facilitates estimation of **presentational styles** in Senate press releases

Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

Describe model that facilitates estimation of **presentational styles** in Senate press releases

- Characterize representation provided to constituents

Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

Describe model that facilitates estimation of **presentational styles** in Senate press releases

- Characterize representation provided to constituents
- Divide attention over a set of topics

Unknown Covariates for Issue Attention: Measuring Attention in Senate Press Releases

Substantive problem:

Senators (representatives) regularly engage the public → presentational style

But we know little about this engagement

Why? **Hard to Measure**

Describe model that facilitates estimation of **presentational styles** in Senate press releases

- Characterize representation provided to constituents
- Divide attention over a set of topics
- Given attention to topics, write press releases

Presentational Styles \rightsquigarrow Objective Function

- $\pi_{itk} \equiv$ Attention senator i allocates to issue k in year t
- $\pi_{itk} \equiv$ Probability press release is about issue k
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Presentational Styles \rightsquigarrow Objective Function

- $\pi_{itk} \equiv$ Attention senator i allocates to issue k in year t
- $\pi_{itk} \equiv$ Probability press release is about issue k
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release j from senator i in year t)

Presentational Styles \rightsquigarrow Objective Function

- $\pi_{itk} \equiv$ Attention senator i allocates to issue k in year t
- $\pi_{itk} \equiv$ Probability press release is about issue k
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release j from senator i in year t)

- **Assume**: Each press release j assigned to one topic.
- Let τ_{ijt} indicate press release j 's topic.

Presentational Styles \rightsquigarrow Objective Function

- $\pi_{itk} \equiv$ Attention senator i allocates to issue k in year t
- $\pi_{itk} \equiv$ Probability press release is about issue k
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release j from senator i in year t)

- **Assume**: Each press release j assigned to one topic.
- Let τ_{ijt} indicate press release j 's topic.

$$\tau_{ijt} \sim \text{Multinomial}(1, \boldsymbol{\pi}_{it})$$

Presentational Styles \rightsquigarrow Objective Function

- $\pi_{itk} \equiv$ Attention senator i allocates to issue k in year t
- $\pi_{itk} \equiv$ Probability press release is about issue k
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release j from senator i in year t)

- **Assume**: Each press release j assigned to one topic.
- Let $\boldsymbol{\tau}_{ijt}$ indicate press release j 's topic.

$$\boldsymbol{\tau}_{ijt} \sim \text{Multinomial}(1, \boldsymbol{\pi}_{it})$$

- Conditional on topic, draw document's content.

Presentational Styles \rightsquigarrow Objective Function

- $\pi_{itk} \equiv$ Attention senator i allocates to issue k in year t
- $\pi_{itk} \equiv$ Probability press release is about issue k
- $\boldsymbol{\pi}_{it} = (\pi_{it1}, \dots, \pi_{it44})$

Press release-level parameters (press release j from senator i in year t)

- **Assume**: Each press release j assigned to one topic.
- Let τ_{ijt} indicate press release j 's topic.

$$\tau_{ijt} \sim \text{Multinomial}(1, \boldsymbol{\pi}_{it})$$

- Conditional on topic, draw document's content.
- If $\tau_{ijtk} = 1$ then

$$\mathbf{x}_{ijt} \sim \text{Multinomial}(n_{ijt}, \boldsymbol{\theta}_k).$$

Priors

Each π_{it} is a draw from one-of- S styles \rightsquigarrow mixture of Dirichlet distributions .

Priors

Each π_{it} is a draw from one-of- S styles \rightsquigarrow mixture of Dirichlet distributions .

$$\sigma_{it} \sim \text{Multinomial}(1, \beta).$$

Priors

Each π_{it} is a draw from one-of- S styles \rightsquigarrow mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s)\end{aligned}$$

Priors

Each π_{it} is a draw from one-of- S styles \rightsquigarrow mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

Priors

Each π_{it} is a draw from one-of- S styles \rightsquigarrow mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

Other priors:

Priors

Each π_{it} is a draw from one-of- S styles \rightsquigarrow mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

Other priors:

$$\theta_k \sim \text{Multinomial}(\lambda)$$

Priors

Each π_{it} is a draw from one-of- S styles \rightsquigarrow mixture of Dirichlet distributions .

$$\begin{aligned}\sigma_{it} &\sim \text{Multinomial}(1, \beta). \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

Other priors:

$$\begin{aligned}\theta_k &\sim \text{Multinomial}(\lambda) \\ \beta &\sim \text{Multinomial}(\mathbf{1})\end{aligned}$$

Presentational Styles \rightsquigarrow Objective Function

Presentation Styles \rightsquigarrow Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\boldsymbol{\lambda}) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1)\end{aligned}$$

Presentational Styles \rightsquigarrow Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta)\end{aligned}$$

Presentational Styles \rightsquigarrow Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta) \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s)\end{aligned}$$

Presentational Styles \rightsquigarrow Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta) \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \tau_{ijt} | \pi_{it} &\sim \text{Multinomial}(1, \pi_{it})\end{aligned}$$

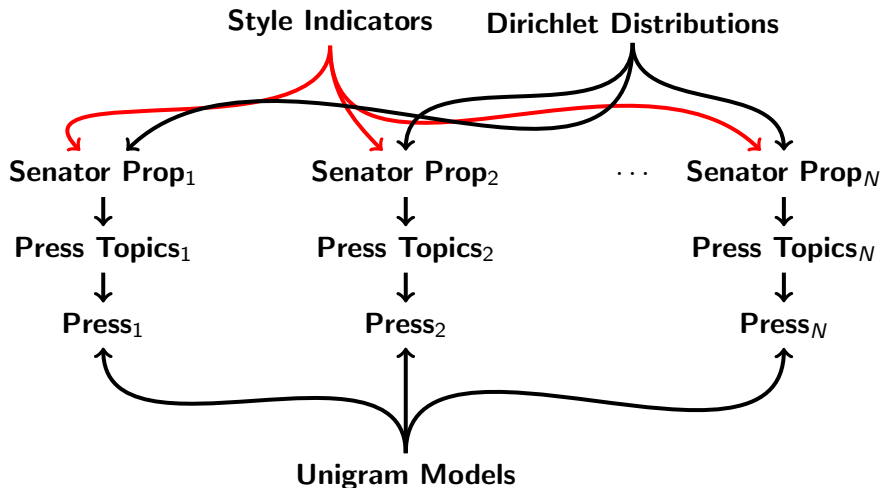
Presentational Styles \rightsquigarrow Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta) \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \tau_{ijt} | \pi_{it} &\sim \text{Multinomial}(1, \pi_{it}) \\ \mathbf{x}_{ijt} | \tau_{ijtk} = 1, \theta_k &\sim \text{Multinomial}(n_{ijt}, \theta_k)\end{aligned}$$

Presentational Styles \rightsquigarrow Objective Function

$$\begin{aligned}\beta &\sim \text{Dirichlet}(\mathbf{1}) \\ \theta_k &\sim \text{Dirichlet}(\lambda) \\ \alpha_{ks} &\sim \text{Gamma}(0.25, 1) \\ \sigma_{it} &\sim \text{Multinomial}(1, \beta) \\ \pi_{it} | \sigma_{its} = 1, \alpha_s &\sim \text{Dirichlet}(\alpha_s) \\ \tau_{ijt} | \pi_{it} &\sim \text{Multinomial}(1, \pi_{it}) \\ \mathbf{x}_{ijt} | \tau_{ijtk} = 1, \theta_k &\sim \text{Multinomial}(n_{ijt}, \theta_k)\end{aligned}$$

Mixture of Styles, Mixture of Topics



Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[\beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[\pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$

Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[\beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[\pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$

1) Estimate with Variational Approximation

Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[\beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[\pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$

- 1) Estimate with Variational Approximation
- 2) Determining number of clusters at top? (Grimmer, Shorey, Wallach, and Zlotnick, In Progress)

Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[\beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[\pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$

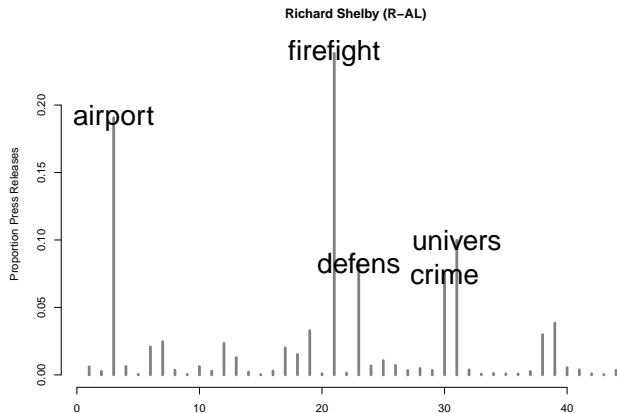
- 1) Estimate with Variational Approximation
- 2) Determining number of clusters at top? (Grimmer, Shorey, Wallach, and Zlotnick, In Progress)
 - Non-parametric model \rightsquigarrow statistical selection

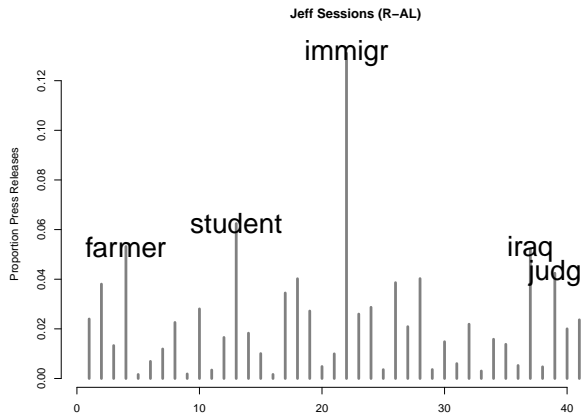
Posterior:

$$p(\alpha, \beta, \theta, \sigma, \pi, \tau | \mathbf{X}) \propto \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times$$

$$\prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[\beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[\pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}}$$

- 1) Estimate with Variational Approximation
- 2) Determining number of clusters at top? (Grimmer, Shorey, Wallach, and Zlotnick, In Progress)
 - Non-parametric model \rightsquigarrow statistical selection
 - Experiments/Coding Exercises to assess





Notions of validity: From Quinn, Monroe, et al (2010)

- **Semantic Validity:** All categories are coherent and meaningful
- **Convergent Construct Validity:** Measures concur with existing measures in critical details.
- **Discriminant Construct Validity:** Measures differ from existing measures in productive ways.
- **Predictive Measure:** Measures from the model corresponds to external events in expected ways.
- **Hypothesis Validity:** Measures generated from the model can be used to test substantive hypotheses.

To establish utility of new measures, demonstrate variety of **validations**

None of these validations are performed using a canned statistic

All: require substantive knowledge on areas (and what we expect!) [

Home Style Measures, Semantic Validity

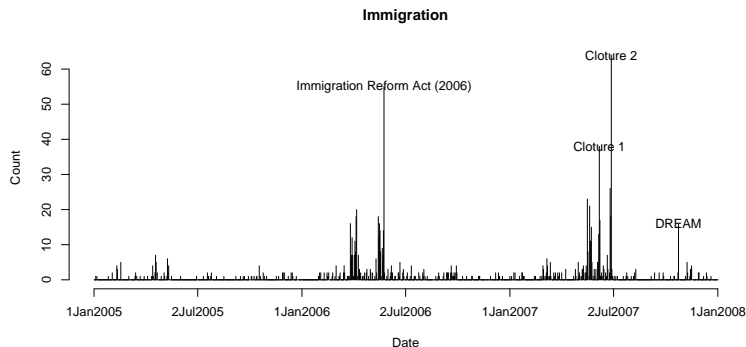
Must: Demonstrate to reader that topics are coherent and semantically meaningful

Description	Stems	%
Honorary	honor,prayer,rememb,fund,tribut	5.0
Transp. Grants	airport,transport,announc,urban,hud	4.8
Iraq	iraq,iraqi,troop,war,sectarian	4.7
DHS Policy	homeland,port,terrorist,dh,fema	4.1
Judicial Nom.	judg,court,suprem,nomin,nomine	3.8
Fire Dept. Grant	firefight,homeland,afgp,award,equip	3.7

How: **examples** in text are also useful.

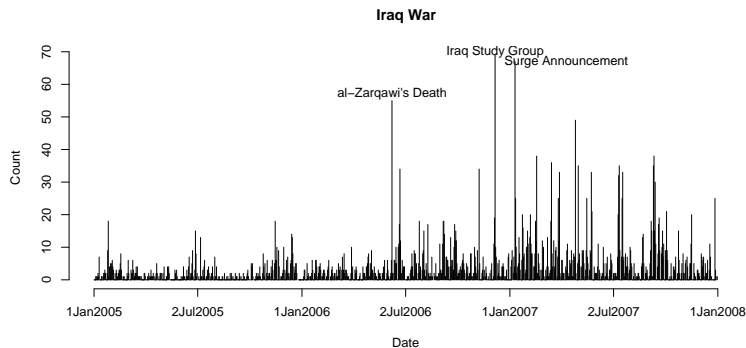
Home Style Measures, Convergent Validity

Over time variation



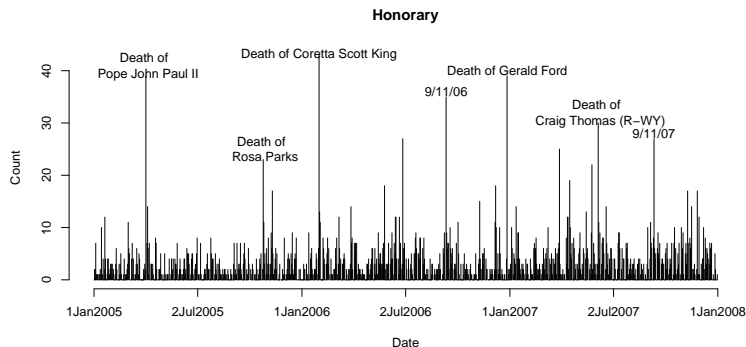
Home Style Measures, Convergent Validity

Over time variation



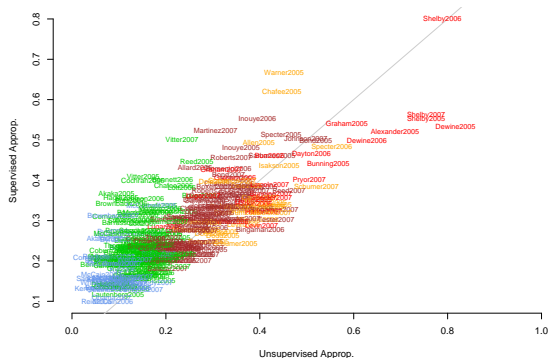
Home Style Measures, Convergent Validity

Over time variation

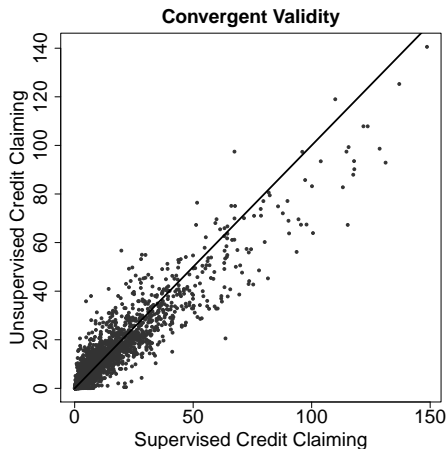


Home Style Measures, Convergent Validity

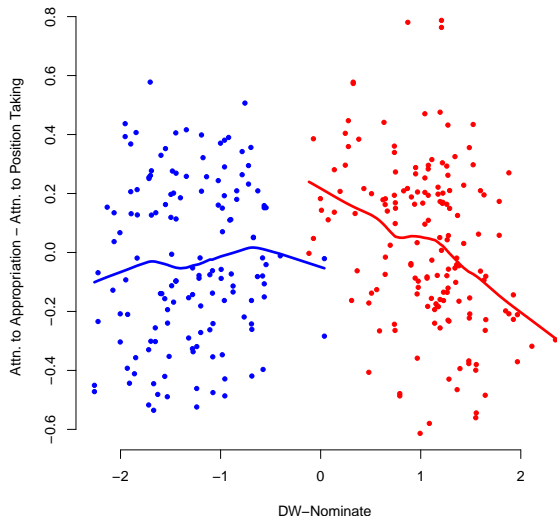
Supervised/Unsupervised Convergence



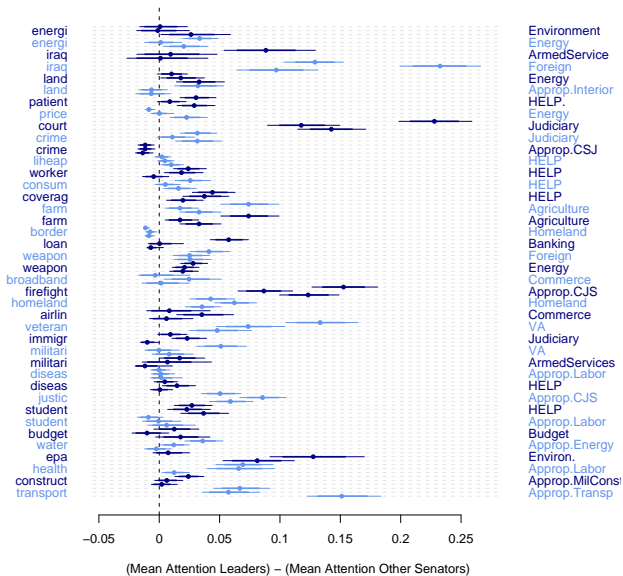
Home Style Measures, Convergent Validity



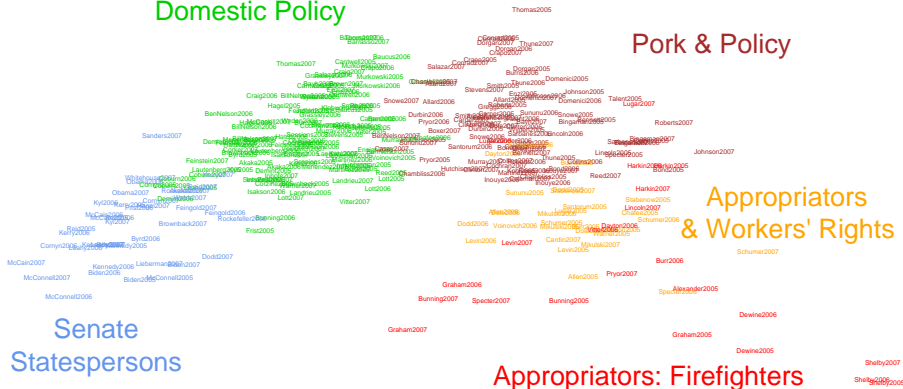
Discriminant Construct Validity



Predictive Validity



Domestic Policy



Hypothesis Validity

Domestic Policy

Pork & Policy

Appropriators & Workers' Rights

Appropriators: Firefighters

Senate Statespersons

Senate Statesperson

- Iraq War
- Intelligence
- Intl.
Relations

Hypothesis Validity

Domestic Policy

Pork & Policy

Appropriators & Workers' Rights

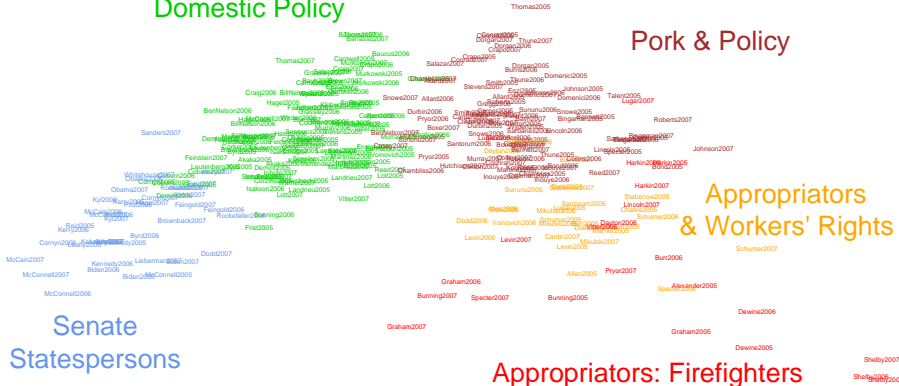
Appropriators: Firefighters

Senate Statespersons

Domestic Policy

- Iraq War
- Intelligence
- Intl.
Relations

- Environment
- Gas prices
- DHS
- Consumer



Hypothesis Validity

Domestic Policy

Pork & Policy

Appropriators & Workers' Rights

Appropriators: Firefighters

Senate Statespersons

Domestic Policy

Pork & Policy

- Iraq War
- Intelligence
- Intl. Relations

- Environment
- Gas prices
- DHS
- Consumer

- WRDA grants
- Farming
- Health Care
- Education

Hypothesis Validity

Domestic Policy

Pork & Policy

Appropriators & Workers' Rights

Appropriators: Firefighters

Senate Statespersons

Domestic Policy

Pork & Policy

Appropriators

- Iraq War
- Intelligence
- Intl. Relations

- Environment
- Gas prices
- DHS
- Consumer

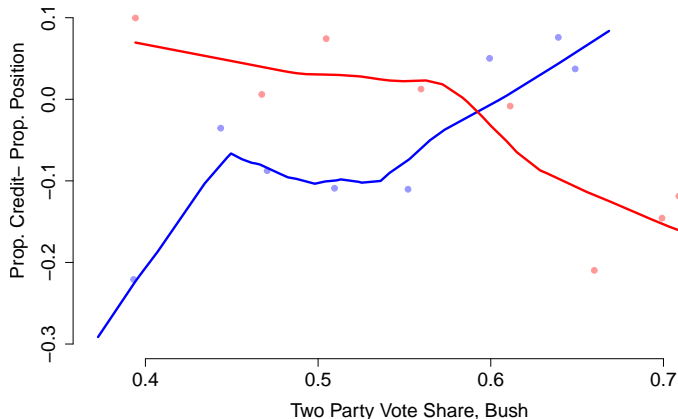
- WRDA grants
- Farming
- Health Care
- Education

- Fire Grants
- Airport Grants
- University
- Money

Hypothesis Validity

Why do senators adopt different styles?

District Fit



What are the right number of topics?

What are the right number of topics?

- Number of topics \rightsquigarrow depends on task at hand

What are the right number of topics?

- Number of topics \rightsquigarrow depends on task at hand
- Coarse \rightsquigarrow broad comparisons, lose distinctions

What are the right number of topics?

- Number of topics \rightsquigarrow depends on task at hand
- Coarse \rightsquigarrow broad comparisons, lose distinctions
- Granular \rightsquigarrow specific insights, lose broader picture

What are the right number of topics?

- Number of topics \rightsquigarrow depends on task at hand
- Coarse \rightsquigarrow broad comparisons, lose distinctions
- Granular \rightsquigarrow specific insights, lose broader picture
- **Hierarchy of topics** \rightsquigarrow Pachinko Allocation, Hierarchies of von-Mises Fisher Distributions

What are the right number of topics?

- Number of topics \rightsquigarrow depends on task at hand
- Coarse \rightsquigarrow broad comparisons, lose distinctions
- Granular \rightsquigarrow specific insights, lose broader picture
- **Hierarchy of topics** \rightsquigarrow Pachinko Allocation, Hierarchies of von-Mises Fisher Distributions

Blaydes, Grimmer, and McQueen 2018 \rightsquigarrow estimate nested topics to explore the **Mirrors for Princes**

The Mirrors Genre (BGM 2017)

26 Christian mirrors

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

21 Islamic texts

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

21 Islamic texts

- Advice on the Art of Governance (1612 CE)

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

Work with translations

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

Work with translations~> little evidence of selection

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

Work with translations~> little evidence of selection

- Collect data on collection of 98 (51 Christian, 47 Islamic, some not translated)

The Mirrors Genre (BGM 2017)

26 Christian mirrors

- The Prince (1513 CE)
- Advice to Justinian (527 CE)
- The Adventures of Telemachus (1699 CE)

21 Islamic texts

- Advice on the Art of Governance (1612 CE)
- Kalila wa Dimna (748 CE)
- The Sultan's Register of Laws (1632-1633 CE)

Work with translations~> little evidence of selection

- Collect data on collection of 98 (51 Christian, 47 Islamic, some not translated)
- No difference on Year/Region

Preprocessing Texts

47 books

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
 - almighty, god

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
 - almighty, god
 - monarch, prince, king, ruler

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
 - almighty, god
 - monarch, prince, king, ruler
 - Lord \neq lord

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
 - almighty, god
 - monarch, prince, king, ruler
 - Lord \neq lord

Result: short segment j in book i is a count vector

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
 - almighty, god
 - monarch, prince, king, ruler
 - Lord \neq lord

Result: short segment j in book i is a count vector

$$\mathbf{x}_{ij} = (x_{ij1}, x_{ij2}, \dots, x_{ij2124})$$

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
 - almighty, god
 - monarch, prince, king, ruler
 - Lord \neq lord

Result: short segment j in book i is a count vector

$$\mathbf{x}_{ij} = (x_{ij1}, x_{ij2}, \dots, x_{ij2124})$$

We work with a normalized version of the documents,

Preprocessing Texts

47 books \rightsquigarrow Each divided into paragraphs

Create feature space

- Bag of words, stem, discard punctuation, stop words
- Translate words left in Arabic (allah) and discard proper nouns
- Identified synonyms
 - almighty, god
 - monarch, prince, king, ruler
 - Lord \neq lord

Result: short segment j in book i is a count vector

$$\mathbf{x}_{ij} = (x_{ij1}, x_{ij2}, \dots, x_{ij2124})$$

We work with a normalized version of the documents,

$$\mathbf{x}_{ij}^* = \frac{\mathbf{x}_{ij}}{\sqrt{\mathbf{x}_{ij}' \mathbf{x}_{ij}}}$$

Measuring Themes in the Mirrors

Model built around two hierarchies:

Measuring Themes in the Mirrors

Model built around two hierarchies:

- 1) Books \rightsquigarrow paragraphs (Blei, Ng, Jordan 2003; Wallach, 2008; Quinn et al 2010; Grimmer 2010; Roberts et al 2014)

Measuring Themes in the Mirrors

Model built around two hierarchies:

- 1) Books \rightsquigarrow paragraphs (Blei, Ng, Jordan 2003; Wallach, 2008; Quinn et al 2010; Grimmer 2010; Roberts et al 2014)
- 2) Coarse topics \rightsquigarrow granular topics (Li and McCallum 2006; Gopal and Yang 2014)

Measuring Themes in the Mirrors

Estimate **four** quantities of interest

Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)

Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)
- 2) Coarse (broad) topics (3)

Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)
- 2) Coarse (broad) topics (3)
 - Each granular topic classified into one coarse topic

Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)
- 2) Coarse (broad) topics (3)
 - Each granular topic classified into one coarse topic
- 3) Each book i 's **themes** _{i}

$$\mathbf{themes}_i = (\text{theme}_{i1}, \text{theme}_{i2}, \dots, \text{theme}_{i60})$$

Measuring Themes in the Mirrors

Estimate **four** quantities of interest

- 1) Granular topics (60)
- 2) Coarse (broad) topics (3)
 - Each granular topic classified into one coarse topic
- 3) Each book *i*'s **themes**;
- 4) Each short segment's granular (and coarse) topic

A Hierarchy of Topics

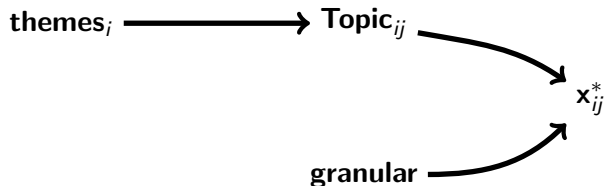
themes;

A Hierarchy of Topics

themes_{*i*} \longrightarrow Topic_{*ij*}

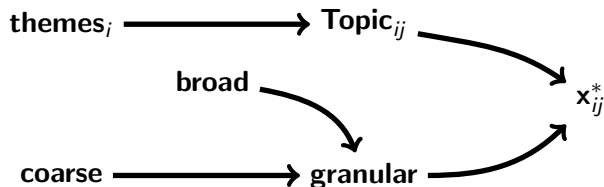
$$\mathbf{Topic}_{ij} \sim \text{Multinomial}(1, \mathbf{themes}_i)$$

A Hierarchy of Topics



$$\begin{aligned} \mathbf{Topic}_{ij} &\sim \text{Multinomial}(1, \mathbf{themes}_i) \\ \mathbf{x}_{ij}^* | \mathbf{Topic}_{ijk} = 1 &\sim \text{vMF}(\kappa, \mathbf{granular}_k) \end{aligned}$$

A Hierarchy of Topics



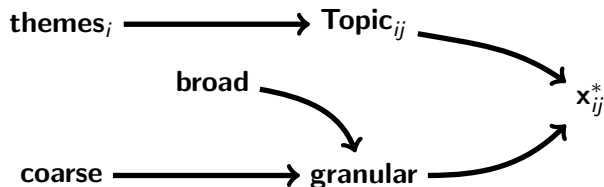
$$\mathbf{Topic}_{ij} \sim \text{Multinomial}(1, \mathbf{themes}_i)$$

$$\mathbf{x}_{ij}^* | \mathbf{Topic}_{ijk} = 1 \sim \text{vMF}(\kappa, \mathbf{granular}_k)$$

$$\mathbf{broad}_k \sim \text{Multinomial}(1, \mathbf{Broad\ Theme\ Prior})$$

$$\mathbf{granular}_k | \mathbf{broad}_{km} = 1 \sim \text{vMF}(\kappa, \mathbf{coarse}_m)$$

A Hierarchy of Topics



$$\mathbf{Topic}_{ij} \sim \text{Multinomial}(1, \mathbf{themes}_i)$$

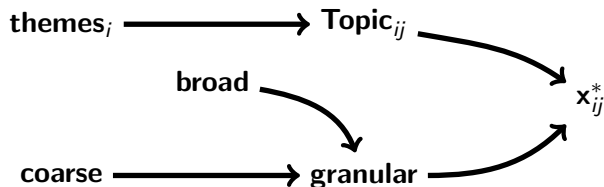
$$\mathbf{x}_{ij}^* | \mathbf{Topic}_{ijk} = 1 \sim \text{vMF}(\kappa, \mathbf{granular}_k)$$

$$\mathbf{broad}_k \sim \text{Multinomial}(1, \mathbf{Broad\ Theme\ Prior})$$

$$\mathbf{granular}_k | \mathbf{broad}_{km} = 1 \sim \text{vMF}(\kappa, \mathbf{coarse}_m)$$

Estimate model with Variational Approximation

A Hierarchy of Topics



$$\mathbf{Topic}_{ij} \sim \text{Multinomial}(1, \mathbf{themes}_i)$$

$$\mathbf{x}_{ij}^* | \mathbf{Topic}_{ijk} = 1 \sim \text{vMF}(\kappa, \mathbf{granular}_k)$$

$$\mathbf{broad}_k \sim \text{Multinomial}(1, \mathbf{Broad Theme Prior})$$

$$\mathbf{granular}_k | \mathbf{broad}_{km} = 1 \sim \text{vMF}(\kappa, \mathbf{coarse}_m)$$

Estimate model with Variational Approximation

Model selection: automatic model fit, qualitative evaluation

Interpreting Unsupervised Models

Two approaches to labeling output

Interpreting Unsupervised Models

Two approaches to labeling output

- 1) **Computational**: identify discriminating words

Interpreting Unsupervised Models

Two approaches to labeling output

- 1) **Computational**: identify discriminating words
- 2) **Manual**: Segments classified to coarse, granular topics. Read, discuss, and label

Interpreting Unsupervised Models

Two approaches to labeling output

- 1) **Computational**: identify discriminating words
- 2) **Manual**: Segments classified to coarse, granular topics. Read, discuss, and label

Unsupervised models **structure** and **guide** our reading

Art of Rulership

Practices and ideals of political rule

Art of Rulership

Practices and ideals of political rule

king

Art of Rulership

Practices and ideals of political rule

king, princ

Art of Rulership

Practices and ideals of political rule

king, princ, citi

Art of Rulership

Practices and ideals of political rule

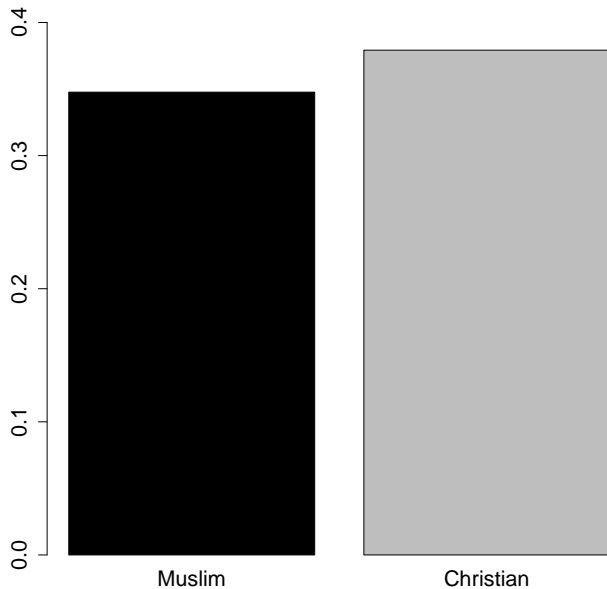
king, princ, citi, great, place, work, emperor, enemi, armi, letter

Art of Rulership

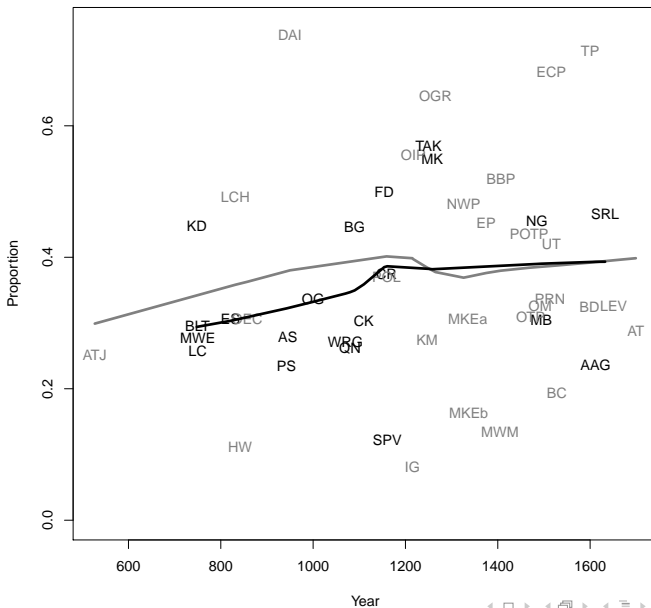
Practices and ideals of political rule

king, princ, citi, great, place, work, emperor, enemi, armi, letter

36.5% of paragraphs



Coarse Topic 1



Religion and Virtue

Connection between religion, virtue, justice and political rule

Religion and Virtue

Connection between religion, virtue, justice and political rule

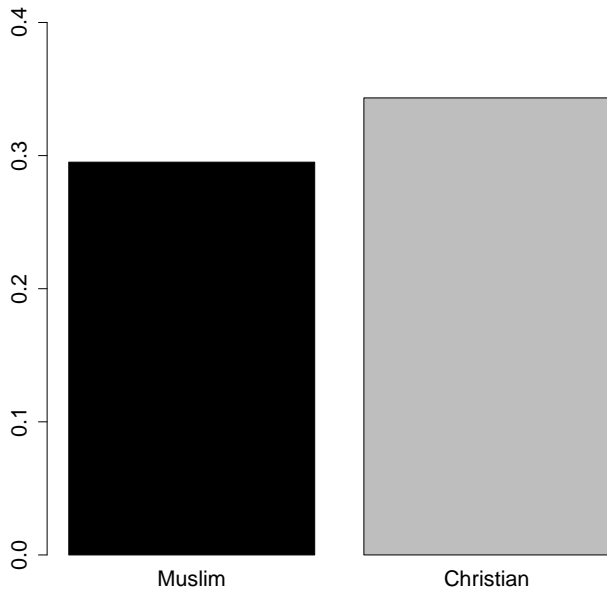
almighti,good,virtu,power,ruler,justic,prayer,rule,prophet,mena

Religion and Virtue

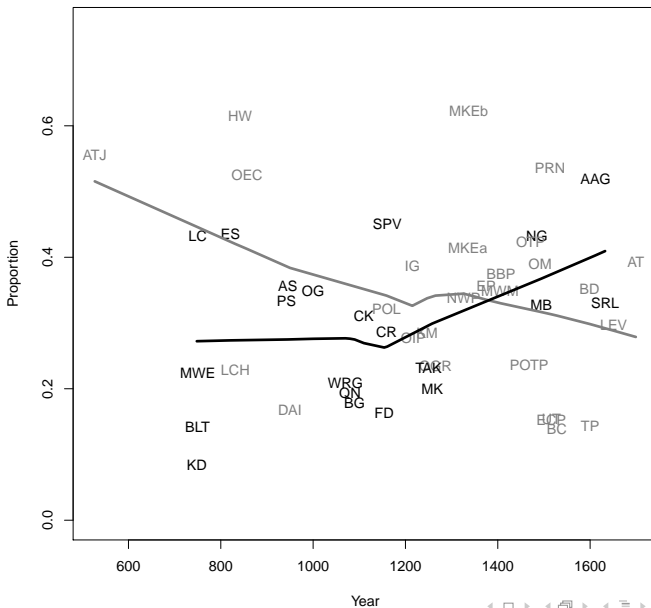
Connection between religion, virtue, justice and political rule

almighti,good,virtu,power,ruler,justic,prayer,rule,prophet,mena

32.2% of paragraphs



Coarse Topic 2



Inner Life of the Ruler

Personal relationships, care for and practices of the self, and ultimate fate of the soul

Inner Life of the Ruler

Personal relationships, care for and practices of the self, and ultimate fate of the soul

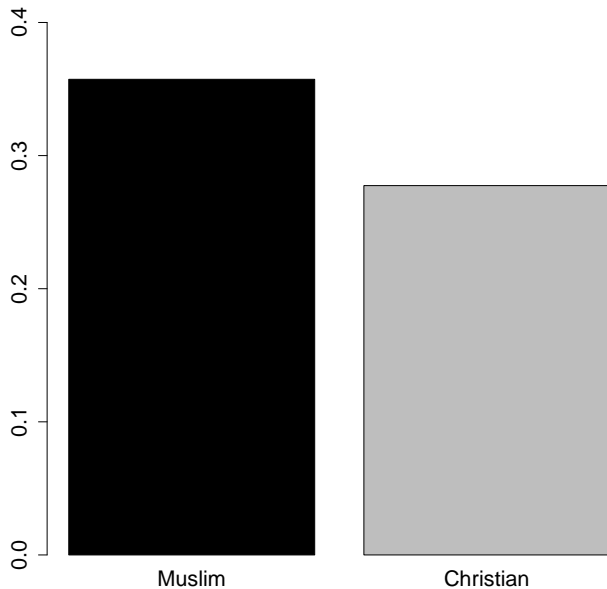
man,land,woman,know,bodi,eye,ladi,love,faculti,old

Inner Life of the Ruler

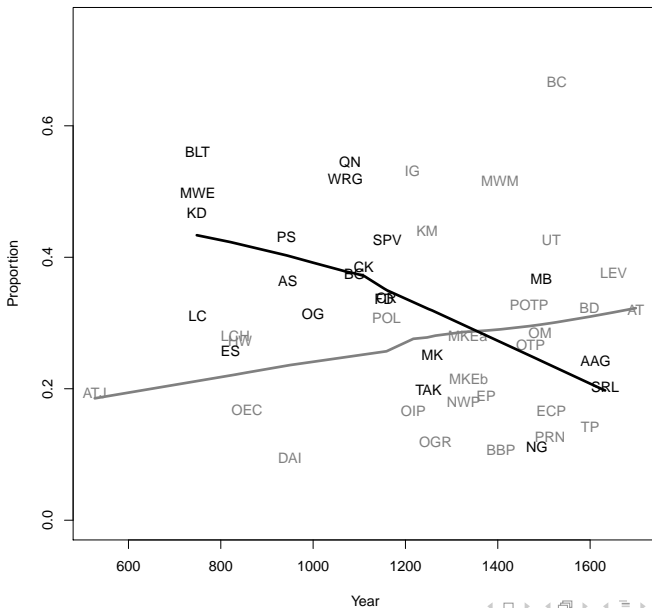
Personal relationships, care for and practices of the self, and ultimate fate of the soul

man,land,woman,know,bodi,eye,ladi,love,faculti,old

31.2% of paragraphs



Coarse Topic 3



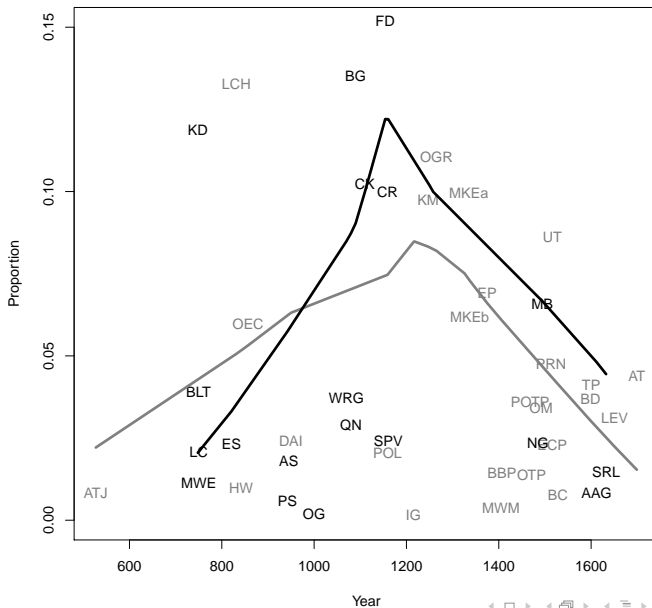
Granular: Best Practices for Ruling

king, princ, citi, great, place, work, emperor, enemi, armi, letter

king, kingdom, royal, minist, reign, father, court, majesti, presenc, war

6.2% of paragraphs

Coarse Topic 1 Granular Topic 1



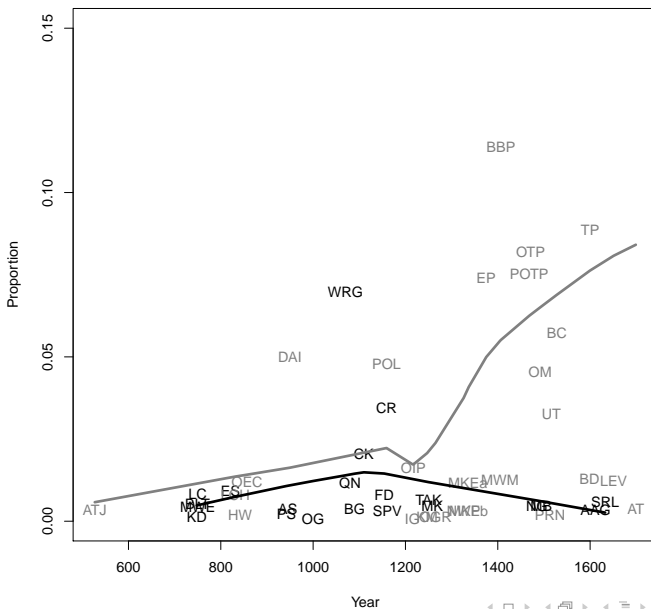
Granular: Characteristics that distinguish Just Ruler from Tyrant

king, princ, citi, great, place, work, emperor, enemi, armi, letter

king, kingdom, royal, minist, reign, father, court, majesti, presenc, war
princ, good, peopl, christian, tyranni, war, mind, ought, state, public

3.1% of paragraphs

Coarse Topic 1 Granular Topic 2



Granular: Religious Virtues and Political Ideals

almighti,good,virtu,power,ruler,justic,prayer,rule,prophet,mena

almighti,bless,grant,peac,messeng,prophet,merci,holi,command,grace

6.9% of paragraphs

Coarse Topic 2 Granular Topic 1

