Text as Data

Justin Grimmer

Associate Professor Department of Political Science University of Chicago

August 23rd, 2017

Discovery and Measurement

What is the research process? (Grimmer, Roberts, and Stewart 2017)

- 1) Discovery: a hypothesis or view of the world
- 2) Measurement according to some organization
- 3) Causal Inference: effect of some intervention

Text as data methods assist at each stage of research process

Measurement

Two approaches to measurement

- 1) Use an existing classification scheme to categorize documents
- 2) Simultaneously discover categories and measure prevalence (repurpose discovery methods)

Clustering

Document → One Cluster

Doc 1

Doc 2

Doc 3

:

Doc N

Cluster 1

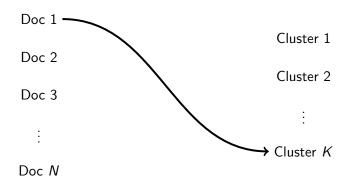
Cluster 2

:

Cluster K

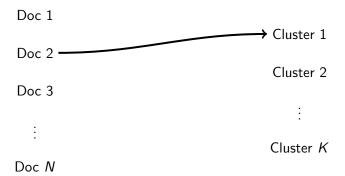
Clustering

Document → One Cluster



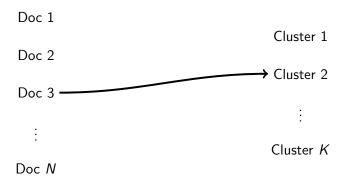
Clustering

Document → One Cluster



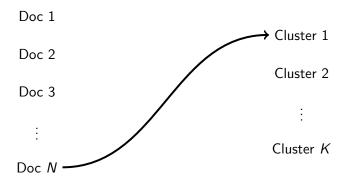
Clustering

Document → One Cluster

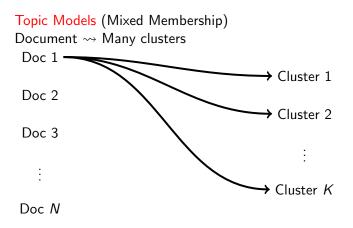


Clustering

Document ~> One Cluster



```
Topic Models (Mixed Membership)
Document → Many clusters
 Doc 1
                                        Cluster 1
 Doc 2
                                        Cluster 2
 Doc 3
                                       Cluster K
Doc N
```



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

required a mere 126 genes. It is 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

"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 genomes," explains Arcady Mushegian, a computational mo-

lecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

Genes

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

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

- Consider document i, (i = 1, 2, ..., N).

- Consider document i, (i = 1, 2, ..., N).
- Suppose there are M_i total words and \mathbf{x}_i is an $M_i \times 1$ vector, where \mathbf{x}_{im} describes the m^{th} word used in the document*.

- Consider document i, (i = 1, 2, ..., N).
- Suppose there are M_i total words and \mathbf{x}_i is an $M_i \times 1$ vector, where \mathbf{x}_{im} describes the m^{th} word used in the document*.

*Notice: this is a different representation than a document-term matrix. x_{im} is a number that says which of the J words are used. The difference is for clarity and we'll this representation is closely related to document-term matrix

- Consider document i, (i = 1, 2, ..., N).
- Suppose there are M_i total words and \mathbf{x}_i is an $M_i \times 1$ vector, where \mathbf{x}_{im} describes the m^{th} word used in the document*.

$$\pi_i | \alpha \sim \mathsf{Dirichlet}(\alpha)$$

- Consider document i, (i = 1, 2, ..., N).
- Suppose there are M_i total words and \mathbf{x}_i is an $M_i \times 1$ vector, where \mathbf{x}_{im} describes the m^{th} word used in the document*.

$$m{\pi}_i | m{lpha} \sim \mathsf{Dirichlet}(m{lpha})$$
 $m{ au}_{im} | m{\pi}_i \sim \mathsf{Multinomial}(1, m{\pi}_i)$

- Consider document i, (i = 1, 2, ..., N).
- Suppose there are M_i total words and \mathbf{x}_i is an $M_i \times 1$ vector, where \mathbf{x}_{im} describes the m^{th} word used in the document*.

$$egin{array}{ll} m{\pi}_i | m{lpha} & \sim & \mathsf{Dirichlet}(m{lpha}) \ m{ au}_{im} | m{\pi}_i & \sim & \mathsf{Multinomial}(1, m{\pi}_i) \ m{ imes}_{im} | m{ heta}_k, au_{imk} = 1 & \sim & \mathsf{Multinomial}(1, m{ heta}_k) \end{array}$$

- Consider document i, (i = 1, 2, ..., N).
- Suppose there are M_i total words and \mathbf{x}_i is an $M_i \times 1$ vector, where \mathbf{x}_{im} describes the m^{th} word used in the document*.

$$m{ heta}_k \sim \mathsf{Dirichlet}(\mathbf{1})$$
 $m{\pi}_i | m{lpha} \sim \mathsf{Dirichlet}(m{lpha})$ $m{ au}_{im} | m{\pi}_i \sim \mathsf{Multinomial}(1, m{\pi}_i)$ $m{x}_{im} | m{ heta}_k, au_{imk} = 1 \sim \mathsf{Multinomial}(1, m{ heta}_k)$

- Consider document i, (i = 1, 2, ..., N).
- Suppose there are M_i total words and \mathbf{x}_i is an $M_i \times 1$ vector, where \mathbf{x}_{im} describes the m^{th} word used in the document*.

$$egin{array}{ll} oldsymbol{ heta}_k & \sim & \mathsf{Dirichlet}(\mathbf{1}) \ lpha_k & \sim & \mathsf{Gamma}(lpha,eta) \ oldsymbol{\pi}_i | oldsymbol{lpha} & \sim & \mathsf{Dirichlet}(oldsymbol{lpha}) \ oldsymbol{ au}_{im} | oldsymbol{\pi}_i & \sim & \mathsf{Multinomial}(1,oldsymbol{\pi}_i) \ oldsymbol{x}_{im} | oldsymbol{ heta}_k, au_{imk} = 1 & \sim & \mathsf{Multinomial}(1,oldsymbol{ heta}_k) \end{array}$$

$$p(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \alpha | \boldsymbol{X}) \propto p(\alpha)p(\boldsymbol{\pi}|\alpha)p(\boldsymbol{T}|\boldsymbol{\pi})p(\boldsymbol{X}|\boldsymbol{\theta}, \boldsymbol{T})$$

$$egin{array}{ll} p(m{\pi},m{T},m{\Theta},lpha|m{X}) & \propto & p(m{lpha})p(m{\pi}|m{lpha})p(m{T}|m{\pi})p(m{X}|m{ heta},m{T}) \ & \propto & p(m{lpha})\prod_{i=1}^N \left[p(m{\pi}_i|m{lpha})\prod_{m=1}^{M_i}p(m{ au}_{im}|m{\pi})p(m{x}_{im}|m{ heta}_k, au_{imk}=1)
ight] \end{array}$$

$$\begin{split} \rho(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \boldsymbol{\alpha} | \boldsymbol{X}) & \propto & \rho(\boldsymbol{\alpha}) \rho(\boldsymbol{\pi} | \boldsymbol{\alpha}) \rho(\boldsymbol{T} | \boldsymbol{\pi}) \rho(\boldsymbol{X} | \boldsymbol{\theta}, \boldsymbol{T}) \\ & \propto & \rho(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\rho(\boldsymbol{\pi}_{i} | \boldsymbol{\alpha}) \prod_{m=1}^{M_{i}} \rho(\boldsymbol{\tau}_{im} | \boldsymbol{\pi}) \rho(\boldsymbol{x}_{im} | \boldsymbol{\theta}_{k}, \boldsymbol{\tau}_{imk} = 1) \right] \\ & \propto & \rho(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\frac{\Gamma(\sum_{k=1}^{K} \alpha_{k})}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \prod_{k=1}^{K} \prod_{m=1}^{M} \prod_{k=1}^{K} \left[\pi_{ik} \prod_{j=1}^{J} \theta_{jk}^{\boldsymbol{x}_{imj}} \right]^{\tau_{ikm}} \right] \end{split}$$

$$\rho(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \boldsymbol{\alpha} | \boldsymbol{X}) \propto \rho(\boldsymbol{\alpha}) \rho(\boldsymbol{\pi} | \boldsymbol{\alpha}) \rho(\boldsymbol{T} | \boldsymbol{\pi}) \rho(\boldsymbol{X} | \boldsymbol{\theta}, \boldsymbol{T}) \\
\propto \rho(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\rho(\boldsymbol{\pi}_{i} | \boldsymbol{\alpha}) \prod_{m=1}^{M_{i}} \rho(\boldsymbol{\tau}_{im} | \boldsymbol{\pi}) \rho(\boldsymbol{x}_{im} | \boldsymbol{\theta}_{k}, \boldsymbol{\tau}_{imk} = 1) \right] \\
\propto \rho(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\frac{\Gamma(\sum_{k=1}^{K} \alpha_{k})}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \prod_{k=1}^{K} \prod_{m=1}^{M} \prod_{k=1}^{K} \left[\pi_{ik} \prod_{j=1}^{J} \theta_{jk}^{\boldsymbol{x}_{imj}} \right]^{\tau_{ikm}} \right] \\$$

$$\begin{split} \rho(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \boldsymbol{\alpha} | \boldsymbol{X}) & \propto & \rho(\boldsymbol{\alpha}) \rho(\boldsymbol{\pi} | \boldsymbol{\alpha}) \rho(\boldsymbol{T} | \boldsymbol{\pi}) \rho(\boldsymbol{X} | \boldsymbol{\theta}, \boldsymbol{T}) \\ & \propto & \rho(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\rho(\boldsymbol{\pi}_{i} | \boldsymbol{\alpha}) \prod_{m=1}^{M_{i}} \rho(\boldsymbol{\tau}_{im} | \boldsymbol{\pi}) \rho(\boldsymbol{x}_{im} | \boldsymbol{\theta}_{k}, \boldsymbol{\tau}_{imk} = 1) \right] \\ & \propto & \rho(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\frac{\Gamma(\sum_{k=1}^{K} \alpha_{k})}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \prod_{k=1}^{K} \prod_{m=1}^{M} \prod_{k=1}^{K} \left[\pi_{ik} \prod_{j=1}^{J} \boldsymbol{\theta}_{jk}^{\boldsymbol{x}_{imj}} \right]^{\tau_{ikm}} \right] \end{split}$$

Together the model implies the following posterior:

$$p(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \boldsymbol{\alpha} | \boldsymbol{X}) \propto p(\boldsymbol{\alpha}) p(\boldsymbol{\pi} | \boldsymbol{\alpha}) p(\boldsymbol{T} | \boldsymbol{\pi}) p(\boldsymbol{X} | \boldsymbol{\theta}, \boldsymbol{T})$$

$$\propto p(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[p(\boldsymbol{\pi}_{i} | \boldsymbol{\alpha}) \prod_{m=1}^{M_{i}} p(\boldsymbol{\tau}_{im} | \boldsymbol{\pi}) p(\boldsymbol{x}_{im} | \boldsymbol{\theta}_{k}, \boldsymbol{\tau}_{imk} = 1) \right]$$

$$\propto p(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\frac{\Gamma(\sum_{k=1}^{K} \alpha_{k})}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \prod_{k=1}^{K} \prod_{m=1}^{M} \prod_{k=1}^{K} \left[\pi_{ik} \prod_{j=1}^{J} \theta_{jk}^{\boldsymbol{x}_{imj}} \right]^{\tau_{ikm}} \right]$$

Optimization:

Together the model implies the following posterior:

$$p(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \boldsymbol{\alpha} | \boldsymbol{X}) \propto p(\boldsymbol{\alpha}) p(\boldsymbol{\pi} | \boldsymbol{\alpha}) p(\boldsymbol{T} | \boldsymbol{\pi}) p(\boldsymbol{X} | \boldsymbol{\theta}, \boldsymbol{T})$$

$$\propto p(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[p(\boldsymbol{\pi}_{i} | \boldsymbol{\alpha}) \prod_{m=1}^{M_{i}} p(\boldsymbol{\tau}_{im} | \boldsymbol{\pi}) p(\boldsymbol{x}_{im} | \boldsymbol{\theta}_{k}, \boldsymbol{\tau}_{imk} = 1) \right]$$

$$\propto p(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\frac{\Gamma(\sum_{k=1}^{K} \alpha_{k})}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \prod_{k=1}^{K} \prod_{m=1}^{M} \prod_{k=1}^{K} \left[\pi_{ik} \prod_{j=1}^{J} \theta_{jk}^{\boldsymbol{x}_{imj}} \right]^{\tau_{ikm}} \right]$$

Optimization:

- Variational Approximation → Find "closest" distribution

Together the model implies the following posterior:

$$\rho(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \boldsymbol{\alpha} | \boldsymbol{X}) \propto \rho(\boldsymbol{\alpha}) \rho(\boldsymbol{\pi} | \boldsymbol{\alpha}) \rho(\boldsymbol{T} | \boldsymbol{\pi}) \rho(\boldsymbol{X} | \boldsymbol{\theta}, \boldsymbol{T}) \\
\propto \rho(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[p(\boldsymbol{\pi}_{i} | \boldsymbol{\alpha}) \prod_{m=1}^{M_{i}} p(\boldsymbol{\tau}_{im} | \boldsymbol{\pi}) p(\boldsymbol{x}_{im} | \boldsymbol{\theta}_{k}, \boldsymbol{\tau}_{imk} = 1) \right] \\
\propto \rho(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\frac{\Gamma(\sum_{k=1}^{K} \alpha_{k})}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \prod_{k=1}^{K} \prod_{m=1}^{M_{i}} \prod_{k=1}^{K} \left[\pi_{ik} \prod_{j=1}^{J} \theta_{jk}^{\boldsymbol{x}_{imj}} \right]^{\tau_{ikm}} \right] \\$$

Optimization:

- Variational Approximation → Find "closest" distribution
- Gibbs sampling \sim MCMC algorithm to approximate posterior

Together the model implies the following posterior:

$$p(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \boldsymbol{\alpha} | \boldsymbol{X}) \propto p(\boldsymbol{\alpha}) p(\boldsymbol{\pi} | \boldsymbol{\alpha}) p(\boldsymbol{T} | \boldsymbol{\pi}) p(\boldsymbol{X} | \boldsymbol{\theta}, \boldsymbol{T})$$

$$\propto p(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[p(\boldsymbol{\pi}_{i} | \boldsymbol{\alpha}) \prod_{m=1}^{M_{i}} p(\boldsymbol{\tau}_{im} | \boldsymbol{\pi}) p(\boldsymbol{x}_{im} | \boldsymbol{\theta}_{k}, \boldsymbol{\tau}_{imk} = 1) \right]$$

$$\propto p(\boldsymbol{\alpha}) \prod_{i=1}^{N} \left[\frac{\Gamma(\sum_{k=1}^{K} \alpha_{k})}{\prod_{k=1}^{K} \Gamma(\alpha_{k})} \prod_{k=1}^{K} \prod_{m=1}^{M} \prod_{k=1}^{K} \left[\pi_{ik} \prod_{j=1}^{J} \theta_{jk}^{\boldsymbol{x}_{imj}} \right]^{\tau_{ikm}} \right]$$

Optimization:

- Variational Approximation → Find "closest" distribution
- Gibbs sampling \infty MCMC algorithm to approximate posterior

Described in the slides appendix

Running a Topic Model with STM

to the STM Code

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

Where's the information for each word's topic? Reconsider document-term matrix

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
:	:	:	٠	:
DocN	0	1		1

Where's the information for each word's topic? Reconsider document-term matrix

	$Word_1$	Word ₂		ر Word
Doc ₁	0	1		0
Doc_2	2	0		3
:	:	÷	٠	:
Doc_{N}	0	1		1

Inner product of Documents (rows): $\mathbf{Doc}_{i}^{'}\mathbf{Doc}_{l}$

Where's the information for each word's topic? Reconsider document-term matrix

	$Word_1$	Word ₂		ر Word
Doc ₁	0	1		0
Doc_2	2	0		3
÷	:	÷	٠	:
$Doc_{\mathcal{N}}$	0	1		1

Inner product of Documents (rows): $\mathbf{Doc}_{i}^{'}\mathbf{Doc}_{l}$

Inner product of Terms (columns): $\mathbf{Word}_{j}'\mathbf{Word}_{k}$

Why does this work → 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
:	:	:	٠	:
$Doc_{\mathcal{N}}$	0	1		1

Inner product of Documents (rows): $\mathbf{Doc}_{i}^{'}\mathbf{Doc}_{l}$

Inner product of Terms (columns): $\mathbf{Word}_{j}'\mathbf{Word}_{k}$

Allows: measure of correlation of term usage across documents (heuristically: partition words, based on usage in documents)

Why does this work → Co-occurrence

Where's the information for each word's topic? Reconsider document-term matrix

	$Word_1$	Word ₂		ر Word
Doc ₁	0	1		0
Doc_2	2	0		3
:	:	:	٠	:
$Doc_{\mathcal{N}}$	0	1		1

Inner product of Documents (rows): $\mathbf{Doc}_{i}^{'}\mathbf{Doc}_{l}$

Inner product of Terms (columns): $\mathbf{Word}_{j}'\mathbf{Word}_{k}$

Allows: measure of correlation of term usage across documents (heuristically: partition words, based on usage in documents)

Latent Semantic Analysis: Reduce information in matrix using linear algebra (provides similar results, difficult to generalize)

4 D > 4 A > 4 E > 4 E > E 9 9 0

Why does this work → 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
÷	:	÷	٠	÷
Doc_{N}	0	1		1

Inner product of Documents (rows): $\mathbf{Doc}_{i}^{'}\mathbf{Doc}_{l}$

Inner product of Terms (columns): Word'_iWord_k

Allows: measure of correlation of term usage across documents (heuristically: partition words, based on usage in documents)

Latent Semantic Analysis: Reduce information in matrix using linear

algebra (provides similar results, difficult to generalize)

Biclustering: Models that partition documents and words simultaneously

$$p(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \alpha | \boldsymbol{X}) \propto p(\alpha) p(\boldsymbol{\pi} | \alpha) p(\boldsymbol{T} | \boldsymbol{\pi}) p(\boldsymbol{X} | \boldsymbol{\theta}, \boldsymbol{T})$$

$$p(\pi, T, \Theta, \alpha | X) \propto p(\alpha)p(\pi | \alpha)p(T | \pi)\underbrace{p(X | \theta, T)}_{1}$$

1) $\theta \leadsto$ Greater weight on terms that occur together

$$p(\pi, T, \Theta, \alpha | X) \propto p(\alpha)p(\pi | \alpha) \underbrace{p(T | \pi)}_{2} \underbrace{p(X | \theta, T)}_{1}$$

- 1) $\theta \leadsto$ Greater weight on terms that occur together
- 2) $\pi \leadsto$ Greater weight on indicators that appear more regularly

$$p(\boldsymbol{\pi}, \boldsymbol{T}, \boldsymbol{\Theta}, \alpha | \boldsymbol{X}) \propto p(\alpha) \underbrace{p(\boldsymbol{\pi} | \alpha)}_{3} \underbrace{p(\boldsymbol{T} | \boldsymbol{\pi})}_{2} \underbrace{p(\boldsymbol{X} | \boldsymbol{\theta}, \boldsymbol{T})}_{1}$$

- 1) $heta \leadsto ext{Greater weight on terms that occur together}$
- 2) $\pi \leadsto$ Greater weight on indicators that appear more regularly
- 3) $\alpha \leadsto \mathsf{Emphasis}$ on π with greater weight

Validation → Topic Intrusion

Wednesday → discussed several validations

- Labeling paragraphs
 - Identify separating words automatically
 - Label topics manually (read!)
- Statistical methods
 - 1) Entropy
 - 2) Exclusivity
 - 3) Cohesiveness
- Experiment Based Methods
 - Word intrusion → topic validity
 - Topic intrusion → model fit

Validation → Topic Intrusion

- 1) Ask research assistant to read paragraph
- 2) Construct experiment
 - For the document, select top three topics
 - Select a fourth topic
 - Show participant, ask her/him to identify intruder

Higher identification → topics are a better model of text

- Why is Japan revising its constitution?
- IR question: why is Japan now willing to engage militaristic foreign action?
- One explanation: election reform in 1993, changed electoral incentives
- To answer well: characterize campaigns across 50 + years
 - That sounds hard
 - That sounds impossible
- Determined (relentless) data collection
- Latent Dirichlet Allocation (on japanese texts)

Japanese Elections:

- Election Administration Commission runs elections \rightarrow district level

- Election Administration Commission runs elections \rightarrow district level
- Required to submit manifestos for all candidates to National Diet

Typical Manifesto:



- Election Administration Commission runs elections \rightarrow district level
- Required to submit manifestos for all candidates to National Diet
- Collected from 1950- 2009

- Election Administration Commission runs elections \rightarrow district level
- Required to submit manifestos for all candidates to National Diet
- Collected from 1950- 2009
 - Available only at district level

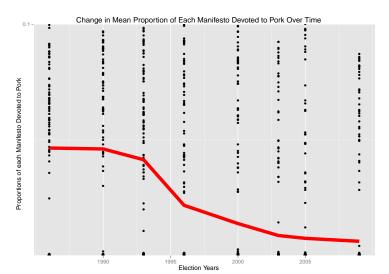
- Election Administration Commission runs elections \rightarrow district level
- Required to submit manifestos for all candidates to National Diet
- Collected from 1950- 2009
 - Available only at district level
 - Until: 2009 national library made texts available on microfilm

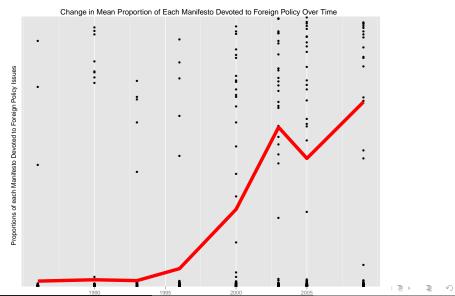
- Election Administration Commission runs elections \rightarrow district level
- Required to submit manifestos for all candidates to National Diet
- Collected from 1950- 2009
 - Available only at district level
 - Until: 2009 national library made texts available on microfilm
- Collected from microfilm, hand transcribed (no OCR worked), used a variety of techniques to create a TDM

- Election Administration Commission runs elections \rightarrow district level
- Required to submit manifestos for all candidates to National Diet
- Collected from 1950- 2009
 - Available only at district level
 - Until: 2009 national library made texts available on microfilm
- Collected from microfilm, hand transcribed (no OCR worked), used a variety of techniques to create a TDM
- Harder for Japanese

- Applies Vanilla LDA (using R Code)
- Output: topics (with Japanese characters)

Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic
改革	年金	推進	X	政治	日本
郵政	円	整備	政策	改革	玉
民営	廃止	図る	地域	国民	外交
小泉	改革	つとめる	まち	企業	国家
構造	兆	社会	鹿児島	自民党	社会
政府	実現	対策	全力	日本	国民
官	無駄	振興	選挙	共産党	保障
推進	日本	充実	国政	献金	安全
民	増税	促進	作り	金権	地域
自民党	削減	安定	横浜	党	拉致
日本	一元化	確立	対策	選挙	経済
制度	政権	企業	中小	禁止	守る
民間	子供	実現	発電	憲法	問題
年金	地域	中小	推進	腐敗	北朝魚
実現	ひと	育成	エネルギー	団体	教育
進める	サラリーマン	制度	企業	X	責任
断行	制度	政治	声	ソ連	力
地方	議員	地域	実現	守る	創る
止める	金	福祉	活性	平和	安心
保障	民主党	事業	自民党	円	目指す
財政	年間	改革	地方	反対	誇り
作る	一掃	確保	尽くす	真	憲法
贊成	郵政	強化	商店	是正	可能
社会	道路	教育	いかす	一掃	道
国民	交代	施設	全国	悪政	未来
公務員	社会保険庁	生活	政党	抜本	ひと
カ	月額	支援	ひと	定数	再生
経済	手当	環境	支援	政党	将来
55	談合	発展	経済	金丸	解決
安心	支援	施策	福祉	改悪	基本
Φ /0	义 恢	爬来	TH 71L	LQC 2005	基 个
Postal privatization	Reducing Wasteful Public Spending	Pork for the District	Policies for the district	Political Reform	Natio





REPRESENTATIONAL STYLE IN CONGRESS

What Legislators Say and Why It Matters

JUSTIN GRIMMER



Example 2: Automated Literature Reviews

Recall: literature reviews are hard to conduct LDA: developed (in part) to help structure JSTOR database Use JSTOR's research service to obtain data to analyze Question: How do scholars use classic text: Home Style Analysis: all articles that cite Home Style in JSTOR's data

Example 2: Automated Literature Reviews

Output: topic estimates

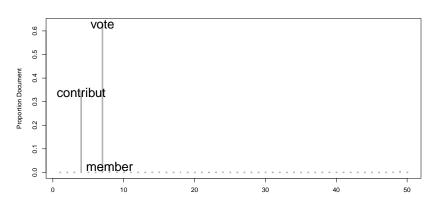
- Obtain $\log \theta_k$ from model
- One method to summarize a topic:
 - $\exp(\log \theta_k)$ (select 10-20 biggest words)
 - $\exp(\log \theta_k)$ Average_{$j \neq k$} $\exp(\log \theta_j)$ (select 10-20 biggest words)

Example 2: Automated Literature Reviews

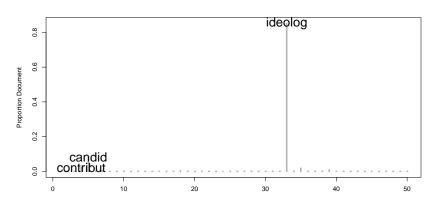
Example topics:

-	Label	Stems	Proportio
Ī	Life Style	member, district, attent, congress, time, cohort, retir	0.03
	Comp.Home	constitu,mp,member,parti,role,local,british	0.02
	Casework	casework, district, constitu, variabl, staff, congression, fiorina	0.03
	Votes	vote, variabl, model, estim, measur, legisl, constitu	0.04
	ld. Shirk	ideolog,vote,shirk,constitu,parti,senat,voter	0.03
	C. letters	mail,govern, activ,respond,commun,offic	0.02

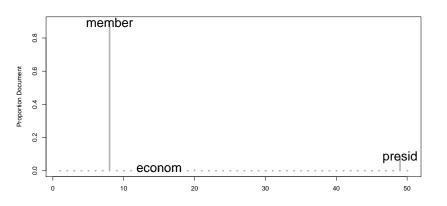
Wawro (2001) "A Panel Probit Analysis of Campaign Contributions and Roll Call Votes"



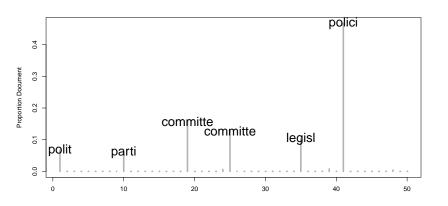
Bender (1996) "Legislator Voting and Shirking A Critical Review of the Literature"

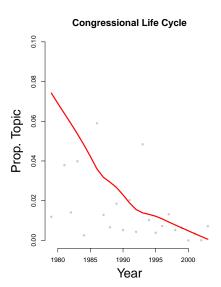


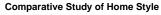
Parker (1980) "Cycles in Congressional District Attention"

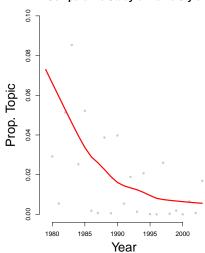


Shepsle (1985) "Policy Consequences of Government by Congressional Subcommittees"

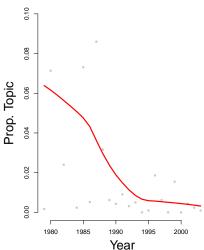




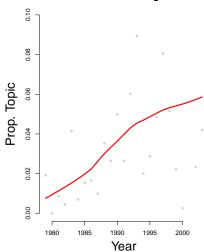


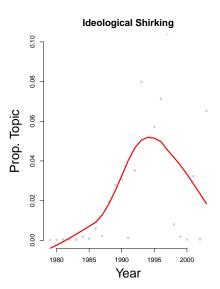


Casework and the Incumbency Advantage



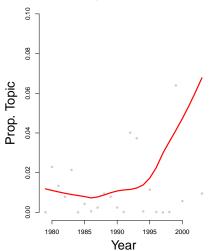
Causes of Roll Call Voting Decisions





History of Home Style

Biases in Congressional Communication



• The IMPRESSION of INFLUENCE

Legislator Communication, Representation, and Democratic Accountability

JUSTIN GRIMMER SEAN J. WESTWOOD SOLOMON MESSING What legislators claim (Grimmer, Westwood, Messing 2014)

Labels Key Words Proportion

	8	
Labels	Key Words	Proportion
Requested appropriations	bill,funding,house,million,appropriations	0.08

"Dave Camp announced today that he was able to secure \$2.5 million for widening M-72 from US-31 easterly 7.2 miles to Old M-72. The bill will now head to the Senate for consideration...We have two more hurdles to clear to make sure the money is in the bill when it hits the President's desk: a vote in the Senate and a conference committee" (Camp, 2005)

	<u> </u>	
Labels	Key Words	Proportion
Requested appropriations	bill,funding,house,million,appropriations	0.08

"Congressman Doc Hastings has boosted federal funding for work on the Columbia Basin water supply for next year. Hastings has added \$400,000 for work on the Odessa Subaquifer, which when combined with the funding in the President's budget request, totals \$1 million for Fiscal Year 2009"... "Hastings' funding for the Odessa Subaquifer and Potholes Reservoir was included in the Fiscal Year 2009 Energy and Water Appropriations bill which was approved today by the full House Appropriations Committee. (Hastings, 2008)"

Labels	Key Words	Proportion
Requested appropriations	bill,funding,house,million,appropriations	0.08
Fire department grants	fire, grant, department, program, fire fighters	0.08

"Maurice Hinchey (D-NY) today announced that the West Endicott Fire Company has been awarded a \$17,051 federal grant to purchase approximately 10 sets of protective clothing, as well as radio equipment and air packs for its volunteer firefighters" (Hinchey, 2008)

	01	
Labels	Key Words	Proportion
Requested appropriations	bill,funding,house,million,appropriations	0.08
Fire department grants	fire, grant, department, program, fire fighters	0.08

"Congressman Pete Visclosky today announced that the Crown Point Fire Department will receive a \$16,550 Department of Homeland Security (DHS) grant to purchase a modular portable video system" (Visclosky, 2008)

Labels	Key Words	Proportion
Requested appropriations	bill,funding,house,million,appropriations	0.08
Fire department grants	fire,grant,department,program,firefighters	0.08
Stimulus	recovery,funding,jobs,information, act,	0.06

	01	
Labels	Key Words	Proportion
Requested appropriations	bill,funding,house,million,appropriations	0.08
Fire department grants	fire,grant,department,program,firefighters	0.08
Stimulus	recovery,funding,jobs,information, act,	0.06
Transportation	transportation, project, airport, transit, million	0.06

Correlated Topic Models

Dirichlet distribution → Assumes negative covariance between topics Logistic Normal Distribution → Allows some positive covariance between topics

$$egin{array}{lll} oldsymbol{ heta}_k & \sim & \mathsf{Dirichlet}(\mathbf{1}) \ oldsymbol{\eta}_i | oldsymbol{\mu}, oldsymbol{\Sigma} & \sim & \mathsf{Multivariate} \; \mathsf{Normal}(oldsymbol{\mu}, oldsymbol{\Sigma}) \ oldsymbol{\pi}_i & = & \dfrac{\exp{(oldsymbol{\eta}_i)}}{\sum_{k=1}^K \exp{(oldsymbol{\eta}_{ik})}} \ oldsymbol{ au}_{im} | oldsymbol{\pi}_i & \sim & \mathsf{Multinomial}(1, oldsymbol{\pi}_i) \ oldsymbol{x}_{im} | oldsymbol{ heta}_k, au_{imk} = 1 & \sim & \mathsf{Multinomial}(1, oldsymbol{ heta}_k) \end{array}$$

LDA Revisited

$$egin{array}{ll} m{ heta}_k & \sim & \mathsf{Dirichlet}(\mathbf{1}) \\ m{\pi}_i | m{lpha} & \sim & \mathsf{Dirichlet}(m{lpha}) \\ m{ au}_{im} | m{\pi}_i & \sim & \mathsf{Multinomial}(\mathbf{1}, m{\pi}_i) \\ m{ imes}_{im} | m{ heta}_k, au_{imk} = \mathbf{1} & \sim & \mathsf{Multinomial}(\mathbf{1}, m{ heta}_k) \end{array}$$

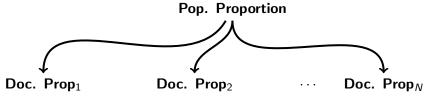
LDA Revisited

```
\begin{array}{lll} \textbf{Unigram Model}_k & \sim & \mathsf{Dirichlet}(\mathbf{1}) \\ & \textbf{Doc. Prop}_i & \sim & \mathsf{Dirichlet}(\textbf{Pop. Proportion}) \\ & \textbf{Word Topic}_{im} & \sim & \mathsf{Multinomial}(1, \textbf{Doc. Prop}_i) \\ & & \mathsf{Word}_{im} & \sim & \mathsf{Multinomial}(1, \textbf{Unigram Model}_k) \end{array}
```

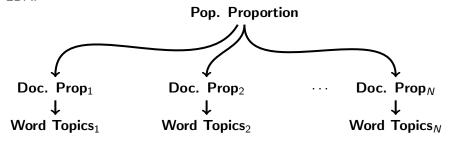
LDA:

Pop. Proportion

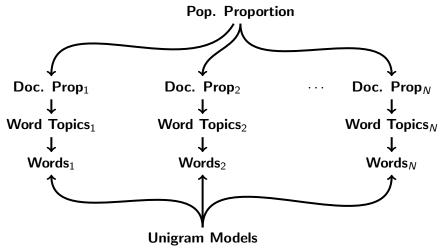
LDA:



LDA:



LDA:



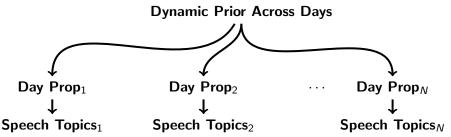
Dynamic Topic Model (Quinn et al 2010)

Dynamic Prior Across Days

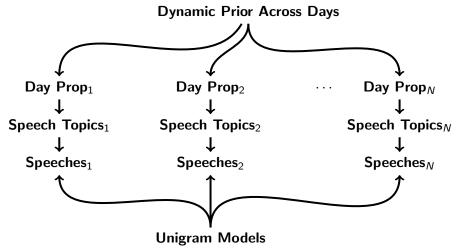
Dynamic Topic Model (Quinn et al 2010)

Day $\operatorname{\mathsf{Prop}}_1$ Day $\operatorname{\mathsf{Prop}}_2$ \cdots Day $\operatorname{\mathsf{Prop}}_N$

Dynamic Topic Model (Quinn et al 2010)



Dynamic Topic Model (Quinn et al 2010)



Expressed Agenda Model (Grimmer 2010)

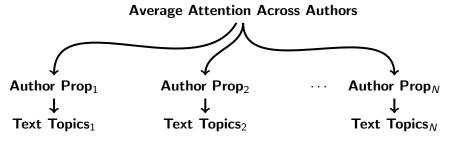
Average Attention Across Authors

Expressed Agenda Model (Grimmer 2010)

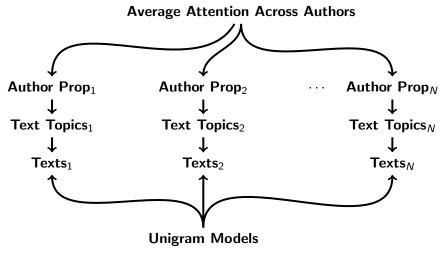
Average Attention Across Authors



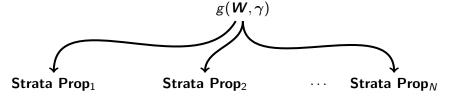
Expressed Agenda Model (Grimmer 2010)

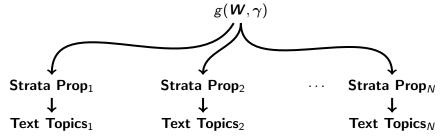


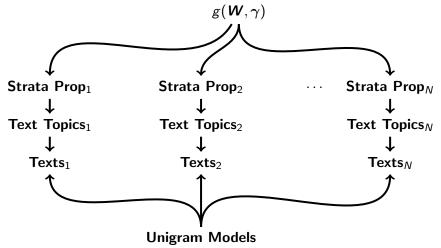
Expressed Agenda Model (Grimmer 2010)

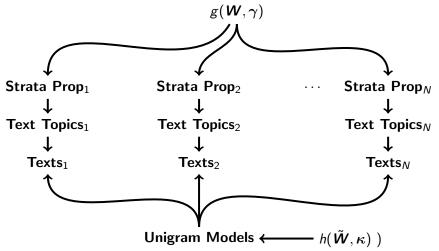


Structural Topic Model (Roberts, Stewart, Airoldi 2014) $g({m W}, {m \gamma})$







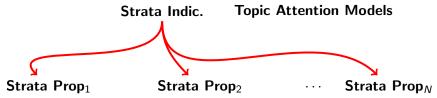


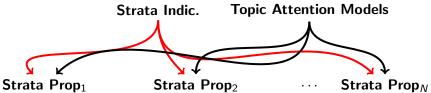
Conditioning on Unknown Covariates → levels of mixtures at proportions (Grimmer 2013; Wallach 2008)

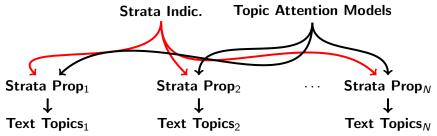
Mixture of Top. Attn. Models

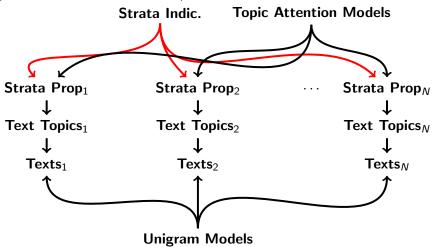
Conditioning on Unknown Covariates → levels of mixtures at proportions (Grimmer 2013; Wallach 2008)

Strata Indic. Topic Attention Models



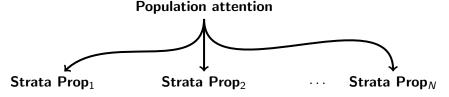


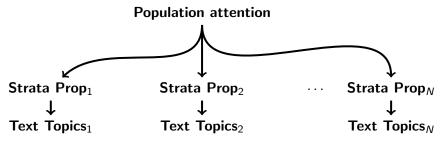


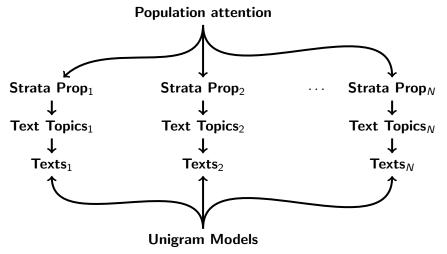


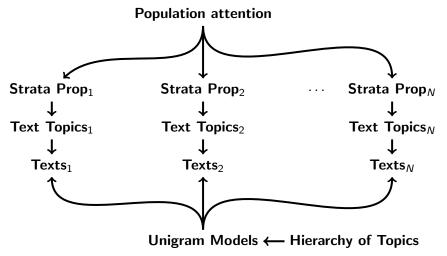
Conditioning on Unknown Covariates for Topics → hierarchy of topics (Li and McCallum 2006; Blaydes, Grimmer, and McQueen 2017)

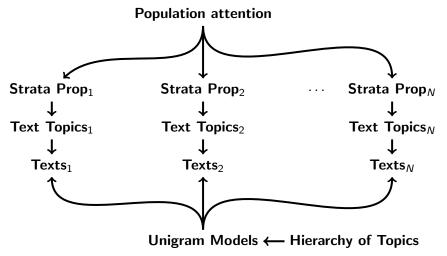
Population attention











- Substantive reasons

- Substantive reasons
 - Additional structure corresponds to substantively interesting content

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

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

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

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

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

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

Plan for the Class

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

Substantive problem:

Substantive problem:

Senators (representatives) regularly engage the public \rightarrow presentational style

But we know little about this engagement

Substantive problem:

Senators (representatives) regularly engage the public \rightarrow presentational style

But we know little about this engagement

Why? Hard to Measure

Substantive problem:

Senators (representatives) regularly engage the public \rightarrow 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

Substantive problem:

Senators (representatives) regularly engage the public \rightarrow 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

Substantive problem:

Senators (representatives) regularly engage the public \rightarrow 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

Substantive problem:

Senators (representatives) regularly engage the public \rightarrow 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

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

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

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

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

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

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

$$oldsymbol{ au}_{ijt} \sim \mathsf{Multinomial}(1, oldsymbol{\pi}_{it})$$

- $\pi_{itk} \equiv$ 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 au_{ijt} indicate press release j's topic.

$$oldsymbol{ au}_{ijt} \sim \mathsf{Multinomial}(1, oldsymbol{\pi}_{it})$$

- Conditional on topic, draw document's content.

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

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

$$oldsymbol{ au}_{\mathit{ijt}} \sim \mathsf{Multinomial}(1, oldsymbol{\pi}_{\mathit{it}})$$

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

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

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

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

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

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

$$egin{aligned} oldsymbol{\sigma}_{it} & \sim & \mathsf{Multinomial}(1,oldsymbol{eta}). \ oldsymbol{\pi}_{it} | \sigma_{its} = 1, lpha_s & \sim & \mathsf{Dirichlet}(lpha_s) \end{aligned}$$

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

$$egin{aligned} oldsymbol{\sigma_{it}} & \sim & \mathsf{Multinomial}(1,oldsymbol{eta}). \ oldsymbol{\pi_{it}} | \sigma_{its} = 1, oldsymbol{lpha_s} & \sim & \mathsf{Dirichlet}(oldsymbol{lpha_s}) \ lpha_{\mathit{ks}} & \sim & \mathsf{Gamma}(0.25,1) \end{aligned}$$

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

$$egin{array}{ll} oldsymbol{\sigma_{it}} & \sim & \mathsf{Multinomial}(1,oldsymbol{eta}). \ oldsymbol{\pi_{it}} | \sigma_{its} = 1, oldsymbol{lpha_s} & \sim & \mathsf{Dirichlet}(oldsymbol{lpha_s}) \ lpha_{\mathit{ks}} & \sim & \mathsf{Gamma}(0.25,1) \end{array}$$

Other priors:

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

$$egin{array}{ll} oldsymbol{\sigma_{it}} & \sim & \mathsf{Multinomial}(1,oldsymbol{eta}). \ oldsymbol{\pi_{it}} | \sigma_{its} = 1, lpha_s & \sim & \mathsf{Dirichlet}(lpha_s) \ lpha_{\mathit{ks}} & \sim & \mathsf{Gamma}(0.25,1) \end{array}$$

Other priors:

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

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

$$egin{array}{ll} oldsymbol{\sigma}_{it} & \sim & \mathsf{Multinomial}(1,eta). \ oldsymbol{\pi}_{it} | \sigma_{its} = 1, lpha_s & \sim & \mathsf{Dirichlet}(lpha_s) \ lpha_{ks} & \sim & \mathsf{Gamma}(0.25,1) \end{array}$$

Other priors:

$$egin{array}{ll} m{ heta}_k & \sim & \mathsf{Multinomial}(m{\lambda}) \\ m{eta} & \sim & \mathsf{Multinomial}(m{1}) \end{array}$$

Presentational Styles -- Objective Function

Presentational Styles Objective Function

```
eta \sim {\sf Dirichlet}({f 1}) \ eta_k \sim {\sf Dirichlet}({f \lambda}) \ lpha_{k{\sf s}} \sim {\sf Gamma}(0.25,1)
```

Presentational Styles --> Objective Function

```
eta \sim 	ext{Dirichlet}(\mathbf{1})
eta_k \sim 	ext{Dirichlet}(oldsymbol{\lambda})
lpha_{ks} \sim 	ext{Gamma}(0.25, 1)
oldsymbol{\sigma}_{it} \sim 	ext{Multinomial}(1, oldsymbol{eta})
```

Presentational Styles Objective Function

$$eta \sim ext{Dirichlet}(\mathbf{1})$$
 $eta_k \sim ext{Dirichlet}(oldsymbol{\lambda})$
 $lpha_{ks} \sim ext{Gamma}(0.25, 1)$
 $oldsymbol{\sigma}_{it} \sim ext{Multinomial}(1, oldsymbol{eta})$
 $oldsymbol{\pi}_{it} | \sigma_{its} = 1, lpha_s \sim ext{Dirichlet}(lpha_s)$

Presentational Styles Objective Function

$$eta \sim ext{Dirichlet}(\mathbf{1})$$
 $eta_k \sim ext{Dirichlet}(oldsymbol{\lambda})$
 $lpha_{ks} \sim ext{Gamma}(0.25, 1)$
 $oldsymbol{\sigma}_{it} \sim ext{Multinomial}(1, oldsymbol{eta})$
 $oldsymbol{\pi}_{it} | \sigma_{its} = 1, lpha_s \sim ext{Dirichlet}(lpha_s)$
 $oldsymbol{ au}_{ijt} | oldsymbol{\pi}_{it} \sim ext{Multinomial}(1, oldsymbol{\pi}_{it})$

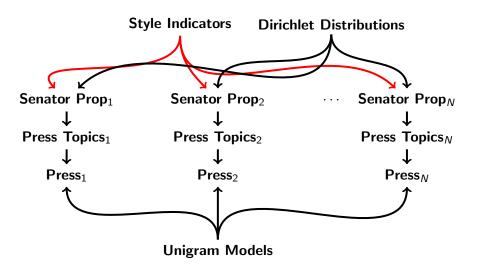
Presentational Styles --> Objective Function

$$eta \sim ext{Dirichlet}(\mathbf{1})$$
 $eta_k \sim ext{Dirichlet}(\lambda)$
 $lpha_{ks} \sim ext{Gamma}(0.25, 1)$
 $eta_{it} \sim ext{Multinomial}(1, eta)$
 $eta_{it} | \sigma_{its} = 1, lpha_s \sim ext{Dirichlet}(lpha_s)$
 $eta_{ijt} | \pi_{it} \sim ext{Multinomial}(1, \pi_{it})$
 $oldsymbol{x}_{ijt} | \tau_{ijtk} = 1, eta_k \sim ext{Multinomial}(n_{ijt}, eta_k)$

Presentational Styles --> Objective Function

$$eta \sim ext{Dirichlet}(\mathbf{1})$$
 $eta_k \sim ext{Dirichlet}(\lambda)$
 $lpha_{ks} \sim ext{Gamma}(0.25, 1)$
 $eta_{it} \sim ext{Multinomial}(1, eta)$
 $eta_{it} | \sigma_{its} = 1, lpha_s \sim ext{Dirichlet}(lpha_s)$
 $eta_{ijt} | \pi_{it} \sim ext{Multinomial}(1, \pi_{it})$
 $oldsymbol{x}_{ijt} | \tau_{ijtk} = 1, eta_k \sim ext{Multinomial}(n_{ijt}, eta_k)$

Mixture of Styles, Mixture of Topics



$$\begin{split} \rho(\alpha,\beta,\theta,\sigma,\pi,\tau|\mathbf{X}) & \quad \propto \quad \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times \\ & \quad \prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[\beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[\pi_{itk} \prod_{w=1}^W \theta_{kw}^{\lambda_{ijtw}}\right]^{\tau_{ijtk}}\right]^{\sigma_{its}} \end{split}$$

$$\begin{split} \rho(\alpha,\beta,\theta,\sigma,\pi,\tau|\textbf{\textit{X}}) & \quad \propto \quad \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_w-1} \times \\ & \quad \prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[\beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^D \prod_{k=1}^K \left[\pi_{itk} \prod_{w=1}^W \theta_{kw}^{x_{ijtw}}\right]^{\tau_{ijtk}}\right]^{\sigma_{its}} \end{split}$$

1) Estimate with Variational Approximation

$$\begin{split} \rho(\boldsymbol{\alpha},\boldsymbol{\beta},\boldsymbol{\theta},\boldsymbol{\sigma},\boldsymbol{\pi},\boldsymbol{\tau}|\boldsymbol{X}) & \propto & \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_W-1} \times \\ & \prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[\beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[\pi_{itk} \prod_{w=1}^W \theta_{kw}^{\lambda_{ijtw}}\right]^{\tau_{ijtk}}\right]^{\sigma_{its}} \end{split}$$

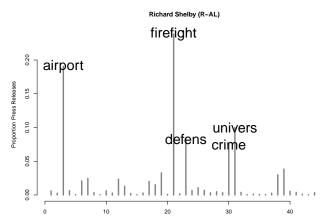
- 1) Estimate with Variational Approximation
- Determining number of clusters at top? (Grimmer, Shorey, Wallach, and Zlotnick, In Progress)

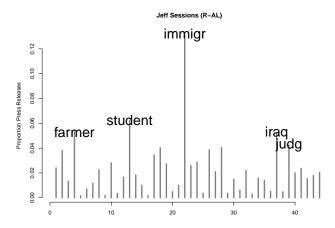
$$\begin{split} \rho(\boldsymbol{\alpha},\boldsymbol{\beta},\boldsymbol{\theta},\boldsymbol{\sigma},\boldsymbol{\pi},\boldsymbol{\tau}|\boldsymbol{X}) & \propto & \prod_{k=1}^{K} \prod_{s=1}^{S} \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^{W} \lambda_{w})}{\prod_{w=1}^{W} \Gamma(\lambda_{w})} \prod_{w=1}^{W} \theta_{k,w}^{\lambda_{w}-1} \times \\ & \prod_{i=1}^{N} \prod_{t=2005}^{2007} \prod_{s=1}^{S} \left[\beta_{s} \frac{\Gamma(\sum_{k=1}^{K} \alpha_{ks})}{\prod_{k=1}^{K} \Gamma(\alpha_{ks})} \prod_{k=1}^{K} \pi_{iks}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^{K} \left[\pi_{itk} \prod_{w=1}^{W} \theta_{kw}^{\lambda_{ijtw}} \right]^{\tau_{ijtk}} \right]^{\sigma_{its}} \end{split}$$

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

$$\begin{split} \rho(\boldsymbol{\alpha},\boldsymbol{\beta},\boldsymbol{\theta},\boldsymbol{\sigma},\boldsymbol{\pi},\boldsymbol{\tau}|\boldsymbol{X}) & \propto & \prod_{k=1}^K \prod_{s=1}^S \frac{\exp(-\frac{\alpha_{ks}}{1/4})}{1/4} \times \frac{\Gamma(\sum_{w=1}^W \lambda_w)}{\prod_{w=1}^W \Gamma(\lambda_w)} \prod_{w=1}^W \theta_{k,w}^{\lambda_W-1} \times \\ & \prod_{i=1}^N \prod_{t=2005}^{2007} \prod_{s=1}^S \left[\beta_s \frac{\Gamma(\sum_{k=1}^K \alpha_{ks})}{\prod_{k=1}^K \Gamma(\alpha_{ks})} \prod_{k=1}^K \pi_{itk}^{\alpha_{ks}-1} \prod_{j=1}^{D_{it}} \prod_{k=1}^K \left[\pi_{itk} \prod_{w=1}^W \theta_{kw}^{\lambda_{ijtw}}\right]^{\tau_{ijtk}}\right]^{\sigma_{its}} \end{split}$$

- Estimate with Variational Approximation
- Determining number of clusters at top? (Grimmer, Shorey, Wallach, and Zlotnick, In Progress)
 - Non-parametric model → 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!) [

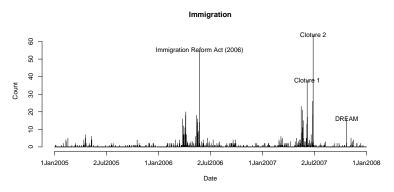
Home Style Measures, Semantic Validity

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

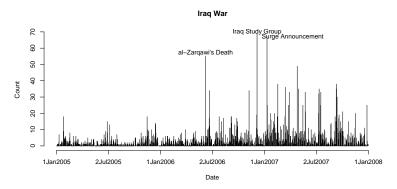
Stems	%
honor, prayer, rememb, fund, tribut	5.0
airport,transport,announc,urban,hud	
iraq,iraqi,troop,war,sectarian	4.7
homeland,port,terrorist,dh,fema	4.1
judg,court,suprem,nomin,nomine	
firefight, homeland, afgp, award, equip	3.7
	honor,prayer,rememb,fund,tribut airport,transport,announc,urban,hud iraq,iraqi,troop,war,sectarian homeland,port,terrorist,dh,fema judg,court,suprem,nomin,nomine

How: examples in text are also useful.

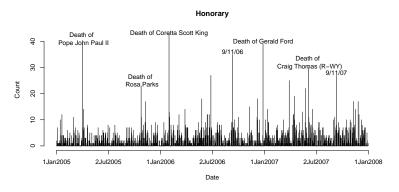
Over time variation



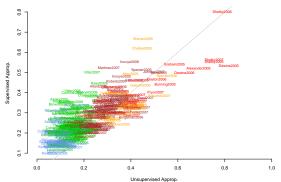
Over time variation

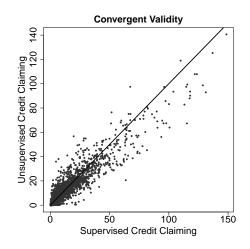


Over time variation

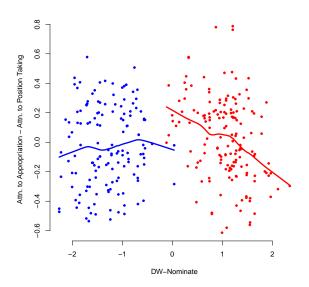


Supervised/Unsupervised Convergence

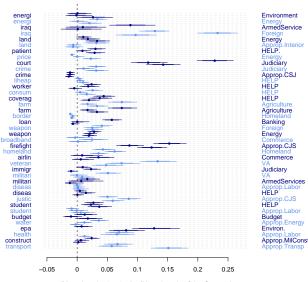




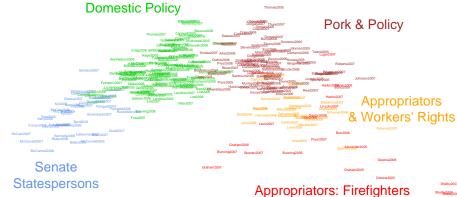
Discriminant Construct Validity

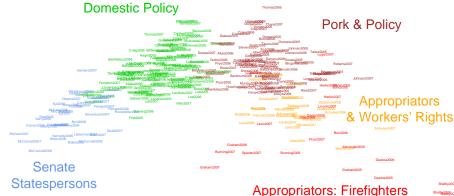


Predictive Validity



(Mean Attention Leaders) - (Mean Attention Other Senators)

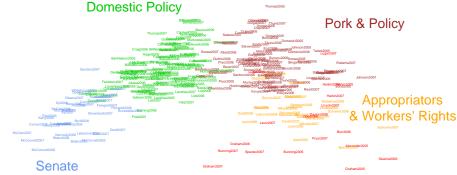




Senate

Statesperson

- Iraq War
- Intelligence
- Intl. Relations



Statespersons

Appropriators: Firefighters

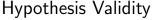
Senate Domestic
Statesperson Policy
- Iraq War - Envir

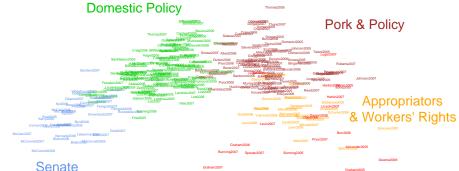
- Environment

Consumer

- Intelligence - Gas prices

- Intl. - DHS Relations Cons



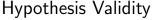


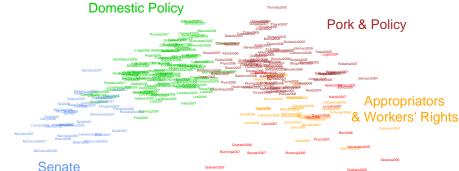
Statespersons

Appropriators: Firefighters Domestic Senate Pork & Policy Statesperson Policy - WRDA - Environment - Iraq War grants - Intelligence - Gas prices - Farming Intl. DHS Health Care

Consumer

Relations



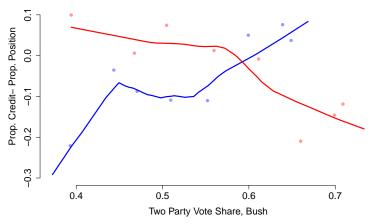


Ctotoonoroon
Statespersons

		Appropriators, Firefighters	
Senate	Domestic	Pork & Policy	Appropriators
Statesperson	Policy	- WRDA	- Fire Grants
- Iraq War	- Environment	grants	- Airport
- Intelligence	- Gas prices	- Farming	Grants
- Intl.	- DHS	- Health Care	- University
Relations	Consumer	- Education 🗗	> ∢ ≣ > Mēnev≣ 900

Consumer

Why do senators adopt different styles? District Fit



- Number of topics → depends on task at hand

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

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

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

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

Blaydes, Grimmer, and McQueen 2017→ estimate nested topics to explore the Mirros for Princes

26 Christian mirrors

26 Christian mirrors

- The Prince (1513 CE)

26 Christian mirrors

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

26 Christian mirrors

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

26 Christian mirrors

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

21 Islamic texts

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)

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)

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)

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

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

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)

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

47 books

47 books → Each divided into paragraphs

47 books → Each divided into paragraphs Create feature space

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

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

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

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

- 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

- 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

47 books → 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

47 books \leadsto 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

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

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,

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,

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

Model built around two hierarchies:

Model built around two hierarchies:

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

Model built around two hierarchies:

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

Estimate four quantities of interest

1) Granular topics (60)

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

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

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

```
themes<sub>i</sub> = (theme<sub>i1</sub>, theme<sub>i2</sub>, ..., theme<sub>i60</sub>)
```

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

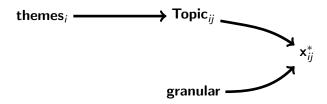
A Hierarchy of Topics

themes;

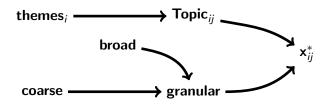
A Hierarchy of Topics

themes_i
$$\longrightarrow$$
 Topic_{ij}

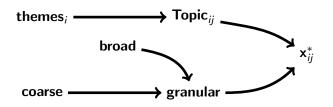
 $\mathsf{Topic}_{ij} \sim \mathsf{Multinomial}(1, \mathsf{themes}_i)$



$$egin{array}{ll} {\sf Topic}_{ij} & \sim & {\sf Multinomial}(1, {\sf themes}_i) \ {\sf x}_{ij}^* | {\sf Topic}_{ijk} = 1 & \sim & {\sf vMF}(\kappa, {\sf granular}_k) \end{array}$$

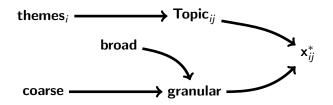


$$\begin{aligned} & \textbf{Topic}_{ij} & \sim & \mathsf{Multinomial}(1, \textbf{themes}_i) \\ & \textbf{x}_{ij}^* | \mathsf{Topic}_{ijk} = 1 & \sim & \mathsf{vMF}(\kappa, \textbf{granular}_k) \\ & \textbf{broad}_k & \sim & \mathsf{Multinomial}(1, \textbf{Broad Theme Prior}) \\ & \textbf{granular}_k | \mathsf{broad}_{km} = 1 & \sim & \mathsf{vMF}(\kappa, \textbf{coarse}_m) \end{aligned}$$



$$\begin{aligned} & \textbf{Topic}_{ij} & \sim & \text{Multinomial}(1, \textbf{themes}_i) \\ & \textbf{x}_{ij}^* | \text{Topic}_{ijk} = 1 & \sim & \text{vMF}(\kappa, \textbf{granular}_k) \\ & \textbf{broad}_k & \sim & \text{Multinomial}(1, \textbf{Broad Theme Prior}) \\ & \textbf{granular}_k | \text{broad}_{km} = 1 & \sim & \text{vMF}(\kappa, \textbf{coarse}_m) \end{aligned}$$

Estimate model with Variational Approximation



$$\begin{aligned} & \textbf{Topic}_{ij} & \sim & \text{Multinomial}(1, \textbf{themes}_i) \\ & \textbf{x}_{ij}^* | \text{Topic}_{ijk} = 1 & \sim & \text{vMF}(\kappa, \textbf{granular}_k) \\ & \textbf{broad}_k & \sim & \text{Multinomial}(1, \textbf{Broad Theme Prior}) \\ & \textbf{granular}_k | \text{broad}_{km} = 1 & \sim & \text{vMF}(\kappa, \textbf{coarse}_m) \end{aligned}$$

Estimate model with Variational Approximation Model selection: automatic model fit, qualitative evaluation

Two approaches to labeling output

Two approaches to labeling output

1) Computational: identify discriminating words

Two approaches to labeling output

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

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

Practices and ideals of political rule

Practices and ideals of political rule

king

Practices and ideals of political rule

king,princ

Practices and ideals of political rule

king,princ,citi

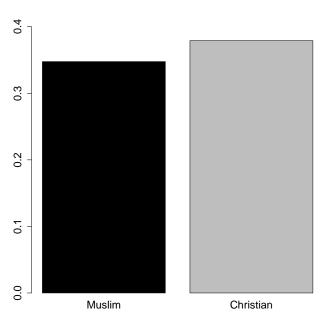
Practices and ideals of political rule

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

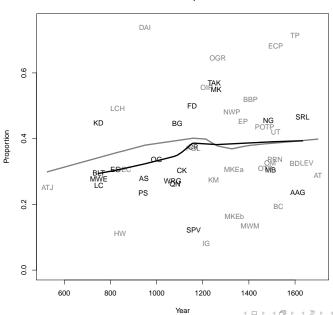
Practices and ideals of political rule

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

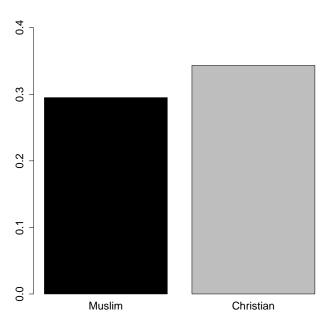
almighti, good, virtu, power, ruler, justic, prayer, rule, prophet, mena

Religion and Virtue

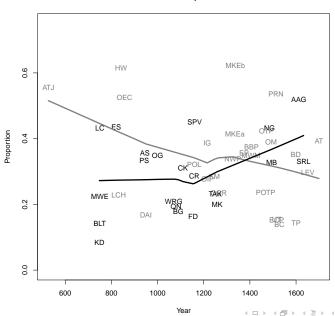
Connection between religion, virtue, justice and political rule

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

32.2% of pargraphs



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

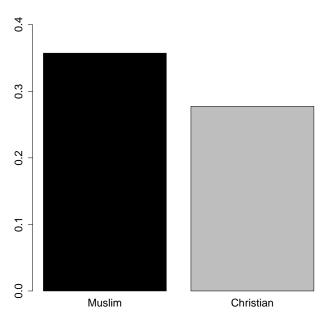
man,land,woman,know,bodi,eye,ladi,love,faculti,old

Inner Life of the Ruler

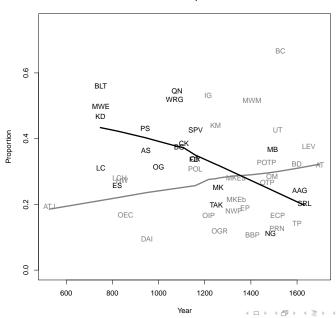
Personal relationships, care for and practices of the self, and ultimate fate of the soul

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

31.2% of paragraphs



Coarse Topic 3

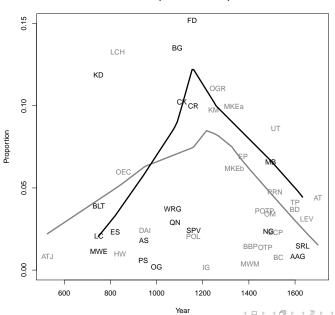


Granular: Best Practices for Ruling

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

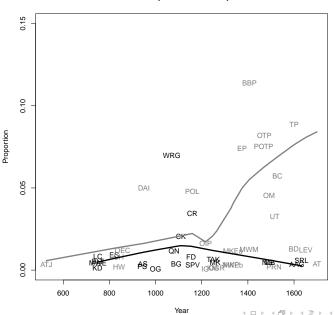


Granular: Characteristics that distinguish Just Ruler from Tyrant

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

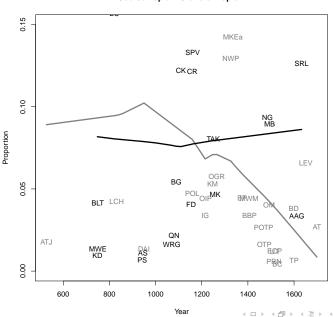


Granular: Religious Virtues and Political Ideals

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



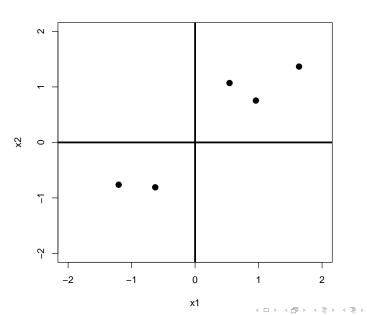
Principal Component Analysis

A Simple Two-Dimensional Example

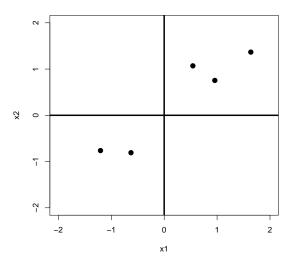
Suppose we have the following observations:

$$x_1 = (0.54, 1.07)$$

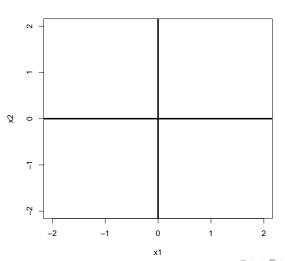
 $x_2 = (-1.20, -0.76)$
 $x_3 = (-0.63, -0.81)$
 $x_4 = (0.96, 0.75)$
 $x_5 = (1.64, 1.37)$



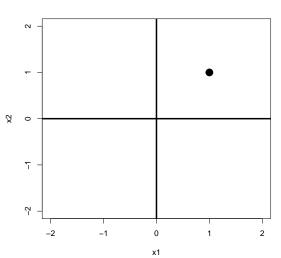
Goal: find line that summarizes bivariate information



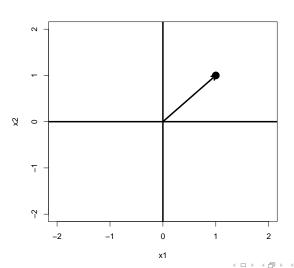
Suppose $\mathbf{w}_1 = (1,1)$



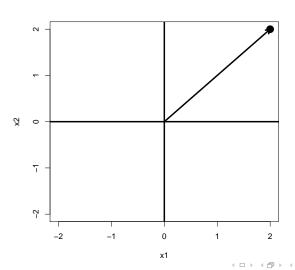
Suppose $\mathbf{w}_1 = (1,1)$



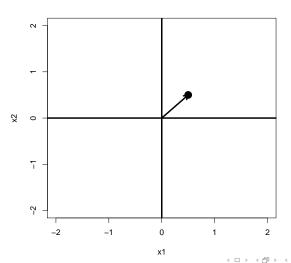
Suppose $\mathbf{w}_1 = (1,1)$



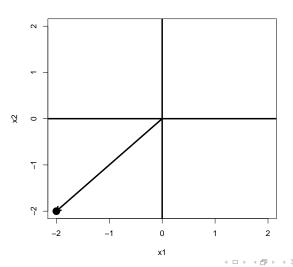
Suppose $\mathbf{w}_1 = (1,1) \ 2\mathbf{w}_1 = (2,2)$



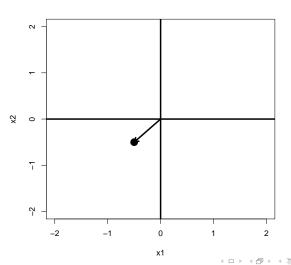
Suppose
$$\mathbf{w}_1 = (1,1) \ \frac{1}{2} \mathbf{w}_1 = (1/2,1/2)$$



Suppose
$$\mathbf{w}_1 = (1,1) -2\mathbf{w}_1 = (-2,-2)$$

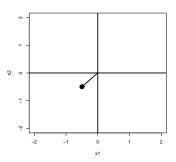


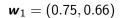
Suppose
$$\mathbf{w}_1 = (1,1)$$
 $-\frac{1}{2}\mathbf{w}_1 = (-1/2, -1/2)$

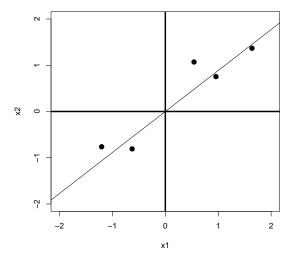


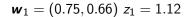
Suppose $\mathbf{w}_1 = (1,1)$

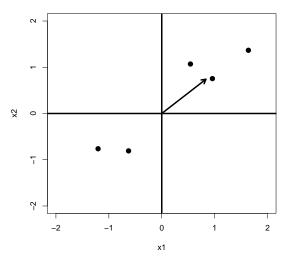
 $z_i = \text{amount we shrink/flip } \mathbf{w}_1 \text{ to approximate point } i.$



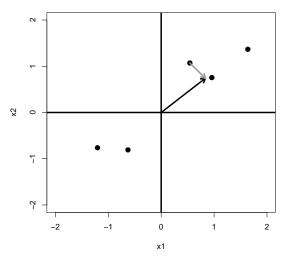




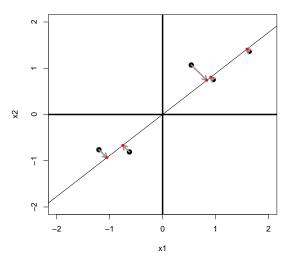




$$\mathbf{w}_1 = (0.75, 0.66) z_1 = 1.12$$



$$\mathbf{w}_1 = (0.75, 0.66) \ z_1 = 1.12$$



$$\mathbf{x}_i = \mathbf{z}_i \mathbf{w}_1 + \mathbf{e}_i$$

$$\mathbf{x}_i = z_i \mathbf{w}_1 + \mathbf{e}_i$$

 $(x_{i1}, x_{i2}) = (z_i w_{11} + e_{i1}, z_i w_{12} + e_{i2})$

$$\mathbf{x}_i = z_i \mathbf{w}_1 + \mathbf{e}_i$$

 $(x_{i1}, x_{i2}) = (z_i w_{11} + e_{i1}, z_i w_{12} + e_{i2})$

Find $\mathbf{w}_1 = (w_{11}, w_{12})$ and z_i to minimize the error

$$x_i = z_i w_1 + e_i$$

 $(x_{i1}, x_{i2}) = (z_i w_{11} + e_{i1}, z_i w_{12} + e_{i2})$

Find $\mathbf{w}_1 = (w_{11}, w_{12})$ and z_i to minimize the error

error =
$$\frac{1}{N} \sum_{i=1}^{N} ((x_{i1}, x_{i2}) - z_i(w_{11}, w_{12}))'((x_{i1}, x_{i2}) - z_i(w_{11}, w_{12}))$$

$$\mathbf{x}_i = z_i \mathbf{w}_1 + \mathbf{e}_i$$

 $(x_{i1}, x_{i2}) = (z_i w_{11} + e_{i1}, z_i w_{12} + e_{i2})$

Find $\mathbf{w}_1 = (w_{11}, w_{12})$ and z_i to minimize the error

error
$$= \frac{1}{N} \sum_{i=1}^{N} ((x_{i1}, x_{i2}) - z_i(w_{11}, w_{12}))'((x_{i1}, x_{i2}) - z_i(w_{11}, w_{12}))$$
$$= \frac{1}{N} \sum_{i=1}^{N} (x_{i1} - z_i w_{11})^2 + (x_{i2} - z_i w_{12})^2$$

Three Dimensional Approximation

$$x_1 = (0.09, -1.02, -0.10)$$

 $x_2 = (0.09, 1.41, 0.67)$
 $x_3 = (-0.81, -1.46, -0.54)$
 $x_4 = (1.43, 0.26, 0.61)$
 $x_5 = (1.23, 0.87, 1.33)$

Find $\mathbf{w}_1 = (w_{11}, w_{12}, w_{13})$ and z_i to provide best one dimensional approximation.

Three-Dimensional Visualization

Three-Dimensional Visualization $\mathbf{w}_1 = (0.48, 0.75, 0.46)$

$$\mathbf{x}_i = \mathbf{z}_i \mathbf{w}_1 + \mathbf{e}_i$$

$$\mathbf{x}_i = z_i \mathbf{w}_1 + \mathbf{e}_i$$

 $(x_{i1}, x_{i2}, x_{i3}) = (z_i w_{11} + e_{i1}, z_i w_{12} + e_{i2}, z_i w_{13} + e_{i3})$

$$\mathbf{x}_i = z_i \mathbf{w}_1 + \mathbf{e}_i$$

 $(x_{i1}, x_{i2}, x_{i3}) = (z_i w_{11} + e_{i1}, z_i w_{12} + e_{i2}, z_i w_{13} + e_{i3})$

Find $\mathbf{w}_1 = (w_{11}, w_{12}, w_{13})$ and z_i to minimize the error

$$\mathbf{x}_i = z_i \mathbf{w}_1 + \mathbf{e}_i$$

 $(x_{i1}, x_{i2}, x_{i3}) = (z_i w_{11} + e_{i1}, z_i w_{12} + e_{i2}, z_i w_{13} + e_{i3})$

Find $\mathbf{w}_1 = (w_{11}, w_{12}, w_{13})$ and z_i to minimize the error

error =
$$\frac{1}{N} \sum_{i=1}^{N} ((x_{i1}, x_{i2}, x_{13}) - z_i(w_{11}, w_{12}, w_{13}))'$$
$$((x_{i1}, x_{i2}, x_{i3}) - z_i(w_{11}, w_{12}, w_{13}))$$

$$\mathbf{x}_i = z_i \mathbf{w}_1 + \mathbf{e}_i$$

 $(x_{i1}, x_{i2}, x_{i3}) = (z_i w_{11} + e_{i1}, z_i w_{12} + e_{i2}, z_i w_{13} + e_{i3})$

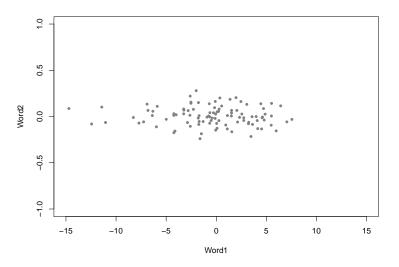
Find $\mathbf{w}_1 = (w_{11}, w_{12}, w_{13})$ and z_i to minimize the error

error
$$= \frac{1}{N} \sum_{i=1}^{N} ((x_{i1}, x_{i2}, x_{13}) - z_i(w_{11}, w_{12}, w_{13}))'$$

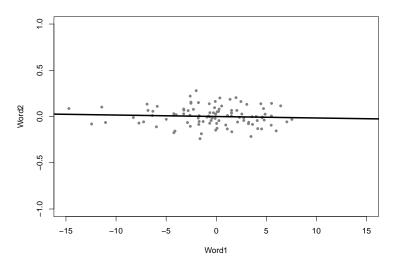
$$((x_{i1}, x_{i2}, x_{i3}) - z_i(w_{11}, w_{12}, w_{13}))$$

$$= \frac{1}{N} \sum_{i=1}^{N} (x_{i1} - z_i w_{11})^2 + (x_{i2} - z_i w_{12})^2 + (x_{i3} - z_i w_{13})^2$$

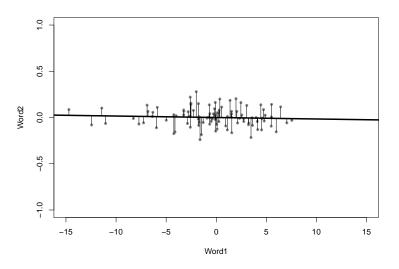
Principal Component Analysis



Principal Component Analysis



Principal Component Analysis



$$\boldsymbol{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$$

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$$

Principal Component Output:

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$$

Principal Component Output:

1) K Principal Components \mathbf{w}_k

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$$

Principal Component Output:

1) K Principal Components \mathbf{w}_k

$$\mathbf{w}_k = (w_{1k}, w_{2k}, \dots, w_{Jk})$$

$$\boldsymbol{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$$

Principal Component Output:

1) K Principal Components \mathbf{w}_k

$$\mathbf{w}_k = (w_{1k}, w_{2k}, \dots, w_{Jk})$$

2) K component vector describing loadings on principal components for each document

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$$

Principal Component Output:

1) K Principal Components \mathbf{w}_k

$$\mathbf{w}_k = (w_{1k}, w_{2k}, \dots, w_{Jk})$$

 K component vector describing loadings on principal components for each document

$$\mathbf{z}_i = (z_{1i}, z_{2i}, \ldots, z_{Ki})$$



An Introduction to Eigenvectors, Values, and Diagonalization

Definition

Suppose **A** is an $N \times N$ matrix and λ is a scalar. If

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$$

Then \mathbf{x} is an eigenvector and λ is the associated eigenvalue

An Introduction to Eigenvectors, Values, and Diagonalization

Definition

Suppose **A** is an $N \times N$ matrix and λ is a scalar. If

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$$

Then \mathbf{x} is an eigenvector and λ is the associated eigenvalue

- **A** stretches the eigenvector **x**

Definition

Suppose **A** is an $N \times N$ matrix and λ is a scalar. If

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$$

Then ${\bf x}$ is an eigenvector and λ is the associated eigenvalue

- **A** stretches the eigenvector **x**
- ${\bf A}$ stretches ${\bf x}$ by λ

Definition

Suppose **A** is an $N \times N$ matrix and λ is a scalar. If

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$$

Then \mathbf{x} is an eigenvector and λ is the associated eigenvalue

- **A** stretches the eigenvector **x**
- ${\bf \it A}$ stretches ${\bf \it x}$ by λ
- To find eigenvectors/values: (eigen in R)

Definition

Suppose **A** is an $N \times N$ matrix and λ is a scalar. If

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$$

Then x is an eigenvector and λ is the associated eigenvalue

- **A** stretches the eigenvector **x**
- ${\bf A}$ stretches ${\bf x}$ by λ
- To find eigenvectors/values: (eigen in R)
 - Find λ that solves $\det(\mathbf{A} \lambda \mathbf{I}) = 0$

Definition

Suppose **A** is an $N \times N$ matrix and λ is a scalar. If

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$$

Then ${m x}$ is an eigenvector and λ is the associated eigenvalue

- **A** stretches the eigenvector **x**
- ${\bf A}$ stretches ${\bf x}$ by λ
- To find eigenvectors/values: (eigen in R)
 - Find λ that solves $\det(\mathbf{A} \lambda \mathbf{I}) = 0$
 - Find vectors in null space of:

Definition

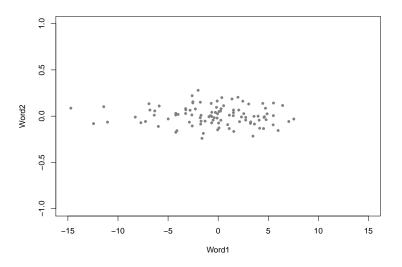
Suppose **A** is an $N \times N$ matrix and λ is a scalar. If

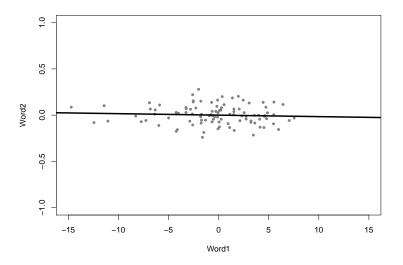
$$\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$$

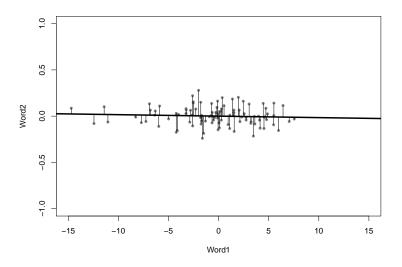
Then x is an eigenvector and λ is the associated eigenvalue

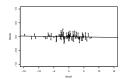
- \boldsymbol{A} stretches the eigenvector \boldsymbol{x}
- ${\bf A}$ stretches ${\bf x}$ by λ
- To find eigenvectors/values: (eigen in R)
 - Find λ that solves $\det(\mathbf{A} \lambda \mathbf{I}) = 0$
 - Find vectors in null space of:

$$(\mathbf{A} - \lambda \mathbf{I}) = 0$$

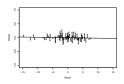






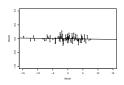


Original data:



Original data:

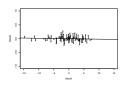
$$\mathbf{x}_i = (x_{i1}, x_{i2})$$



Original data:

$$\mathbf{x}_i = (x_{i1}, x_{i2})$$

Which we approximate with



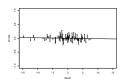
Original data:

$$\mathbf{x}_i = (x_{i1}, x_{i2})$$

Which we approximate with

$$\tilde{\boldsymbol{x}}_{i} = z_{i} \boldsymbol{w}_{1} \\
= z_{i} (w_{11}, w_{12})$$





Original data $\mathbf{x}_i \in \Re^J$

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$$

Which we approximate with $L \leq J$ weights z_{il} and vectors $\boldsymbol{w}_l \in \Re^J$

$$\tilde{\mathbf{x}}_i = \mathbf{z}_{i1}\mathbf{w}_1 + \mathbf{z}_{i2}\mathbf{w}_2 + \ldots + \mathbf{z}_{iL}\mathbf{w}_L$$

Define
$$\theta = (\underbrace{Z}_{N \times L}, \underbrace{W_L}_{I \times I})$$



Consider 1-dimensional case (L=1), centered data, and $||\mathbf{w}_1||=1$.

Consider 1-dimensional case (L=1), centered data, and $||{m w}_1||=1$.

$$f(\theta, \mathbf{X}) = \frac{1}{N} \sum_{i=1}^{N} ||\mathbf{x}_i - z_{i1} \mathbf{w}_1||^2$$

Consider 1-dimensional case (L=1), centered data, and $||\mathbf{w}_1||=1$.

$$f(\theta, \mathbf{X}) = \frac{1}{N} \sum_{i=1}^{N} ||\mathbf{x}_i - z_{i1} \mathbf{w}_1||^2$$
$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_i - z_{i1} \mathbf{w}_1)' (\mathbf{x}_i - z_{i1} \mathbf{w}_1)$$

Consider 1-dimensional case (L=1), centered data, and $||\mathbf{w}_1|| = 1$.

$$f(\theta, \mathbf{X}) = \frac{1}{N} \sum_{i=1}^{N} ||\mathbf{x}_{i} - z_{i1} \mathbf{w}_{1}||^{2}$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - z_{i1} \mathbf{w}_{1})' (\mathbf{x}_{i} - z_{i1} \mathbf{w}_{1})$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i}' \mathbf{x}_{i} - 2z_{i1} \mathbf{w}_{1}' \mathbf{x}_{i} + z_{i1}^{2})$$

Consider 1-dimensional case (L=1), centered data, and $||{m w}_1||=1$.

$$f(\theta, \mathbf{X}) = \frac{1}{N} \sum_{i=1}^{N} ||\mathbf{x}_{i} - z_{i1} \mathbf{w}_{1}||^{2}$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - z_{i1} \mathbf{w}_{1})' (\mathbf{x}_{i} - z_{i1} \mathbf{w}_{1})$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i}' \mathbf{x}_{i} - 2z_{i1} \mathbf{w}_{1}' \mathbf{x}_{i} + z_{i1}^{2})$$

$$\mathbf{w}_{1}^{'}\mathbf{w}_{1}=1$$

$$\frac{\partial f(\boldsymbol{\theta}, \boldsymbol{X})}{\partial z_{i1}} = -\frac{2\boldsymbol{w}_1'\boldsymbol{x}_i + 2z_{i1}}{N}$$

$$\frac{\partial f(\boldsymbol{\theta}, \boldsymbol{X})}{\partial z_{i1}} = -\frac{2\boldsymbol{w}_{1}'\boldsymbol{x}_{i} + 2z_{i1}}{N}$$
$$0 = -\frac{2\boldsymbol{w}_{1}'\boldsymbol{x}_{i} + 2z_{i1}^{*}}{N}$$

$$\frac{\partial f(\boldsymbol{\theta}, \boldsymbol{X})}{\partial z_{i1}} = -\frac{2\boldsymbol{w}_{1}'\boldsymbol{x}_{i} + 2z_{i1}}{N}$$

$$0 = -\frac{2\boldsymbol{w}_{1}'\boldsymbol{x}_{i} + 2z_{i1}^{*}}{N}$$

$$z_{i1}^{*} = \boldsymbol{w}_{1}'\boldsymbol{x}_{i}$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - z_{i1}^{*} \mathbf{w}_{1})' (\mathbf{x}_{i} - z_{i1}^{*} \mathbf{w}_{1})$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - z_{i1}^{*} \mathbf{w}_{1})' (\mathbf{x}_{i} - z_{i1}^{*} \mathbf{w}_{1})$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\underbrace{\mathbf{x}_{i}' \mathbf{x}_{i}}_{\text{Constant}} -2z_{i1}^{*} \underbrace{\mathbf{w}_{1}' \mathbf{x}_{i}}_{z_{i1}^{*}} + (z_{i1}^{*})^{2} \underbrace{\mathbf{w}_{1}' \mathbf{w}_{1}}_{1})$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \mathbf{z}_{i1}^{*} \mathbf{w}_{1})' (\mathbf{x}_{i} - \mathbf{z}_{i1}^{*} \mathbf{w}_{1})$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\underbrace{\mathbf{x}_{i}' \mathbf{x}_{i}}_{\text{Constant}} - 2\mathbf{z}_{i1}^{*} \underbrace{\mathbf{w}_{1}' \mathbf{x}_{i}}_{\mathbf{z}_{i1}^{*}} + (\mathbf{z}_{i1}^{*})^{2} \underbrace{\mathbf{w}_{1}' \mathbf{w}_{1}}_{1})$$

$$= -\frac{1}{N} \sum_{i=1}^{N} (\mathbf{z}_{i1}^{*})^{2} + c$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \mathbf{z}_{i1}^{*} \mathbf{w}_{1})' (\mathbf{x}_{i} - \mathbf{z}_{i1}^{*} \mathbf{w}_{1})$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\underbrace{\mathbf{x}_{i}' \mathbf{x}_{i}}_{\text{Constant}} - 2\mathbf{z}_{i1}^{*} \underbrace{\mathbf{w}_{1}' \mathbf{x}_{i}}_{\mathbf{z}_{i1}^{*}} + (\mathbf{z}_{i1}^{*})^{2} \underbrace{\mathbf{w}_{1}' \mathbf{w}_{1}}_{1})$$

$$= -\frac{1}{N} \sum_{i=1}^{N} (\mathbf{z}_{i1}^{*})^{2} + c$$

$$= -\frac{1}{N} \sum_{i=1}^{N} \mathbf{w}_{1}' \mathbf{x}_{i} \mathbf{x}_{i}' \mathbf{w}_{1}$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \mathbf{z}_{i1}^{*} \mathbf{w}_{1})' (\mathbf{x}_{i} - \mathbf{z}_{i1}^{*} \mathbf{w}_{1})$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\underbrace{\mathbf{x}_{i}' \mathbf{x}_{i}}_{\text{Constant}} - 2\mathbf{z}_{i1}^{*} \underbrace{\mathbf{w}_{1}' \mathbf{x}_{i}}_{\mathbf{z}_{i1}^{*}} + (\mathbf{z}_{i1}^{*})^{2} \underbrace{\mathbf{w}_{1}' \mathbf{w}_{1}}_{1})$$

$$= -\frac{1}{N} \sum_{i=1}^{N} (\mathbf{z}_{i1}^{*})^{2} + c$$

$$= -\frac{1}{N} \sum_{i=1}^{N} \mathbf{w}_{1}' \mathbf{x}_{i} \mathbf{x}_{i}' \mathbf{w}_{1}$$

$$= -\mathbf{w}_{1}' \mathbf{\Sigma} \mathbf{w}_{1}$$

$$= -\mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1$$

$$= -\mathbf{w}_1^{\prime} \mathbf{\Sigma} \mathbf{w}_1$$

$$= -\mathbf{w}_1^{\prime} \mathbf{\Sigma} \mathbf{w}_1$$

where Σ is the :

- Empirical covariance matrix $\leftrightarrow \frac{1}{N} \boldsymbol{X}' \boldsymbol{X}$

$$= -\mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1$$

- Empirical covariance matrix $\rightsquigarrow \frac{1}{N} \boldsymbol{X}' \boldsymbol{X}$
- Variance of the projected data. Define

$$= -\mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1$$

- Empirical covariance matrix $\rightsquigarrow \frac{1}{N} \mathbf{X}' \mathbf{X}$
- Variance of the projected data. Define

$$\mathbf{z}_1 = (\mathbf{w}_1 \mathbf{x}_1, \mathbf{w}_1 \mathbf{x}_2, \dots, \mathbf{w}_1 \mathbf{x}_N)$$

$$= -\mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1$$

- Empirical covariance matrix $\rightarrow \frac{1}{N} \boldsymbol{X}' \boldsymbol{X}$
- Variance of the projected data. Define

$$z_1 = (w_1x_1, w_1x_2, ..., w_1x_N)$$

 $var(z_1) = E[z_1^2] - E[z_1]^2$

$$= -\mathbf{w}_1^{'}\mathbf{\Sigma}\mathbf{w}_1$$

- Empirical covariance matrix $\rightsquigarrow \frac{1}{N} \boldsymbol{X}' \boldsymbol{X}$
- Variance of the projected data. Define

$$z_1 = (w_1x_1, w_1x_2, ..., w_1x_N)$$

 $var(z_1) = E[z_1^2] - E[z_1]^2$
 $= \frac{1}{N} \sum_{i=1}^{N} z_{i1}^2 - 0$

$$= -\mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1$$

- Empirical covariance matrix $\rightsquigarrow \frac{1}{N} \boldsymbol{X}' \boldsymbol{X}$
- Variance of the projected data. Define

$$egin{array}{lcl} m{z}_1 &=& (m{w}_1m{x}_1, m{w}_1m{x}_2, \dots, m{w}_1m{x}_N) \\ ext{var}(m{z}_1) &=& E[m{z}_1^2] - E[m{z}_1]^2 \\ &=& rac{1}{N} \sum_{i=1}^N z_{i1}^2 - 0 \\ &=& rac{1}{N} \sum_{i=1}^N m{w}_1' m{x}_i m{x}_i' m{w}_1 = m{w}_1' m{\Sigma} m{w}_1 \end{array}$$

$$= -\mathbf{w}_1^{'}\mathbf{\Sigma}\mathbf{w}_1$$

where Σ is the :

- Empirical covariance matrix $\rightsquigarrow \frac{1}{N} \boldsymbol{X}' \boldsymbol{X}$
- Variance of the projected data. Define

$$z_{1} = (w_{1}x_{1}, w_{1}x_{2}, ..., w_{1}x_{N})$$

$$var(z_{1}) = E[z_{1}^{2}] - E[z_{1}]^{2}$$

$$= \frac{1}{N} \sum_{i=1}^{N} z_{i1}^{2} - 0$$

$$= \frac{1}{N} \sum_{i=1}^{N} w'_{1}x_{i}x'_{i}w_{1} = w'_{1}\Sigma w_{1}$$

Minimize reconstruction error

4ロト 4回ト 4 き ト 4 き ト き め 9 0 0

$$= -\mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1$$

where Σ is the :

- Empirical covariance matrix $\rightsquigarrow \frac{1}{N} \boldsymbol{X}' \boldsymbol{X}$
- Variance of the projected data. Define

$$egin{array}{lcl} m{z}_1 &=& (m{w}_1m{x}_1, m{w}_1m{x}_2, \dots, m{w}_1m{x}_N) \\ ext{var}(m{z}_1) &=& E[m{z}_1^2] - E[m{z}_1]^2 \\ &=& rac{1}{N} \sum_{i=1}^N z_{i1}^2 - 0 \\ &=& rac{1}{N} \sum_{i=1}^N m{w}_1' m{x}_i m{x}_i' m{w}_1 = m{w}_1' m{\Sigma} m{w}_1 \end{array}$$

Minimize reconstruction error \rightsquigarrow maximize variance of projected data

$$g(z^*, w_1, X) = w_1' \Sigma w_1 - \lambda_1 (w_1' w_1 - 1)$$

$$g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) = \mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}_1' \mathbf{w}_1 - 1)$$

$$\frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} = 2\mathbf{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1$$

$$g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) = \mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}_1' \mathbf{w}_1 - 1)$$

$$\frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} = 2\mathbf{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1$$

$$\mathbf{\Sigma} \mathbf{w}_1^* = \lambda_1 \mathbf{w}_1^*$$

Maximize variance, subject to constraints

$$g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) = \mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}_1' \mathbf{w}_1 - 1)$$

$$\frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} = 2\mathbf{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1$$

$$\mathbf{\Sigma} \mathbf{w}_1^* = \lambda_1 \mathbf{w}_1^*$$

 $\mathbf{w}_1^* = \text{Eigenvector of } \mathbf{\Sigma}$

Maximize variance, subject to constraints

$$g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) = \mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}_1' \mathbf{w}_1 - 1)$$

$$\frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} = 2\mathbf{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1$$

$$\mathbf{\Sigma} \mathbf{w}_1^* = \lambda_1 \mathbf{w}_1^*$$

 $\mathbf{w}_1^* = \mathsf{Eigenvector} \ \mathsf{of} \ \mathsf{\Sigma} \ (!!!!!!)$

Maximize variance, subject to constraints

$$g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) = \mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}_1' \mathbf{w}_1 - 1)$$

$$\frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} = 2\mathbf{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1$$

$$\mathbf{\Sigma} \mathbf{w}_1^* = \lambda_1 \mathbf{w}_1^*$$

 $\mathbf{w}_1^* = \mathsf{Eigenvector} \ \mathsf{\Sigma} \ (!!!!!!)$

We want \mathbf{w}_1 to maximize variance and

Maximize variance, subject to constraints

$$g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) = \mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}_1' \mathbf{w}_1 - 1)$$

$$\frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} = 2\mathbf{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1$$

$$\mathbf{\Sigma} \mathbf{w}_1^* = \lambda_1 \mathbf{w}_1^*$$

 $\mathbf{w}_1^* = \text{Eigenvector of } \mathbf{\Sigma} (!!!!!!)$

We want \mathbf{w}_1 to maximize variance and

$${\pmb w}_1^{'} {\pmb \Sigma} {\pmb w}_1 = \lambda_1$$

Maximize variance, subject to constraints

$$g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X}) = \mathbf{w}_1' \mathbf{\Sigma} \mathbf{w}_1 - \lambda_1 (\mathbf{w}_1' \mathbf{w}_1 - 1)$$

$$\frac{\partial g(\mathbf{z}^*, \mathbf{w}_1, \mathbf{X})}{\partial \mathbf{w}_1} = 2\mathbf{\Sigma} \mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1$$

$$\mathbf{\Sigma} \mathbf{w}_1^* = \lambda_1 \mathbf{w}_1^*$$

 $\mathbf{w}_1^* = \text{Eigenvector of } \mathbf{\Sigma} (!!!!!!)$

We want \mathbf{w}_1 to maximize variance and

$$oldsymbol{w}_1^{'} oldsymbol{\Sigma} oldsymbol{w}_1 = \lambda_1$$

So ${m w}_1$ is eigenvector associated with the largest eigenvalue λ_1

An Introduction to Eigenvectors, Values, and Diagonalization

Theorem

Suppose **A** is an invertible $N \times N$ matrix with N linearly independent eigenvectors. Then we can write **A** as,

$$\mathbf{A} = \mathbf{W}' \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_N \end{pmatrix} \mathbf{W}$$

where $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N)$ is an $N \times N$ matrix with the N eigenvectors as column vectors.

An Introduction to Eigenvectors, Values, and Diagonalization

Definition

Suppose A is a covariance matrix. Then, we can write A as

$$\mathbf{A} = \mathbf{W}' \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_N \end{pmatrix} \mathbf{W}$$

Where $\lambda_1 > \lambda_2 > \ldots > \lambda_N \geq 0$.

We will call \mathbf{w}_1 the first eigenvector, \mathbf{w}_2 the second eigenvector, ..., \mathbf{w}_j the i^{th} eigenvector.

Theorem

Suppose we want to approximate N observations $\mathbf{x}_i \in \mathbb{R}^J$ with L < J orthogonal-unit length vectors $\mathbf{w}_I \in \mathbb{R}^J$ with associated scores z_{il} to minimize reconstruction error:

Theorem

Suppose we want to approximate N observations $\mathbf{x}_i \in \mathbb{R}^J$ with L < J orthogonal-unit length vectors $\mathbf{w}_I \in \mathbb{R}^J$ with associated scores z_{il} to minimize reconstruction error:

$$f(\mathbf{X}, \theta) = \frac{1}{N} \sum_{i=1}^{N} ||\mathbf{x}_i - \sum_{l=1}^{L} z_{il} \mathbf{w}_l||^2$$

Theorem

Suppose we want to approximate N observations $\mathbf{x}_i \in \mathbb{R}^J$ with L < J orthogonal-unit length vectors $\mathbf{w}_I \in \mathbb{R}^J$ with associated scores z_{il} to minimize reconstruction error:

$$f(\mathbf{X}, \theta) = \frac{1}{N} \sum_{i=1}^{N} ||\mathbf{x}_i - \sum_{l=1}^{L} z_{il} \mathbf{w}_l||^2$$

The optimal solution sets each \mathbf{w}_l to be the l^{th} eigenvector of the empirical covariance matrix.

Theorem

Suppose we want to approximate N observations $\mathbf{x}_i \in \mathbb{R}^J$ with L < J orthogonal-unit length vectors $\mathbf{w}_I \in \mathbb{R}^J$ with associated scores z_{il} to minimize reconstruction error:

$$f(\mathbf{X}, \theta) = \frac{1}{N} \sum_{i=1}^{N} ||\mathbf{x}_i - \sum_{l=1}^{L} z_{il} \mathbf{w}_l||^2$$

The optimal solution sets each \mathbf{w}_{l} to be the l^{th} eigenvector of the empirical covariance matrix. Further $z_{il}^{*} = \mathbf{w}_{l}^{'} \mathbf{x}_{i}$ so that the L dimensional representation is:

Theorem

Suppose we want to approximate N observations $\mathbf{x}_i \in \mathbb{R}^J$ with L < J orthogonal-unit length vectors $\mathbf{w}_I \in \mathbb{R}^J$ with associated scores z_{il} to minimize reconstruction error:

$$f(\mathbf{X}, \theta) = \frac{1}{N} \sum_{i=1}^{N} ||\mathbf{x}_i - \sum_{l=1}^{L} z_{il} \mathbf{w}_l||^2$$

The optimal solution sets each \mathbf{w}_{l} to be the l^{th} eigenvector of the empirical covariance matrix. Further $z_{il}^{*} = \mathbf{w}_{l}^{'} \mathbf{x}_{i}$ so that the L dimensional representation is:

$$\mathbf{x}_{i}^{L} = (\mathbf{w}_{1}^{'}\mathbf{x}_{i}, \mathbf{w}_{2}^{'}\mathbf{x}_{i}, \ldots, \mathbf{w}_{L}^{'}\mathbf{x}_{i})$$

◆ロ > ◆昼 > ◆差 > ◆差 > 差 り < ②</p>

Consider press releases from 2005 US Senators

Consider press releases from 2005 US Senators Define $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ as the rate senator i uses J words.

Consider press releases from 2005 US Senators Define $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ as the rate senator i uses J words.

$$x_{ij} = \frac{\text{No. Times } i \text{ uses word } j}{\text{No. words } i \text{ uses}}$$

Consider press releases from 2005 US Senators Define $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ as the rate senator i uses J words.

$$x_{ij} = \frac{\text{No. Times } i \text{ uses word } j}{\text{No. words } i \text{ uses}}$$

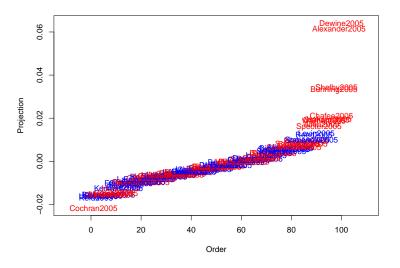
dtm: 100×2796 matrix containing word rates for senators

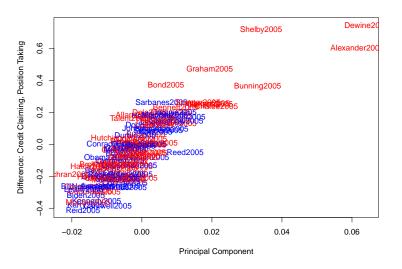
Consider press releases from 2005 US Senators Define $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ as the rate senator i uses J words.

$$x_{ij} = \frac{\text{No. Times } i \text{ uses word } j}{\text{No. words } i \text{ uses}}$$

dtm: 100×2796 matrix containing word rates for senators prcomp(dtm) applies principal components

```
load("SenateTDM.RData")
dtm<- t(tdm)
for(z in 1:100){
dtm[z,]<- dtm[z,]/sum(dtm[z,])
}
store<- prcomp(dtm, scale = F)
scores<- store$x[,1]</pre>
```





Probabilistic Principal Components (Tipping and Bishop 1999)

$$m{x} | m{w} \sim ext{Multivariate Normal}(m{Z}m{W} + m{\mu}, \sigma^2 m{I})$$
 $m{w} \sim ext{Multivariate Normal}(m{0}, m{I})$
 $m{x} \sim ext{Multivariate Normal}(m{\mu}, m{\Sigma})$
 $m{\Sigma} = m{W}m{W}' + \sigma^2 m{I}$

- Log-likelihood → straightforward
- 2) Optimization via EM-Algorithm
- 3) Corresponds to traditional PCA is $\lim_{\sigma^2} \to 0$
- 4) Closely related to Factor analysis.

How do we select the number of dimensions $L? \rightsquigarrow \mathsf{Model}$ We want to minimize reconstruction error

We want to minimize reconstruction error → how well did we do?

We want to minimize reconstruction error → how well did we do?

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \sum_{l=1}^{L} z_{il} w_l||^2$$

We want to minimize reconstruction error → how well did we do?

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \sum_{l=1}^{L} z_{il} w_l||^2$$

Simplifying:

We want to minimize reconstruction error → how well did we do?

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \sum_{l=1}^{L} z_{il} \mathbf{w}_l||^2$$

Simplifying:

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})' (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})$$

We want to minimize reconstruction error → how well did we do?

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \sum_{l=1}^{L} z_{il} w_l||^2$$

Simplifying:

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})' (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})$$

Four types of terms:

We want to minimize reconstruction error → how well did we do?

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \sum_{l=1}^{L} z_{il} \mathbf{w}_l||^2$$

Simplifying:

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})' (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})$$

Four types of terms: 1) $\mathbf{x}_{i}^{'}\mathbf{x}_{i}$

1)
$$\mathbf{x}_{i}^{'}\mathbf{x}_{i}$$

We want to minimize reconstruction error → how well did we do?

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \sum_{l=1}^{L} z_{il} \mathbf{w}_l||^2$$

Simplifying:

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})' (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})$$

Four types of terms:

- 1) $\mathbf{x}_{i}^{'}\mathbf{x}_{i}$
- 2) $z_{ij}z_{ik}\mathbf{w}_{i}^{\prime}\mathbf{w}_{k}=z_{ij}z_{ik}0=0$

We want to minimize reconstruction error → how well did we do?

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \sum_{l=1}^{L} z_{il} \mathbf{w}_l||^2$$

Simplifying:

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})' (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})$$

Four types of terms:

1)
$$\mathbf{x}_{i}^{'}\mathbf{x}_{i}$$

2)
$$z_{ij}z_{ik}\mathbf{w}_{j}^{\prime}\mathbf{w}_{k} = z_{ij}z_{ik}0 = 0$$

3) $z_{ij}z_{ij}\mathbf{w}_{i}^{\prime}\mathbf{w}_{j} = z_{ii}^{2}$

3)
$$z_{ij}z_{ij}\mathbf{w}_{j}^{\prime}\mathbf{w}_{j}=z_{ij}^{2}$$

We want to minimize reconstruction error \rightsquigarrow how well did we do?

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \sum_{l=1}^{L} z_{il} \mathbf{w}_l||^2$$

Simplifying:

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})' (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})$$

Four types of terms:

1)
$$\mathbf{x}_{i}^{\prime}\mathbf{x}_{i}$$

2)
$$z_{ij}z_{ik}\mathbf{w}'_{j}\mathbf{w}_{k} = z_{ij}z_{ik}0 = 0$$

3) $z_{ij}z_{ij}\mathbf{w}'_{i}\mathbf{w}_{j} = z_{ij}^{2}$

3)
$$z_{ij}z_{ij}\boldsymbol{w}_{j}^{'}\boldsymbol{w}_{j}=z_{ij}^{2}$$

4)
$$\mathbf{x}'_{i} \sum_{l=1}^{L} z_{il} \mathbf{w}_{l} = \sum_{l=1}^{L} z_{il}^{2}$$

We want to minimize reconstruction error \rightsquigarrow how well did we do?

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} ||x_i - \sum_{l=1}^{L} z_{il} w_l||^2$$

Simplifying:

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})' (\mathbf{x}_{i} - \sum_{l=1}^{L} z_{il} \mathbf{w}_{l})$$

= $\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i}' \mathbf{x}_{i} - \sum_{l=1}^{L} z_{il}^{2})$

Four types of terms:

1)
$$\mathbf{x}_{i}^{\prime}\mathbf{x}_{i}$$

2)
$$z_{ij}z_{ik}\mathbf{w}_{j}'\mathbf{w}_{k} = z_{ij}z_{ik}0 = 0$$

3) $z_{ij}z_{ij}\mathbf{w}_{j}'\mathbf{w}_{j} = z_{ij}^{2}$

3)
$$z_{ij}z_{ij}\boldsymbol{w}_{j}^{'}\boldsymbol{w}_{j}=z_{ij}^{2}$$

4)
$$\mathbf{x}'_{i} \sum_{l=1}^{L} z_{il} \mathbf{w}_{l} = \sum_{l=1}^{L} z_{il}^{2}$$

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} z_{il}^{2} \right)$$

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} z_{il}^{2} \right)$$
$$= \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} \mathbf{w}_{l} \mathbf{x}_{i} \mathbf{x}_{i}^{'} \mathbf{w}_{l} \right)$$

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} z_{il}^{2} \right)$$

= $\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} \mathbf{w}_{l} \mathbf{x}_{i} \mathbf{x}_{i}^{'} \mathbf{w}_{l} \right)$
= $\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \frac{1}{N} \sum_{l=1}^{L} \sum_{l=1}^{N} \mathbf{w}_{i}^{'} \mathbf{x}_{i} \mathbf{x}_{i}^{'} \mathbf{w}_{i}$

$$error(L) = \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} z_{il}^{2} \right)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} \mathbf{w}_{l} \mathbf{x}_{i} \mathbf{x}_{i}^{'} \mathbf{w}_{l} \right)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \frac{1}{N} \sum_{l=1}^{L} \sum_{l=1}^{N} \mathbf{w}_{i}^{'} \mathbf{x}_{i} \mathbf{x}_{i}^{'} \mathbf{w}_{i}$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \sum_{l=1}^{L} \mathbf{w}_{l}^{'} \mathbf{\Sigma} \mathbf{w}_{l}$$

error(L) =
$$\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} z_{il}^{2} \right)$$

= $\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} \mathbf{w}_{l} \mathbf{x}_{i} \mathbf{x}_{i}^{'} \mathbf{w}_{l} \right)$
= $\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \frac{1}{N} \sum_{l=1}^{L} \sum_{l=1}^{N} \mathbf{w}_{i}^{'} \mathbf{x}_{i} \mathbf{x}_{i}^{'} \mathbf{w}_{i}$
= $\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \sum_{l=1}^{L} \mathbf{w}_{l}^{'} \mathbf{\Sigma} \mathbf{w}_{l}$
= $\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \sum_{l=1}^{L} \lambda_{l} \mathbf{w}_{l}^{'} \mathbf{w}_{l}$

$$error(L) = \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} \mathbf{z}_{il}^{2} \right)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} - \sum_{l=1}^{L} \mathbf{w}_{l} \mathbf{x}_{i} \mathbf{x}_{i}^{'} \mathbf{w}_{l} \right)$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \frac{1}{N} \sum_{l=1}^{L} \sum_{l=1}^{N} \mathbf{w}_{i}^{'} \mathbf{x}_{i} \mathbf{x}_{i}^{'} \mathbf{w}_{i}$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \sum_{l=1}^{L} \mathbf{w}_{l}^{'} \mathbf{\Sigma} \mathbf{w}_{l}$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \sum_{l=1}^{L} \lambda_{l} \mathbf{w}_{l}^{'} \mathbf{w}_{l}$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}^{'} \mathbf{x}_{i} \right) - \sum_{l=1}^{L} \lambda_{l}$$

error(J) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i}^{'} \mathbf{x}_{i}) - \sum_{l=1}^{J} \lambda_{l} = 0$$

error(J) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i}' \mathbf{x}_{i}) - \sum_{l=1}^{J} \lambda_{l} = 0$$

So for L < J,

error(J) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i}^{'} \mathbf{x}_{i}) - \sum_{l=1}^{J} \lambda_{l} = 0$$

So for L < J.

$$0 = \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}' \mathbf{x}_{i} \right) - \left(\sum_{l=1}^{L} \lambda_{l} + \sum_{j=L+1}^{J} \lambda_{l} \right)$$

error(J) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i}^{'} \mathbf{x}_{i}) - \sum_{l=1}^{J} \lambda_{l} = 0$$

So for L < J,

$$0 = \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}' \mathbf{x}_{i} \right) - \left(\sum_{l=1}^{L} \lambda_{l} + \sum_{j=L+1}^{J} \lambda_{l} \right)$$
$$\sum_{j=L+1}^{J} \lambda_{l} = \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}' \mathbf{x}_{i} \right) - \sum_{l=1}^{L} \lambda_{l}$$

error(J) =
$$\frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_{i}' \mathbf{x}_{i}) - \sum_{l=1}^{J} \lambda_{l} = 0$$

So for L < J.

$$0 = \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}' \mathbf{x}_{i} \right) - \left(\sum_{l=1}^{L} \lambda_{l} + \sum_{j=L+1}^{J} \lambda_{l} \right)$$

$$\sum_{j=L+1}^{J} \lambda_{l} = \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{x}_{i}' \mathbf{x}_{i} \right) - \sum_{l=1}^{L} \lambda_{l}$$

$$\sum_{j=L+1}^{J} \lambda_{l} = \text{error}(L)$$

$$\sum_{j=L+1}^{J} \lambda_{I} = \text{error}(L)$$

$$\sum_{j=L+1}^{J} \lambda_{I} = \text{error}(L)$$

- Error = Sum of "remaining" eigenvalues

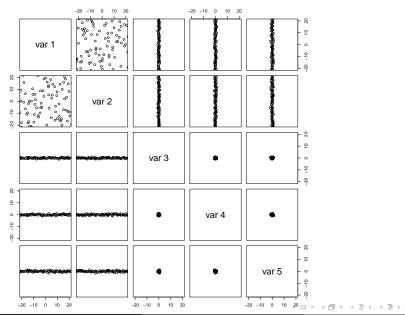
$$\sum_{j=L+1}^J \lambda_I = \operatorname{error}(L)$$

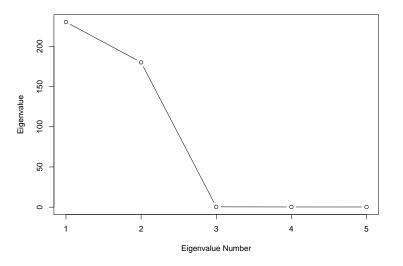
- Error = Sum of "remaining" eigenvalues
- Total variance explained = (sum of included eigenvalues)/(sum of all eigenvalues)

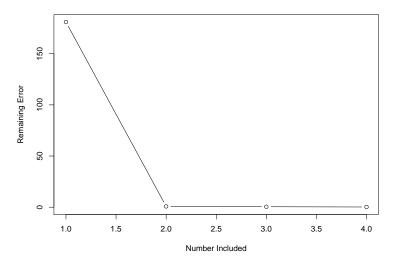
$$\sum_{j=L+1}^{J} \lambda_{I} = \text{error}(L)$$

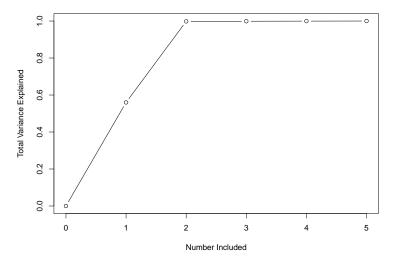
- Error = Sum of "remaining" eigenvalues
- Total variance explained = (sum of included eigenvalues)/(sum of all eigenvalues)

Recommendation >>> look for Elbow









What is the true underlying dimensionality of X?

What is the true underlying dimensionality of X? J

What is the true underlying dimensionality of X? J(!!!!!)

What is the true underlying dimensionality of X? J(!!!!!)

 Attempts to assess dimensionality require a model → some way to tradeoff accuracy of reconstruction with simplicity

What is the true underlying dimensionality of X? J(!!!!!)

- Attempts to assess dimensionality require a model → some way to tradeoff accuracy of reconstruction with simplicity
- Any answer (no matter how creatively obtained) supposes you have the right function to measure tradeoff

What is the true underlying dimensionality of X? J(!!!!!)

- Attempts to assess dimensionality require a model → some way to tradeoff accuracy of reconstruction with simplicity
- Any answer (no matter how creatively obtained) supposes you have the right function to measure tradeoff
- The "right" number of dimensions depends on the task you have in mind

What is the true underlying dimensionality of X? J(!!!!!)

- Attempts to assess dimensionality require a model → some way to tradeoff accuracy of reconstruction with simplicity
- Any answer (no matter how creatively obtained) supposes you have the right function to measure tradeoff
- The "right" number of dimensions depends on the task you have in mind

Mathematical model → insufficient to make modeling decision

Appendix

Define a Kernel $(N \times N)$ matrix as:

$$K = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_N) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \dots & k(x_N, x_N) \end{pmatrix}$$

where $k(\cdot, \cdot)$ is a function that behaves like a similarity function.

Define a Kernel $(N \times N)$ matrix as:

$$K = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_N) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \dots & k(x_N, x_N) \end{pmatrix}$$

where $k(\cdot,\cdot)$ is a function that behaves like a similarity function. Where we suppose this matrix emerges from applying $\phi: \Re^J \to \Re^M$ to the data and then taking the inner product:

Define a Kernel $(N \times N)$ matrix as:

$$K = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_N) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \dots & k(x_N, x_N) \end{pmatrix}$$

where $k(\cdot,\cdot)$ is a function that behaves like a similarity function. Where we suppose this matrix emerges from applying $\phi: \Re^J \to \Re^M$ to the data and then taking the inner product:

$$K = \Phi \Phi'$$
 (The inner product matrix)

Define a Kernel $(N \times N)$ matrix as:

$$K = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_N) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \dots & k(x_N, x_N) \end{pmatrix}$$

where $k(\cdot,\cdot)$ is a function that behaves like a similarity function. Where we suppose this matrix emerges from applying $\phi: \Re^J \to \Re^M$ to the data and then taking the inner product:

$$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}' \text{ (The inner product matrix)}$$

$$= \begin{pmatrix} \phi(\mathbf{x}_1)'\phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)'\phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_1)'\phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)'\phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)'\phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_2)'\phi(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_N)'\phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)'\phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_N)'\phi(\mathbf{x}_N) \end{pmatrix}$$

Define a Kernel $(N \times N)$ matrix as:

$$K = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_N) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \dots & k(x_N, x_N) \end{pmatrix}$$

where $k(\cdot,\cdot)$ is a function that behaves like a similarity function. Where we suppose this matrix emerges from applying $\phi: \Re^J \to \Re^M$ to the data and then taking the inner product:

$$\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}' \text{ (The inner product matrix)}$$

$$= \begin{pmatrix} \phi(\mathbf{x}_1)'\phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)'\phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_1)'\phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)'\phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)'\phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_2)'\phi(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(\mathbf{x}_N)'\phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)'\phi(\mathbf{x}_2) & \dots & \phi(\mathbf{x}_N)'\phi(\mathbf{x}_N) \end{pmatrix}$$

Compute PCA of Φ from $\Phi\Phi'$

←ロト ←団 ト ← 重 ト ・ 重 ・ り へ ○

Kernel PCA PCA of **X**

PCA of ${\pmb X}$ Eigenvectors of ${\pmb X}'{\pmb X}$ ($\frac{1}{N}$ doesn't affect eigenvectors)

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 .

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

$$(\boldsymbol{X}\boldsymbol{X}')\boldsymbol{u}_1 = \lambda_1\boldsymbol{u}_1$$

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

$$(\boldsymbol{X}\boldsymbol{X}')\boldsymbol{u}_1 = \lambda_1\boldsymbol{u}_1 (\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{X}'\boldsymbol{u}_1) = \lambda_1(\boldsymbol{X}'\boldsymbol{u}_1)$$

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

$$(\mathbf{X}\mathbf{X}')\mathbf{u}_1 = \lambda_1\mathbf{u}_1$$

 $(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) = \lambda_1(\mathbf{X}'\mathbf{u}_1)$
 $= \lambda_1\mathbf{v}_1$

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

$$(\mathbf{X}\mathbf{X}')\mathbf{u}_1 = \lambda_1\mathbf{u}_1$$

 $(\mathbf{X}'\mathbf{X})(\mathbf{X}'\mathbf{u}_1) = \lambda_1(\mathbf{X}'\mathbf{u}_1)$
 $= \lambda_1\mathbf{v}_1$

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

$$(\boldsymbol{X}\boldsymbol{X}')\boldsymbol{u}_1 = \lambda_1\boldsymbol{u}_1$$

 $(\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{X}'\boldsymbol{u}_1) = \lambda_1(\boldsymbol{X}'\boldsymbol{u}_1)$
 $= \lambda_1\boldsymbol{v}_1$

$$||\textbf{\textit{v}}_1||^2=\textbf{\textit{v}}_1^{'}\textbf{\textit{v}}_1$$

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

$$(\boldsymbol{X}\boldsymbol{X}')\boldsymbol{u}_1 = \lambda_1\boldsymbol{u}_1$$

 $(\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{X}'\boldsymbol{u}_1) = \lambda_1(\boldsymbol{X}'\boldsymbol{u}_1)$
 $= \lambda_1\boldsymbol{v}_1$

$$||\mathbf{v}_{1}||^{2} = \mathbf{v}_{1}^{'}\mathbf{v}_{1}$$

$$= \mathbf{u}_{1}^{'}\mathbf{X}\mathbf{X}^{'}\mathbf{u}_{1}$$

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

$$(\boldsymbol{X}\boldsymbol{X}')\boldsymbol{u}_1 = \lambda_1\boldsymbol{u}_1$$

 $(\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{X}'\boldsymbol{u}_1) = \lambda_1(\boldsymbol{X}'\boldsymbol{u}_1)$
 $= \lambda_1\boldsymbol{v}_1$

$$||\mathbf{v}_1||^2 = \mathbf{v}_1' \mathbf{v}_1$$

$$= \mathbf{u}_1' \mathbf{X} \mathbf{X}' \mathbf{u}_1$$

$$= \lambda_1 \mathbf{u}_1' \mathbf{u}_1 = \lambda_1$$

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

$$(\boldsymbol{X}\boldsymbol{X}')\boldsymbol{u}_1 = \lambda_1\boldsymbol{u}_1$$

 $(\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{X}'\boldsymbol{u}_1) = \lambda_1(\boldsymbol{X}'\boldsymbol{u}_1)$
 $= \lambda_1\boldsymbol{v}_1$

But \mathbf{v}_1 needs unit length, and

$$||\mathbf{v}_1||^2 = \mathbf{v}_1' \mathbf{v}_1$$

$$= \mathbf{u}_1' \mathbf{X} \mathbf{X}' \mathbf{u}_1$$

$$= \lambda_1 \mathbf{u}_1' \mathbf{u}_1 = \lambda_1$$

So first eigenvector of $\boldsymbol{X}'\boldsymbol{X}$ is

PCA of \boldsymbol{X} Eigenvectors of $\boldsymbol{X}'\boldsymbol{X}$ ($\frac{1}{N}$ doesn't affect eigenvectors) Suppose \boldsymbol{u}_1 is an eigenvector for $\boldsymbol{X}\boldsymbol{X}'$, with value λ_1 . Then

$$(\boldsymbol{X}\boldsymbol{X}')\boldsymbol{u}_1 = \lambda_1\boldsymbol{u}_1$$

 $(\boldsymbol{X}'\boldsymbol{X})(\boldsymbol{X}'\boldsymbol{u}_1) = \lambda_1(\boldsymbol{X}'\boldsymbol{u}_1)$
 $= \lambda_1\boldsymbol{v}_1$

But \mathbf{v}_1 needs unit length, and

$$||\mathbf{v}_1||^2 = \mathbf{v}_1' \mathbf{v}_1$$

$$= \mathbf{u}_1' \mathbf{X} \mathbf{X}' \mathbf{u}_1$$

$$= \lambda_1 \mathbf{u}_1' \mathbf{u}_1 = \lambda_1$$

So first eigenvector of $\mathbf{X}'\mathbf{X}$ is

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \mathbf{X}' \mathbf{u}_1$$



 $K = \Phi \Phi'$ (assume Φ is mean-centered, for now)

 $\pmb{K} = \pmb{\Phi} \pmb{\Phi}^{'}$ (assume $\pmb{\Phi}$ is mean-centered, for now) We can obtain \pmb{u}_1 and λ_1 from \pmb{K} .

 $\pmb{K} = \pmb{\Phi} \pmb{\Phi}'$ (assume $\pmb{\Phi}$ is mean-centered, for now) We can obtain \pmb{u}_1 and λ_1 from \pmb{K} . We know that

 $\pmb{K} = \pmb{\Phi} \pmb{\Phi}'$ (assume $\pmb{\Phi}$ is mean-centered, for now) We can obtain \pmb{u}_1 and λ_1 from \pmb{K} . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\mathsf{Unknown}} \mathbf{u}_1$$

 $\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$ (assume $\mathbf{\Phi}$ is mean-centered, for now) We can obtain \mathbf{u}_1 and λ_1 from \mathbf{K} . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}^{'}}_{\mathsf{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point $\phi(\mathbf{x}_i)$, then

 $\pmb{K} = \pmb{\Phi} \pmb{\Phi}'$ (assume $\pmb{\Phi}$ is mean-centered, for now) We can obtain \pmb{u}_1 and λ_1 from \pmb{K} . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\mathsf{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point $\phi(\mathbf{x}_i)$, then

$$\phi(\mathbf{x}_i)'\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}}\phi(\mathbf{x}_i)'\mathbf{\Phi}'\mathbf{u}_1$$

 $\pmb{K} = \pmb{\Phi} \pmb{\Phi}'$ (assume $\pmb{\Phi}$ is mean-centered, for now) We can obtain \pmb{u}_1 and λ_1 from \pmb{K} . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\mathsf{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point $\phi(x_i)$, then

$$\phi(\mathbf{x}_i)'\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}}\phi(\mathbf{x}_i)'\mathbf{\Phi}'\mathbf{u}_1$$

$$\phi(\mathbf{x}_i)'\mathbf{\Phi}' = \left[\phi(\mathbf{x}_i)'\phi(\mathbf{x}_1), \phi(\mathbf{x}_i)'\phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)'\phi(\mathbf{x}_N)\right]$$

 $\pmb{K} = \pmb{\Phi} \pmb{\Phi}'$ (assume $\pmb{\Phi}$ is mean-centered, for now) We can obtain \pmb{u}_1 and λ_1 from \pmb{K} . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\mathsf{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point $\phi(x_i)$, then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

$$\phi(\mathbf{x}_i)' \mathbf{\Phi}' = \left[\phi(\mathbf{x}_i)' \phi(\mathbf{x}_1), \phi(\mathbf{x}_i)' \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)' \phi(\mathbf{x}_N) \right]$$

$$= \left[k(\mathbf{x}_i, \mathbf{x}_1), k(\mathbf{x}_i, \mathbf{x}_2), \dots, k(\mathbf{x}_i, \mathbf{x}_N) \right]$$

 $\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$ (assume $\mathbf{\Phi}$ is mean-centered, for now) We can obtain \mathbf{u}_1 and λ_1 from \mathbf{K} . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\mathsf{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point $\phi(x_i)$, then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

$$\phi(\mathbf{x}_i)' \mathbf{\Phi}' = \left[\phi(\mathbf{x}_i)' \phi(\mathbf{x}_1), \phi(\mathbf{x}_i)' \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)' \phi(\mathbf{x}_N) \right]$$

$$= \left[k(\mathbf{x}_i, \mathbf{x}_1), k(\mathbf{x}_i, \mathbf{x}_2), \dots, k(\mathbf{x}_i, \mathbf{x}_N) \right] = \mathbf{k}(\mathbf{x}_i, *)$$

 $\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$ (assume $\mathbf{\Phi}$ is mean-centered, for now) We can obtain \mathbf{u}_1 and λ_1 from \mathbf{K} . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\mathsf{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point $\phi(x_i)$, then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

$$\phi(\mathbf{x}_i)' \mathbf{\Phi}' = \left[\phi(\mathbf{x}_i)' \phi(\mathbf{x}_1), \phi(\mathbf{x}_i)' \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)' \phi(\mathbf{x}_N) \right]$$

$$= \left[k(\mathbf{x}_i, \mathbf{x}_1), k(\mathbf{x}_i, \mathbf{x}_2), \dots, k(\mathbf{x}_i, \mathbf{x}_N) \right] = \mathbf{k}(\mathbf{x}_i, *)$$

Then, we can obtain projection for observation i using Kernel with

 $\mathbf{K} = \mathbf{\Phi}\mathbf{\Phi}'$ (assume $\mathbf{\Phi}$ is mean-centered, for now) We can obtain \mathbf{u}_1 and λ_1 from \mathbf{K} . We know that

$$\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \underbrace{\mathbf{\Phi}'}_{\mathsf{Unknown}} \mathbf{u}_1$$

But suppose we want to project a point $\phi(\mathbf{x}_i)$, then

$$\phi(\mathbf{x}_i)' \mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}} \phi(\mathbf{x}_i)' \mathbf{\Phi}' \mathbf{u}_1$$

$$\phi(\mathbf{x}_i)' \mathbf{\Phi}' = \left[\phi(\mathbf{x}_i)' \phi(\mathbf{x}_1), \phi(\mathbf{x}_i)' \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_i)' \phi(\mathbf{x}_N) \right]$$

$$= \left[k(\mathbf{x}_i, \mathbf{x}_1), k(\mathbf{x}_i, \mathbf{x}_2), \dots, k(\mathbf{x}_i, \mathbf{x}_N) \right] = \mathbf{k}(\mathbf{x}_i, *)$$

Then, we can obtain projection for observation i using Kernel with

$$\phi(\mathbf{x}_i)'\mathbf{w}_1 = \frac{1}{\sqrt{\lambda_1}}\mathbf{k}(\mathbf{x}_i,*)\mathbf{u}_1$$

4 D > 4 D > 4 E > 4 E > E 990

Center **K**? Use centering matrix **H**

$$H = I_N - \frac{(\mathbf{1}_N \mathbf{1}_N')}{N}$$
 $K_{center} = HKH$

Spirling (2013): model Treaties between US and Native Americans Why?

- American political development

- American political development
- IR Theories of Treaties and Treaty Violations

- American political development
- IR Theories of Treaties and Treaty Violations
- Comparative studies of indigenous/colonialist interaction

- American political development
- IR Theories of Treaties and Treaty Violations
- Comparative studies of indigenous/colonialist interaction
- Political Science question: how did Native Americans lose land so quickly?

How do we preserve word order and semantic language?

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us are identical.

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

How do we preserve word order and semantic language? After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order broad application

Consider documents x_i and x_j , where we have preserved order, punctuation, and all else.

Consider documents x_i and x_j , where we have preserved order, punctuation, and all else.

We say $\mathbf{x}_i \in \mathcal{X}$

Consider documents x_i and x_j , where we have preserved order, punctuation, and all else.

We say $\mathbf{x}_i \in \mathcal{X}$

Spirling examines 5-character strings, $s \in \mathcal{A}$

Consider documents x_i and x_j , where we have preserved order, punctuation, and all else.

We say $\mathbf{x}_i \in \mathcal{X}$

Spirling examines 5-character strings, $s \in \mathcal{A}$

Define:

Consider documents x_i and x_j , where we have preserved order, punctuation, and all else.

We say $\mathbf{x}_i \in \mathcal{X}$

Spirling examines 5-character strings, $s \in \mathcal{A}$

Define:

 $\phi_s: \mathcal{X} \to \Re$ as a function that counts the number of times string s occurs in document x.

Consider documents x_i and x_j , where we have preserved order, punctuation, and all else.

We say $\mathbf{x}_i \in \mathcal{X}$

Spirling examines 5-character strings, $s \in \mathcal{A}$

Define:

 $\phi_s: \mathcal{X} \to \Re$ as a function that counts the number of times string s occurs in document x.

Define string kernel to be,

Consider documents x_i and x_j , where we have preserved order, punctuation, and all else.

We say $\mathbf{x}_i \in \mathcal{X}$

Spirling examines 5-character strings, $s \in \mathcal{A}$

Define:

 $\phi_s: \mathcal{X} \to \Re$ as a function that counts the number of times string s occurs in document x.

Define string kernel to be,

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sum_{s \in \mathcal{A}} w_s \phi_s(\mathbf{x}_i) \phi_s(\mathbf{x}_j)$$

Consider documents x_i and x_j , where we have preserved order, punctuation, and all else.

We say $\mathbf{x}_i \in \mathcal{X}$

Spirling examines 5-character strings, $s \in \mathcal{A}$

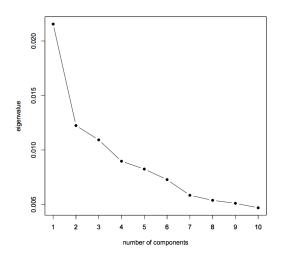
Define:

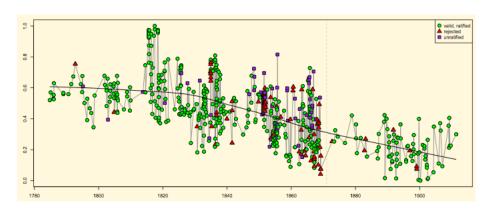
 $\phi_s: \mathcal{X} \to \Re$ as a function that counts the number of times string s occurs in document x.

Define string kernel to be,

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sum_{s \in \mathcal{A}} w_s \phi_s(\mathbf{x}_i) \phi_s(\mathbf{x}_j)$$

 $\phi(\emph{\textbf{x}}_i) pprox {32 \choose 5}$ element long count vector





- Suppose we have actor i (i = 1, 2, 3, ..., N)

- Suppose we have actor i $(i=1,2,3,\ldots,N)$
- Actor has ideal point $oldsymbol{ heta}_i \in \Re^M$

- Suppose we have actor i (i = 1, 2, 3, ..., N)
- Actor has ideal point $\theta_i \in \Re^M$
- We describe actor i's utility from proposal $\mathbf{p} \in \Re^{M}$ with utility function

- Suppose we have actor i (i = 1, 2, 3, ..., N)
- Actor has ideal point $\theta_i \in \Re^M$
- We describe actor i's utility from proposal $\mathbf{p} \in \Re^{M}$ with utility function

$$u_i(\boldsymbol{\theta}_i, \boldsymbol{p}) = -d(\boldsymbol{\theta}_i, \boldsymbol{p})$$

- Suppose we have actor i (i = 1, 2, 3, ..., N)
- Actor has ideal point $\theta_i \in \Re^M$
- We describe actor i's utility from proposal $\mathbf{p} \in \Re^{M}$ with utility function

$$u_i(\theta_i, \mathbf{p}) = -d(\theta_i, \mathbf{p})$$

$$= -\sum_{l=1}^{L} (\underbrace{\theta_{il}}_{\text{ideal policy}} - p_l)^2$$

- Suppose we have actor i (i = 1, 2, 3, ..., N)
- Actor has ideal point $\theta_i \in \Re^M$
- We describe actor i's utility from proposal $\mathbf{p} \in \Re^{M}$ with utility function

$$u_i(\theta_i, \mathbf{p}) = -d(\theta_i, \mathbf{p})$$

= $-\sum_{l=1}^{L} (\underbrace{\theta_{il}}_{\text{ideal policy}} -p_l)^2$

Estimation goal: $\widehat{\boldsymbol{\theta}}_i$

- Suppose we have actor i (i = 1, 2, 3, ..., N)
- Actor has ideal point $\theta_i \in \Re^M$
- We describe actor i's utility from proposal $\mathbf{p} \in \Re^{M}$ with utility function

$$u_i(\theta_i, \mathbf{p}) = -d(\theta_i, \mathbf{p})$$

= $-\sum_{l=1}^{L} (\underbrace{\theta_{il}}_{\text{ideal policy}} -p_l)^2$

Estimation goal: $\widehat{\theta}_i$ Scaling \leadsto placing actors in low-dimensional space (like principal components!)

US Congress and Roll Call

- Poole and Rosenthal voteview

- Poole and Rosenthal voteview
 - Roll Call Data → 1789-Present

- Poole and Rosenthal voteview
 - Roll Call Data → 1789-Present
 - NOMINATE methods → place legislators on one dimension, estimate of ideology

- Poole and Rosenthal voteview
 - Roll Call Data → 1789-Present
 - NOMINATE methods → place legislators on one dimension, estimate of ideology
- Wildly successful:

- Poole and Rosenthal voteview
 - Roll Call Data → 1789-Present
 - NOMINATE methods → place legislators on one dimension, estimate of ideology
- Wildly successful:
 - Estimates are accurate: face validity Congressional scholars agree upon

- Poole and Rosenthal voteview
 - Roll Call Data → 1789-Present
 - NOMINATE methods → place legislators on one dimension, estimate of ideology
- Wildly successful:
 - Estimates are accurate: face validity Congressional scholars agree upon
 - Insightful→ unidimensional Congress

- Poole and Rosenthal voteview
 - Roll Call Data → 1789-Present
 - NOMINATE methods → place legislators on one dimension, estimate of ideology
- Wildly successful:
 - Estimates are accurate: face validity Congressional scholars agree upon
 - Insightful → unidimensional Congress
 - Extensible: insight of IRT allows model to be embedded in many forms

- Poole and Rosenthal voteview
 - Roll Call Data → 1789-Present
 - NOMINATE methods → place legislators on one dimension, estimate of ideology
- Wildly successful:
 - Estimates are accurate: face validity Congressional scholars agree upon
 - Insightful → unidimensional Congress
 - Extensible: insight of IRT allows model to be embedded in many forms
 - Widely used: hard to write a paper on American political institutions with ideal points

Two Limitations with the NOMINATE project:

1) US Congress is distinct

Two Limitations with the NOMINATE project:

1) US Congress is distinct → roll call votes fail to measure ideology in other settings

- US Congress is distinct → roll call votes fail to measure ideology in other settings
 - Weak party pressure → individual discretion on votes

- US Congress is distinct → roll call votes fail to measure ideology in other settings
 - Weak party pressure → individual discretion on votes
 - Parliamentary systems → no discretion, no variation.

- US Congress is distinct → roll call votes fail to measure ideology in other settings
 - Weak party pressure → individual discretion on votes
 - Parliamentary systems → no discretion, no variation.
 - Spirling and Quinn (2011) → mixture model like models for blocs in UK Parliament

- US Congress is distinct → roll call votes fail to measure ideology in other settings
 - Weak party pressure → individual discretion on votes
 - Parliamentary systems→ no discretion, no variation.
 - Spirling and Quinn (2011) → mixture model like models for blocs in UK Parliament
- 2) Not everyone votes!

- US Congress is distinct → roll call votes fail to measure ideology in other settings
 - Weak party pressure → individual discretion on votes
 - Parliamentary systems → no discretion, no variation.
 - Spirling and Quinn (2011) → mixture model like models for blocs in UK Parliament
- 2) Not everyone votes!
 - Voters → survey responses (but problems with that)

- US Congress is distinct → roll call votes fail to measure ideology in other settings
 - Weak party pressure → individual discretion on votes
 - Parliamentary systems → no discretion, no variation.
 - Spirling and Quinn (2011) → mixture model like models for blocs in UK Parliament
- 2) Not everyone votes!
 - Voters → survey responses (but problems with that)
 - Challengers → NPAT surveys (but they don't fill those out anymore)

- US Congress is distinct → roll call votes fail to measure ideology in other settings
 - Weak party pressure → individual discretion on votes
 - Parliamentary systems → no discretion, no variation.
 - Spirling and Quinn (2011) → mixture model like models for blocs in UK Parliament
- 2) Not everyone votes!
 - Voters→ survey responses (but problems with that)
 - Challengers → NPAT surveys (but they don't fill those out anymore)
 - Bonica (2013, 2014) → estimate ideology from donations (but not everyone donates)

But Everyone talks!

- If we could scale based on conversation, we can measure ideology anywhere

- If we could scale based on conversation, we can measure ideology anywhere
- Much of the literature → relies upon intuition from US Congress

- If we could scale based on conversation, we can measure ideology anywhere
- Much of the literature → relies upon intuition from US Congress
 - Hard **not** to find ideology

- If we could scale based on conversation, we can measure ideology anywhere
- Much of the literature→ relies upon intuition from US Congress
 - Hard not to find ideology
 - Behavior that is primarily ideological

- If we could scale based on conversation, we can measure ideology anywhere
- Much of the literature→ relies upon intuition from US Congress
 - Hard not to find ideology
 - Behavior that is primarily ideological
- Reality: scaling is much more difficult than roll call voting examples

- If we could scale based on conversation, we can measure ideology anywhere
- Much of the literature → relies upon intuition from US Congress
 - Hard not to find ideology
 - Behavior that is primarily ideological
- Reality: scaling is much more difficult than roll call voting examples
 - Hard to find ideology

- If we could scale based on conversation, we can measure ideology anywhere
- Much of the literature → relies upon intuition from US Congress
 - Hard not to find ideology
 - Behavior that is primarily ideological
- Reality: scaling is much more difficult than roll call voting examples
 - Hard to find ideology
 - Much of political speech reveals little about position on ideological spectrum → advertising, regional

But Everyone talks!

- If we could scale based on conversation, we can measure ideology anywhere
- Much of the literature → relies upon intuition from US Congress
 - Hard not to find ideology
 - Behavior that is primarily ideological
- Reality: scaling is much more difficult than roll call voting examples
 - Hard to find ideology
 - Much of political speech reveals little about position on ideological spectrum
 → advertising, regional

Healthy skepticism!

Our plan

Our plan

1) Wordscores ~~ "supervised" scaling

Our plan

- 1) Wordscores --- "supervised" scaling
- 2) Wordfish → single dimension

Wordscores → Big in Europe

Wordscores, Like the Hoff→ Big in Europe



For each legislator i, suppose we observe D_i documents.

For each legislator i, suppose we observe D_i documents.

Define:

For each legislator i, suppose we observe D_i documents. Define:

$$\mathbf{x}_i = \sum_{l=1}^{D_i} \mathbf{x}_{il}$$

For each legislator i, suppose we observe D_i documents. Define:

$$\mathbf{x}_{i} = \sum_{l=1}^{D_{i}} \mathbf{x}_{il}$$

$$= \sum_{l=1}^{D_{i}} (x_{il1}, x_{il2}, \dots, x_{ilJ})$$

For each legislator i, suppose we observe D_i documents. Define:

$$\mathbf{x}_{i} = \sum_{l=1}^{D_{i}} \mathbf{x}_{il}$$

$$= \sum_{l=1}^{D_{i}} (x_{il1}, x_{il2}, \dots, x_{ilJ})$$

 $x_i \rightsquigarrow$ aggregation across documents.

Choose two legislators as exemplars

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$
- For example, might select Elizabeth Warren

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$
- For example, might select Elizabeth Warren
- Legislator $C \in \{1, 2, \dots, N\}$ is Conservative. $Y_C = 1$

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$
- For example, might select Elizabeth Warren
- Legislator $C \in \{1, 2, ..., N\}$ is Conservative. $Y_C = 1$
- For example, might select Ted Cruz

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$
- For example, might select Elizabeth Warren
- Legislator $C \in \{1, 2, \dots, N\}$ is Conservative. $Y_C = 1$
- For example, might select Ted Cruz

For each word j we can define:

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$
- For example, might select Elizabeth Warren
- Legislator $C \in \{1, 2, \dots, N\}$ is Conservative. $Y_C = 1$
- For example, might select Ted Cruz

For each word j we can define:

 P_{jL} = Probability of word from Liberal

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$
- For example, might select Elizabeth Warren
- Legislator $C \in \{1, 2, \dots, N\}$ is Conservative. $Y_C = 1$
- For example, might select Ted Cruz

For each word j we can define:

 P_{jL} = Probability of word from Liberal

 P_{jC} = Probability of word from Conservative

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$
- For example, might select Elizabeth Warren
- Legislator $C \in \{1, 2, \dots, N\}$ is Conservative. $Y_C = 1$
- For example, might select Ted Cruz

For each word j we can define:

 P_{iL} = Probability of word from Liberal

 P_{jC} = Probability of word from Conservative

Define the score for word j

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$
- For example, might select Elizabeth Warren
- Legislator $C \in \{1, 2, \dots, N\}$ is Conservative. $Y_C = 1$
- For example, might select Ted Cruz

For each word j we can define:

 P_{jL} = Probability of word from Liberal

 P_{jC} = Probability of word from Conservative

Define the score for word j

$$S_j = Y_C P_{jC} + Y_L P_{jL}$$

Choose two legislators as exemplars

- Legislator $L \in \{1, 2, \dots, N\}$ is liberal. $Y_L = -1$
- For example, might select Elizabeth Warren
- Legislator $C \in \{1, 2, \dots, N\}$ is Conservative. $Y_C = 1$
- For example, might select Ted Cruz

For each word j we can define:

 P_{jL} = Probability of word from Liberal

 P_{jC} = Probability of word from Conservative

Define the score for word j

$$S_j = Y_C P_{jC} + Y_L P_{jL}$$
$$= P_{jC} - P_{jL}$$

Scale other legislators:

Scale other legislators:

$$N_i = \sum_{j=1}^J x_i$$

Scale other legislators:

$$N_i = \sum_{j=1}^J x_i$$

 $\hat{ heta}_i$ is

Scale other legislators:

$$N_i = \sum_{j=1}^J x_i$$

 $\hat{ heta}_i$ is

$$\hat{\theta}_i = \sum_{j=1}^J \left(\frac{x_{ij}}{N_i}\right) S_j$$

Scale other legislators:

$$N_i = \sum_{j=1}^J x_i$$

 $\hat{ heta}_i$ is

$$\hat{\theta}_{i} = \sum_{j=1}^{J} \left(\frac{x_{ij}}{N_{i}}\right) S_{j}$$
$$= \frac{\mathbf{x}_{i}^{'}}{N_{i}} \mathbf{S}$$

$$N_{L} = \sum_{m=1}^{J} x_{mL}$$

$$N_{C} = \sum_{m=1}^{J} x_{mC}$$

$$N_{L} = \sum_{m=1}^{J} x_{mL}$$

$$N_{C} = \sum_{m=1}^{J} x_{mC}$$

$$N_{L} = \sum_{m=1}^{J} x_{mL}$$

$$N_{C} = \sum_{m=1}^{J} x_{mC}$$

$$P_{jL} = \frac{\frac{x_{jL}}{N_L}}{\frac{x_{jL}}{N_L} + \frac{x_{jC}}{N_C}}$$

$$N_{L} = \sum_{m=1}^{J} x_{mL}$$

$$N_{C} = \sum_{m=1}^{J} x_{mC}$$

$$P_{jL} = \frac{\frac{x_{jL}}{N_L}}{\frac{x_{jL}}{N_L} + \frac{x_{jC}}{N_C}}$$

$$P_{jC} = 1 - P_{jL} = \frac{\frac{x_{jC}}{N_C}}{\frac{x_{jL}}{N_L} + \frac{x_{jC}}{N_C}}$$

$$N_{L} = \sum_{m=1}^{J} x_{mL}$$

$$N_{C} = \sum_{m=1}^{J} x_{mC}$$

$$P_{jL} = \frac{\frac{x_{jL}}{N_L}}{\frac{x_{jL}}{N_L} + \frac{x_{jC}}{N_C}}$$

$$P_{jC} = 1 - P_{jL} = \frac{\frac{x_{jC}}{N_C}}{\frac{x_{jL}}{N_L} + \frac{x_{jC}}{N_C}}$$

$$S_j = P_{jC} - P_{jL}$$

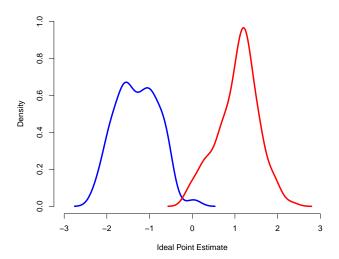
Applied to the Senate Press Releases

L = Ted Kennedy

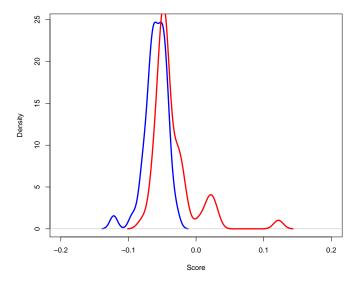
C = Tom Coburn

Apply to other senators.

Applying to Senate Press Releases → Gold Standard Scaling from NOMINATE



Applying to Senate Press Releases → WordScores



Wordscores is one example of supervised embedding:

Wordscores is one example of supervised embedding:

1) Linear Discriminant Analysis (LDA)→ the Federalist Papers

Wordscores is one example of supervised embedding:

- 1) Linear Discriminant Analysis (LDA) → the Federalist Papers
- 2) Multinomial Inverse Regression

Wordscores is one example of supervised embedding:

- 1) Linear Discriminant Analysis (LDA) → the Federalist Papers
- 2) Multinomial Inverse Regression
- 3) Any method that separates words

Wordscores is one example of supervised embedding:

- 1) Linear Discriminant Analysis (LDA) → the Federalist Papers
- 2) Multinomial Inverse Regression
- 3) Any method that separates words

The findings will depend heavily on supervision

Wordscores is one example of supervised embedding:

- 1) Linear Discriminant Analysis (LDA) → the Federalist Papers
- 2) Multinomial Inverse Regression
- 3) Any method that separates words

The findings will depend heavily on supervision

- Sensitive to who is chosen

Wordscores is one example of supervised embedding:

- 1) Linear Discriminant Analysis (LDA) → the Federalist Papers
- 2) Multinomial Inverse Regression
- 3) Any method that separates words

The findings will depend heavily on supervision

- Sensitive to who is chosen

Lowe (2008): Discusses potentially problematic wordscores properties

Lowe (2008): Discusses potentially problematic wordscores properties

1) Each word is weighted equally (fixable with different scoring procedure)

Lowe (2008): Discusses potentially problematic wordscores properties

- 1) Each word is weighted equally (fixable with different scoring procedure)
- Unique words are conflated with centrist (fixable with MCQ fightin' words style algorithm)

Lowe (2008): Discusses potentially problematic wordscores properties

- 1) Each word is weighted equally (fixable with different scoring procedure)
- Unique words are conflated with centrist (fixable with MCQ fightin' words style algorithm)
- 3) General problem: hard to interpret → lack of a model

Lowe (2008): Discusses potentially problematic wordscores properties

- 1) Each word is weighted equally (fixable with different scoring procedure)
- Unique words are conflated with centrist (fixable with MCQ fightin' words style algorithm)
- 3) General problem: hard to interpret → lack of a model

To be fair: fast, nonparametric, and novel [trailblazing] method for scoring documents (starts conversation)

Unsupervised Embedding

Basic idea:

Unsupervised Embedding

Basic idea:

- Actors have underlying latent position

Unsupervised Embedding

Basic idea:

- Actors have underlying latent position
- Actors articulate that latent position in their speech

Unsupervised Embedding

Basic idea:

- Actors have underlying latent position
- Actors articulate that latent position in their speech
- This is associated with word usage, so high discriminating words correspond to ideological speech

Unsupervised Embedding

Basic idea:

- Actors have underlying latent position
- Actors articulate that latent position in their speech
- This is associated with word usage, so high discriminating words correspond to ideological speech
- Some words discriminate better than others
 ⇔ encode that in our model

Unsupervised Embedding

Basic idea:

- Actors have underlying latent position
- Actors articulate that latent position in their speech
- This is associated with word usage, so high discriminating words correspond to ideological speech
- Some words discriminate better than others → encode that in our model

Simplest model: Principal Components

Consider press releases from 2005 US Senators

Consider press releases from 2005 US Senators Define $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ as the rate senator i uses J words.

Consider press releases from 2005 US Senators Define $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ as the rate senator i uses J words.

$$x_{ij} = \frac{\text{No. Times } i \text{ uses word } j}{\text{No. words } i \text{ uses}}$$

Consider press releases from 2005 US Senators Define $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ as the rate senator i uses J words.

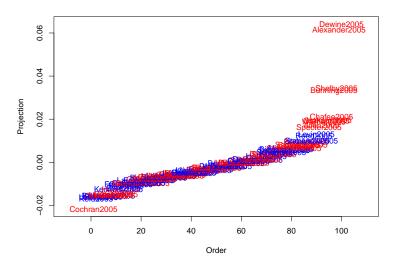
$$x_{ij} = \frac{\text{No. Times } i \text{ uses word } j}{\text{No. words } i \text{ uses}}$$

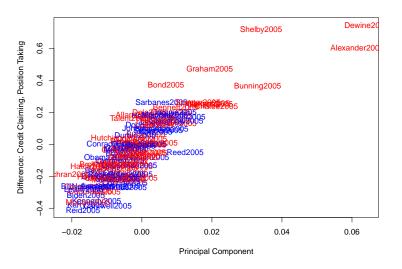
dtm: 100×2796 matrix containing word rates for senators

Consider press releases from 2005 US Senators Define $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ as the rate senator i uses J words.

$$x_{ij} = \frac{\text{No. Times } i \text{ uses word } j}{\text{No. words } i \text{ uses}}$$

dtm: 100×2796 matrix containing word rates for senators prcomp(dtm) applies principal components





Principal components is powerful

Principal components is powerful \leadsto statistical model for unsupervised scaling

Principal components is powerful → statistical model for unsupervised scaling

Principal components is powerful → statistical model for unsupervised scaling

Item Response Theory (IRT)

- Origins: educational testing

Principal components is powerful → statistical model for unsupervised scaling

- Origins: educational testing
- Jackman (2002), Clinton, Jackman, and Rivers (2004) apply to roll call voting

Principal components is powerful \leadsto statistical model for unsupervised scaling

- Origins: educational testing
- Jackman (2002), Clinton, Jackman, and Rivers (2004) apply to roll call voting
- Power of IRT:

Principal components is powerful → statistical model for unsupervised scaling

- Origins: educational testing
- Jackman (2002), Clinton, Jackman, and Rivers (2004) apply to roll call voting
- Power of IRT:
 - a) Estimate ideal points with few observations

Principal components is powerful → statistical model for unsupervised scaling

- Origins: educational testing
- Jackman (2002), Clinton, Jackman, and Rivers (2004) apply to roll call voting
- Power of IRT:
 - a) Estimate ideal points with few observations
 - b) Makes clear how to extend models

Principal components is powerful → statistical model for unsupervised scaling

- Origins: educational testing
- Jackman (2002), Clinton, Jackman, and Rivers (2004) apply to roll call voting
- Power of IRT:
 - a) Estimate ideal points with few observations
 - b) Makes clear how to extend models
 - c) Tuesday IRT and topic models to scale with votes + text

Principal components is powerful → statistical model for unsupervised scaling

- Origins: educational testing
- Jackman (2002), Clinton, Jackman, and Rivers (2004) apply to roll call voting
- Power of IRT:
 - a) Estimate ideal points with few observations
 - b) Makes clear how to extend models
 - c) Tuesday IRT and topic models to scale with votes + text
- Clinton, Jackman, and Rivers (2004) → intuition about IRT

Principal components is powerful → statistical model for unsupervised scaling

- Origins: educational testing
- Jackman (2002), Clinton, Jackman, and Rivers (2004) apply to roll call voting
- Power of IRT:
 - a) Estimate ideal points with few observations
 - b) Makes clear how to extend models
 - c) Tuesday IRT and topic models to scale with votes + text
- Clinton, Jackman, and Rivers (2004) → intuition about IRT
- Rivers (2002) → Identification conditions

Principal components is powerful → statistical model for unsupervised scaling

- Origins: educational testing
- Jackman (2002), Clinton, Jackman, and Rivers (2004) apply to roll call voting
- Power of IRT:
 - a) Estimate ideal points with few observations
 - b) Makes clear how to extend models
 - c) Tuesday IRT and topic models to scale with votes + text
- Clinton, Jackman, and Rivers (2004) → intuition about IRT
- Rivers (2002) → Identification conditions
- Bonica (2014a, 2014b) → uses IRT (like the one we're about to use) to scale donors

Principal components is powerful → statistical model for unsupervised scaling

Item Response Theory (IRT)

- Origins: educational testing
- Jackman (2002), Clinton, Jackman, and Rivers (2004) apply to roll call voting
- Power of IRT:
 - a) Estimate ideal points with few observations
 - b) Makes clear how to extend models
 - c) Tuesday IRT and topic models to scale with votes + text
- Clinton, Jackman, and Rivers (2004) → intuition about IRT
- Rivers (2002) → Identification conditions
- Bonica (2014a, 2014b) → uses IRT (like the one we're about to use) to scale donors

Monroe and Maeda (2005) and Slopkin and Proksch (2008) develop similar algorithms

Suppose we have legislator i.

Suppose we have legislator i.

$$x_{ij} \sim \mathsf{Poisson}(\lambda_{ij})$$

Suppose we have legislator i.

$$x_{ij} \sim \text{Poisson}(\lambda_{ij})$$

 $\lambda_{ij} = \exp(\alpha_i + \psi_j + \beta_j \times \theta_i)$

Suppose we have legislator i.

$$x_{ij} \sim \text{Poisson}(\lambda_{ij})$$

 $\lambda_{ij} = \exp(\alpha_i + \psi_j + \beta_j \times \theta_i)$

Where,

Suppose we have legislator i.

$$x_{ij} \sim \text{Poisson}(\lambda_{ij})$$

 $\lambda_{ij} = \exp(\alpha_i + \psi_j + \beta_j \times \theta_i)$

Where,

 λ_{ij} = Rate individual *i* uses word *j*

Suppose we have legislator i.

$$x_{ij} \sim \text{Poisson}(\lambda_{ij})$$

 $\lambda_{ij} = \exp(\alpha_i + \psi_j + \beta_j \times \theta_i)$

Where,

 λ_{ij} = Rate individual *i* uses word *j* α_i = Individual *i* loquaciousness

Suppose we have legislator i.

$$x_{ij} \sim \text{Poisson}(\lambda_{ij})$$

 $\lambda_{ij} = \exp(\alpha_i + \psi_j + \beta_j \times \theta_i)$

Where,

 $\lambda_{ij} = \text{Rate individual } i \text{ uses word } j$ $\alpha_i = \text{Individual } i \text{ loquaciousness}$ $\psi_j = \text{Word } j \text{'s frequency}$

Suppose we have legislator i.

$$x_{ij} \sim \text{Poisson}(\lambda_{ij})$$

 $\lambda_{ij} = \exp(\alpha_i + \psi_j + \beta_j \times \theta_i)$

Where,

 λ_{ij} = Rate individual i uses word j

 α_i = Individual *i* loquaciousness

 $\psi_j = \operatorname{Word} j$'s frequency

 β_j = Word j's discrimination

Suppose we have legislator i.

$$x_{ij} \sim \text{Poisson}(\lambda_{ij})$$

 $\lambda_{ij} = \exp(\alpha_i + \psi_j + \beta_j \times \theta_i)$

Where,

 λ_{ij} = Rate individual *i* uses word *j*

 α_i = Individual *i* loquaciousness

 ψ_j = Word j's frequency

 β_i = Word j's discrimination

 θ_i = Legislator i's latent positions

Suppose we have legislator i.

$$x_{ij} \sim \text{Poisson}(\lambda_{ij})$$

 $\lambda_{ij} = \exp(\alpha_i + \psi_j + \beta_j \times \theta_i)$

Where,

 λ_{ij} = Rate individual *i* uses word *j*

 α_i = Individual *i* loquaciousness

 ψ_i = Word j's frequency

 β_i = Word j's discrimination

 θ_i = Legislator i's latent positions

"Regression" of x_{ij} on ideal points θ_i , where we have to learn θ_i

Implies the following posterior distribution:

Implies the following posterior distribution:

$$p(\theta, \alpha, \psi, \beta) \propto p(\alpha)p(\beta)p(\psi)p(\theta) \times \prod_{i=1}^{N} \prod_{j=1}^{J} \frac{\exp\left[-(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i})\right](\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i})^{x_{ij}}}{x_{ij}!}$$

Implies the following posterior distribution:

$$p(\theta, \alpha, \psi, \beta) \propto p(\alpha)p(\beta)p(\psi)p(\theta) \times \prod_{i=1}^{N} \prod_{j=1}^{J} \frac{\exp\left[-\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)\right]\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)^{x_{ij}}}{x_{ij}!}$$

Estimate parameters:

Implies the following posterior distribution:

$$p(\theta, \alpha, \psi, \beta) \propto p(\alpha)p(\beta)p(\psi)p(\theta) \times \prod_{i=1}^{N} \prod_{j=1}^{J} \frac{\exp\left[-\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)\right]\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)^{x_{ij}}}{x_{ij}!}$$

Estimate parameters:

- EM-algorithm

Implies the following posterior distribution:

$$p(\theta, \alpha, \psi, \beta) \propto p(\alpha)p(\beta)p(\psi)p(\theta) \times \prod_{i=1}^{N} \prod_{j=1}^{J} \frac{\exp\left[-\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)\right]\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)^{x_{ij}}}{x_{ij}!}$$

Estimate parameters:

- EM-algorithm
- MCMC algorithm

Implies the following posterior distribution:

$$p(\theta, \alpha, \psi, \beta) \propto p(\alpha)p(\beta)p(\psi)p(\theta) \times \prod_{i=1}^{N} \prod_{j=1}^{J} \frac{\exp\left[-\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)\right]\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)^{x_{ij}}}{x_{ij}!}$$

Estimate parameters:

- EM-algorithm
- MCMC algorithm
- Variational Approximation

Implies the following posterior distribution:

$$p(\theta, \alpha, \psi, \beta) \propto p(\alpha)p(\beta)p(\psi)p(\theta) \times \prod_{i=1}^{N} \prod_{j=1}^{J} \frac{\exp\left[-\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)\right]\left(\alpha_{i} + \psi_{j} + \beta_{j} \times \theta_{i}\right)^{x_{ij}}}{x_{ij}!}$$

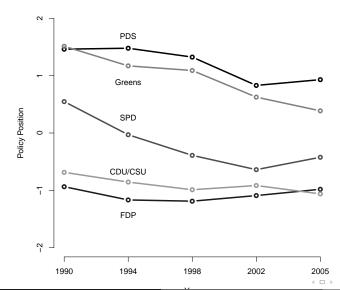
Estimate parameters:

- EM-algorithm
- MCMC algorithm
- Variational Approximation

Wordfish package in R

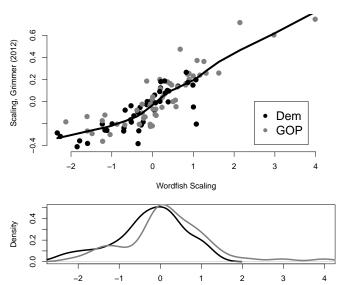
Applications: German Party Manifestos

Wordfish and German Platforms



Applications: German Party Manifestos

Wordfish and Senate Press Releases



The Problem with Text-Based Scaling

What does validation mean?

- 1) Replicate NOMINATE, DIME, or other gold standards?
- 2) Agreement with experts
- 3) Prediction of other behavior

Must answer this to make progress on pure text scaling