

# Text Data in Business and Economics

Basel University – Autumn 2023

## 5. Topic Models

# Different Goals, Different Methods

- ▶ Supervised Learning (later)
  - ▶ pursuing a known goal, e.g., predicting whether a political speech is from a Democrat or a Republican.
  - ▶ machine learns to replicate labels for new data points

# Different Goals, Different Methods

- ▶ Supervised Learning (later)
  - ▶ pursuing a known goal, e.g., predicting whether a political speech is from a Democrat or a Republican.
  - ▶ machine learns to replicate labels for new data points
- ▶ Unsupervised Learning
  - ▶ algorithm discovers themes/patterns in text (or other high-dimensional data)
  - ▶ human interprets the results (e.g. inspect content of topics or clusters)

# Different Goals, Different Methods

- ▶ Supervised Learning (later)
  - ▶ pursuing a known goal, e.g., predicting whether a political speech is from a Democrat or a Republican.
  - ▶ machine learns to replicate labels for new data points
- ▶ Unsupervised Learning
  - ▶ algorithm discovers themes/patterns in text (or other high-dimensional data)
  - ▶ human interprets the results (e.g. inspect content of topics or clusters)
- ▶ Both strategies amplify human effort, each in different ways.

# Different Goals, Different Methods

- ▶ Supervised Learning (later)
  - ▶ pursuing a known goal, e.g., predicting whether a political speech is from a Democrat or a Republican.
  - ▶ machine learns to replicate labels for new data points
- ▶ Unsupervised Learning
  - ▶ algorithm discovers themes/patterns in text (or other high-dimensional data)
  - ▶ human interprets the results (e.g. inspect content of topics or clusters)
- ▶ Both strategies amplify human effort, each in different ways.
- ▶ Distinctions are not clear-cut:
  - ▶ supervised learning models can be used to discover themes/patterns
  - ▶ unsupervised learning models can be used in service of prediction or known goals.

# Objectives: Social-Science Research using Unsupervised Learning

1. **What is the research question?**

# Objectives: Social-Science Research using Unsupervised Learning

1. **What is the research question?**
2. Corpus and Data:
  - ▶ obtain, clean, preprocess, and link.
  - ▶ Produce descriptive visuals and statistics on the text and metadata

# Objectives: Social-Science Research using Unsupervised Learning

1. **What is the research question?**
2. Corpus and Data:
  - ▶ obtain, clean, preprocess, and link.
  - ▶ Produce descriptive visuals and statistics on the text and metadata
3. Unsupervised learning:
  - ▶ **What are we trying to measure?**



# Objectives: Social-Science Research using Unsupervised Learning

1. **What is the research question?**
2. Corpus and Data:
  - ▶ obtain, clean, preprocess, and link.
  - ▶ Produce descriptive visuals and statistics on the text and metadata
3. Unsupervised learning:
  - ▶ **What are we trying to measure?**
  - ▶ Select a model and train it.
  - ▶ Probe sensitivity to hyperparameters.
  - ▶ Validate that the model is measuring what we want.

# Objectives: Social-Science Research using Unsupervised Learning

1. **What is the research question?**
2. Corpus and Data:
  - ▶ obtain, clean, preprocess, and link.
  - ▶ Produce descriptive visuals and statistics on the text and metadata
3. Unsupervised learning:
  - ▶ **What are we trying to measure?**
  - ▶ Select a model and train it.
  - ▶ Probe sensitivity to hyperparameters.
  - ▶ Validate that the model is measuring what we want.
4. Empirical analysis
  - ▶ Produce statistics or predictions with the trained model.
  - ▶ **Answer the research question.**

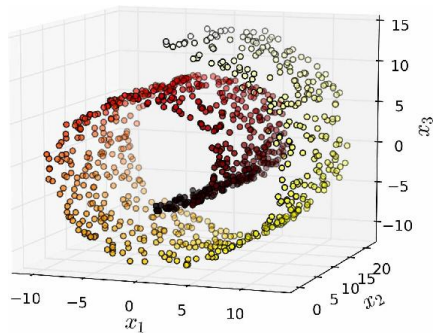
# Outline

Dimensionality Reduction

Topic Models

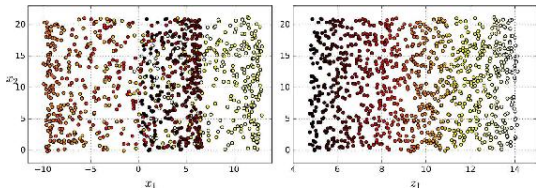
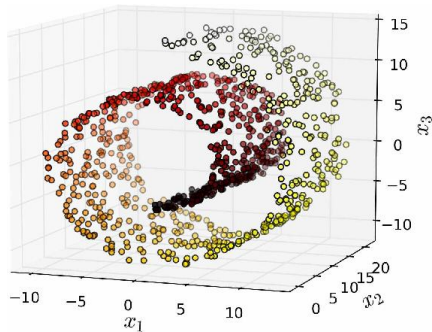
- ▶ Datasets are not distributed uniformly across the feature space.
- ▶ They have a lower-dimensional latent structure – a **manifold** – that can be learned.

“The Swiss Roll”



## “The Swiss Roll”

- ▶ Datasets are not distributed uniformly across the feature space.
- ▶ They have a lower-dimensional latent structure – a **manifold** – that can be learned.



- ▶ **Dimensionality reduction** makes data more interpretable – for example by projecting down to two dimensions for visualization.
- ▶ improves computational tractability.
- ▶ can improve model performance.

What dimension reductions have you already tried in this class?

# The Document-Term Matrix is high-dimensional

The **document-term matrix**  $\mathbf{X}$ :

- ▶ each row  $d$  represents a **document**, while each column  $w$  represents a word (or term more generally, e.g. n-grams).
- ▶ A matrix entry  $\mathbf{X}_{[d,w]}$  quantifies the strength of association between a document and a word, generally its count or frequency

# The Document-Term Matrix is high-dimensional

The **document-term matrix**  $\mathbf{X}$ :

- ▶ each row  $d$  represents a **document**, while each column  $w$  represents a word (or term more generally, e.g. n-grams).
  - ▶ A matrix entry  $\mathbf{X}_{[d,w]}$  quantifies the strength of association between a document and a word, generally its count or frequency
- ▶ each document/row  $\mathbf{X}_{[d,:]}$  is a distribution over terms
  - ▶ term vocabularies can be in the hundreds of thousands



# The Document-Term Matrix is high-dimensional

The **document-term matrix**  $\mathbf{X}$ :

- ▶ each row  $d$  represents a **document**, while each column  $w$  represents a word (or term more generally, e.g. n-grams).
  - ▶ A matrix entry  $\mathbf{X}_{[d,w]}$  quantifies the strength of association between a document and a word, generally its count or frequency
- ▶ each document/row  $\mathbf{X}_{[d,:]}$  is a distribution over terms
  - ▶ term vocabularies can be in the hundreds of thousands
- ▶ each word/column  $\mathbf{X}_{[:,w]}$  is a distribution over documents.
  - ▶ many interesting corpora have millions of documents

# The Document-Term Matrix is high-dimensional

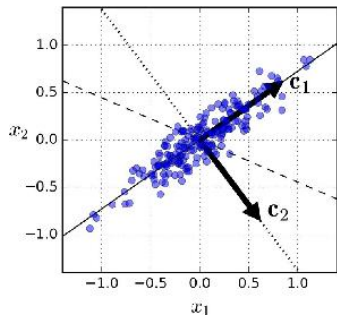
The **document-term matrix**  $\mathbf{X}$ :

- ▶ each row  $d$  represents a **document**, while each column  $w$  represents a word (or term more generally, e.g. n-grams).
  - ▶ A matrix entry  $\mathbf{X}_{[d,w]}$  quantifies the strength of association between a document and a word, generally its count or frequency
- ▶ each document/row  $\mathbf{X}_{[d,:]}$  is a distribution over terms
  - ▶ term vocabularies can be in the hundreds of thousands
- ▶ each word/column  $\mathbf{X}_{[:,w]}$  is a distribution over documents.
  - ▶ many interesting corpora have millions of documents

→  $\mathbf{X}$  often has billions of cells.

# PCA (principal component analysis)

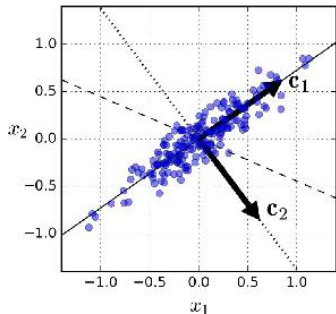
## PCA (principal component analysis)



- PCA computes the dimension in data explaining most variance.

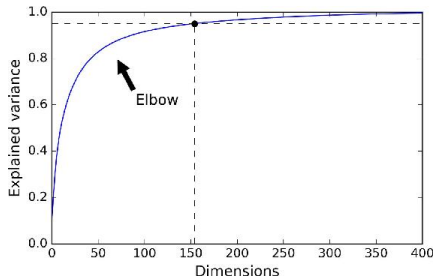
```
from sklearn.decomposition import PCA  
pca = PCA(n_components=10)  
X_train_pca = pca.fit_transform(X_train)
```

# PCA (principal component analysis)



- PCA computes the dimension in data explaining most variance.

```
from sklearn.decomposition import PCA  
pca = PCA(n_components=10)  
X_train_pca = pca.fit_transform(X_train)
```



- after the first component, subsequent components learn the (orthogonal) dimensions explaining most variance in dataset after projecting out first component.

# PCA and LSA

The document-term matrix  $\mathbf{X}$  can be reduced by projecting down to first principal component dimensions.

- ▶ This is known as “latent semantic analysis”
- ▶ Distance metrics between observations (e.g. cosine similarity) are approximately preserved.

# PCA and LSA

The document-term matrix  $\mathbf{X}$  can be reduced by projecting down to first principal component dimensions.

- ▶ This is known as “latent semantic analysis”
- ▶ Distance metrics between observations (e.g. cosine similarity) are approximately preserved.
- ▶ PCA factors are not interpretable.
  - ▶ Hoffman (1999) fixes this and puts LSA on firmer foundations by assuming a generative model of text – the word counts in a document are generated by a multinomial distribution.
  - ▶ For non-negative data (e.g. counts or frequencies), **Non-negative Matrix Factorization (NMF)** provides more interpretable factors than PCA.

# Outline

Dimensionality Reduction

Topic Models



# Topic Models in Social Science

- ▶ Core methods for topic models were developed in computer science and statistics
  - ▶ summarize unstructured text
  - ▶ use words within document to infer subject
  - ▶ useful for dimension reduction

# Topic Models in Social Science

- ▶ Core methods for topic models were developed in computer science and statistics
  - ▶ summarize unstructured text
  - ▶ use words within document to infer subject
  - ▶ useful for dimension reduction
- ▶ Social scientists use topics as a form of measurement
  - ▶ how observed covariates drive trends in language
  - ▶ tell a story not just about what, but how and why

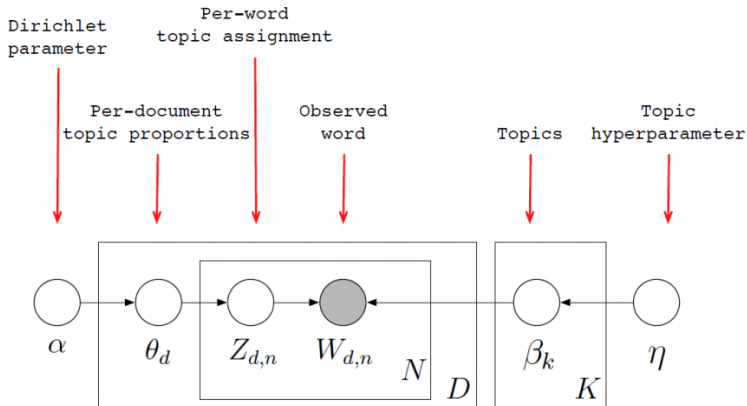
# Topic Models in Social Science

- ▶ Core methods for topic models were developed in computer science and statistics
  - ▶ summarize unstructured text
  - ▶ use words within document to infer subject
  - ▶ useful for dimension reduction
- ▶ Social scientists use topics as a form of measurement
  - ▶ how observed covariates drive trends in language
  - ▶ tell a story not just about what, but how and why
  - ▶ **topic models are more interpretable** than other dimension reduction methods, such as PCA.

- ▶ Latent Dirichlet Allocation (LDA):
  - ▶ Each topic is a distribution over words.
  - ▶ Each document is a distribution over topics.

- ▶ Latent Dirichlet Allocation (LDA):
  - ▶ Each topic is a distribution over words.
  - ▶ Each document is a distribution over topics.
- ▶ Input:  $N \times M$  document-term count matrix  $X$
- ▶ Assume: there are  $K$  topics (tunable hyperparameter, use coherence).
- ▶ Like PCA or NMF, LDA works by factorizing  $X$  into:
  - ▶ an  $N \times K$  document-topic matrix
  - ▶ an  $K \times M$  topic-term matrix.

- ▶ Latent Dirichlet Allocation (LDA):
  - ▶ Each topic is a distribution over words.
  - ▶ Each document is a distribution over topics.
- ▶ Input:  $N \times M$  document-term count matrix  $X$
- ▶ Assume: there are  $K$  topics (tunable hyperparameter, use coherence).
- ▶ Like PCA or NMF, LDA works by factorizing  $X$  into:
  - ▶ an  $N \times K$  document-topic matrix
  - ▶ an  $K \times M$  topic-term matrix.



# **ldagibbs: A command for Topic Modeling in Stata using Latent Dirichlet Allocation**

Carlo Schwarz

University of Warwick

Coventry, United Kingdom

c.r.schwarz@warwick.ac.uk

```
ldagibbs varname [ , ttopics(integer) burnin_iter(integer) alpha(real)  
  beta(real) samples(integer) sampling_iter(integer) seed(integer)  
  likelihood min_char(integer) stopwords(string) name_new_var(string)  
  normalize mat_save path(string) ]
```

# A statistical highlighter

## Seeking Life's Bare (Genetic) Necessities

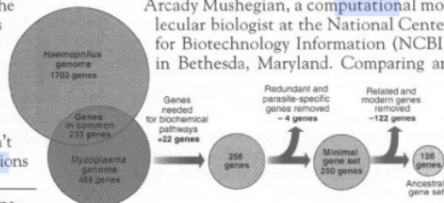
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known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

Although the numbers don't match precisely, those predictions

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

“are not all that far apart,” especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. “It may be a way of organizing any newly sequenced genome,” explains

Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



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



## Using an LDA Model

Once trained, can easily get topic proportions for a corpus.

- ▶ for any document – doesn't have to be in training corpus.
- ▶ main topic is the highest-probability topic

## Using an LDA Model

Once trained, can easily get topic proportions for a corpus.

- ▶ for any document – doesn't have to be in training corpus.
- ▶ main topic is the highest-probability topic
- ▶ documents with highest share in a topic work as representative documents for the topic.

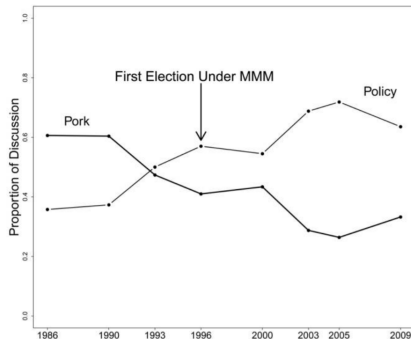
# Using an LDA Model

Once trained, can easily get topic proportions for a corpus.

- ▶ for any document – doesn't have to be in training corpus.
- ▶ main topic is the highest-probability topic
- ▶ documents with highest share in a topic work as representative documents for the topic.

Can then use the topic proportions as variables in a social science analysis.

- ▶ e.g., Catalinac (2016) shows that after a Japanese political reform that reduced intraparty competition, candidate platforms reduced pork-barrel policies and increased national ones.



**TABLE 1 A Summary of Common Assumptions and Relative Costs Across Different Methods of Discrete Text Categorization**

	Method				
	<i>Reading</i>	<i>Human Coding</i>	<i>Dictionaries</i>	<i>Supervised Learning</i>	<i>Topic Model</i>
<b>A. Assumptions</b>					
<i>Categories are known</i>	No	Yes	Yes	Yes	No
<i>Category nesting, if any, is known</i>	No	Yes	Yes	Yes	No
<i>Relevant text features are known</i>	No	No	Yes	Yes	Yes
<i>Mapping is known</i>	No	No	Yes	No	No
<i>Coding can be automated</i>	No	No	Yes	Yes	Yes
<b>B. Costs</b>					
Preanalysis Costs					
<i>Person-hours spent conceptualizing</i>	Low	High	High	High	Low
<i>Level of substantive knowledge</i>	Moderate/High	High	High	High	Low
Analysis Costs					
<i>Person hours spent per text</i>	High	High	Low	Low	Low
<i>Level of substantive knowledge</i>	Moderate/High	Moderate	Low	Low	Low
Postanalysis Costs					
<i>Person-hours spent interpreting</i>	High	Low	Low	Low	Moderate
<i>Level of substantive knowledge</i>	High	High	High	High	High

Recommended: read this part of Quinn, Monroe, Colaresi, Crespin, and Radev (2010).

# Structural Topic Model = LDA + Metadata

Roberts, Stewart, and Tingley

STM provides two ways to include contextual information:

- ▶ Topic prevalence can vary by metadata
  - ▶ e.g. Republicans talk about military issues more than Democrats

# Structural Topic Model = LDA + Metadata

Roberts, Stewart, and Tingley

STM provides two ways to include contextual information:

- ▶ Topic prevalence can vary by metadata
  - ▶ e.g. Republicans talk about military issues more than Democrats
- ▶ Topic content can vary by metadata
  - ▶ e.g. Republicans talk about military issues more patriotically than Democrats.

# Structural Topic Model = LDA + Metadata

Roberts, Stewart, and Tingley

STM provides two ways to include contextual information:

- ▶ Topic prevalence can vary by metadata
  - ▶ e.g. Republicans talk about military issues more than Democrats
- ▶ Topic content can vary by metadata
  - ▶ e.g. Republicans talk about military issues more patriotically than Democrats.
- ▶ Structural topic model is not a prediction model:
  - ▶ it will tell you which topics or features correlate with an outcome, but it will not provide an in-sample or out-of-sample prediction for an outcome

# Structural Topic Model = LDA + Metadata

Roberts, Stewart, and Tingley

STM provides two ways to include contextual information:

- ▶ Topic prevalence can vary by metadata
  - ▶ e.g. Republicans talk about military issues more than Democrats
- ▶ Topic content can vary by metadata
  - ▶ e.g. Republicans talk about military issues more patriotically than Democrats.
- ▶ Structural topic model is not a prediction model:
  - ▶ it will tell you which topics or features correlate with an outcome, but it will not provide an in-sample or out-of-sample prediction for an outcome
- ▶ The main implementation is in R. gensim has a light-weight version called “author topic model”.