

# Unstructured Data for Economics

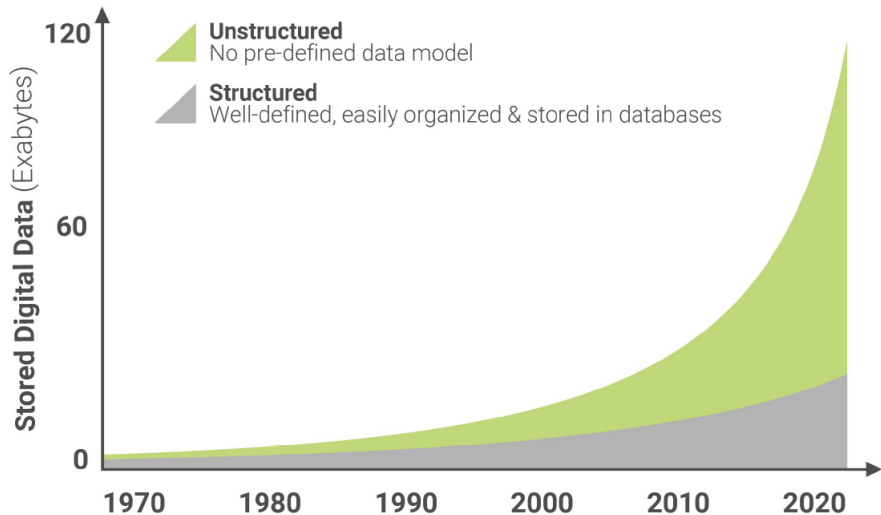
## Lecture 1: Introduction and Bag-of-Words Model

Stephen Hansen  
University College London



FONDAZIONE LUIGI EINAUDI  
PER STUDI DI POLITICA ECONOMIA E STORIA

# Trends in Data Types



# Unstructured Data

Unstructured data does not come organized in a traditional relational database.

Extracting relevant information and separating it from irrelevant information is a primary challenge.

Examples:

1. Text
2. Audio
3. Images
4. Videos

# Happenstance Data<sup>1</sup>

Traditional economic data is constructed with a particular measurement in mind, e.g. GDP statistics.

Most data generated in the private sector is happenstance, and arises via the everyday activities of agents (“digital exhaust”).

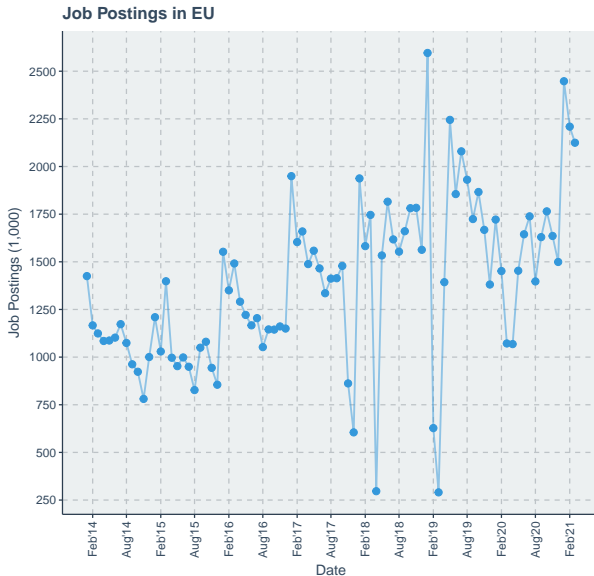
Statistical challenge is that data is not collected with a consistent, representative sample frame.

Organizational challenge is that data access arrangements have yet to be normalized.

---

<sup>1</sup>Discussed more fully in <https://rs-delve.github.io/reports/2020/11/24/data-readiness-lessons-from-an-emergency.html>

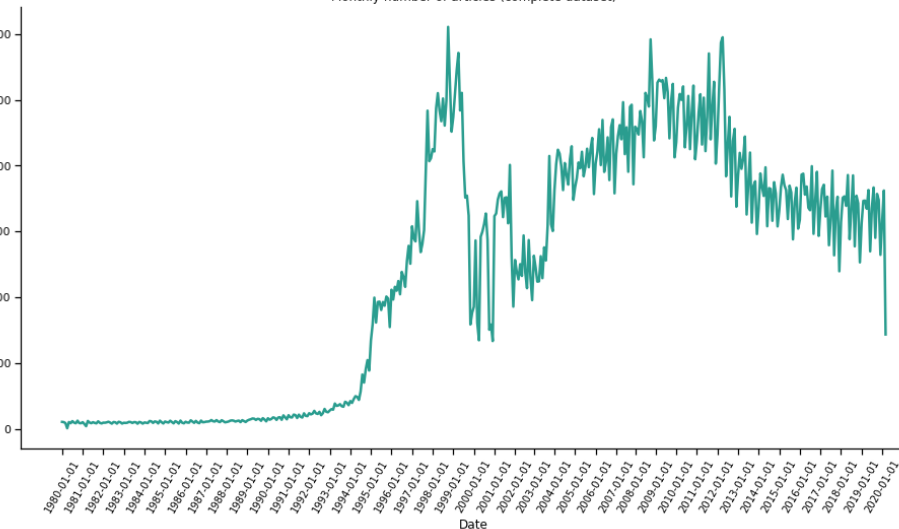
# Monthly Online Job Postings



Source: Web-scraped Job Ads from EU27 Countries, provided by Burning Glass Technologies.

# Monthly Newswire Postings

Monthly number of articles (complete dataset)



# Unstructured vs Happenstance Data

	Administrative	Happenstance
Structured	Traditional Economic Data	Credit Card Transactions Amazon product ratings
Unstructured	10-K Filings FOMC press conferences	Tweets Online Job Postings

# What is the Value of Unstructured Data?

The main application of unstructured data in economics and related disciplines has been to **measure** important phenomena.

Can **complement existing measures**: e.g. build more granular versions of official data.

Or **create entirely new measures**: economic policy uncertainty, media slant, central bank communication.

Makes information retrieval methods useful in a wide variety of fields.

**Emerging application area**: input into econometric models for inference and structural estimation.



# This Course

1. Bag-of-words model + topic models
2. Word embeddings
3. Sequence embeddings and large language models
4. Finetuning LLMs
5. Econometrics with unstructured data

Also applications to non-textual data.

# What is Text?

At an abstract level, text is simply a string of characters.

Some of these may be from the Latin alphabet—‘a’, ‘A’, ‘p’ and so on—but there may also be:

1. Decorated Latin letters (e.g. ö)
2. Non-Latin alphabetic characters (e.g. Chinese and Arabic)
3. Punctuation (e.g. ‘!’)
4. White spaces, tabs, newlines
5. Numbers
6. Non-alphanumeric characters (e.g. ‘@’)

**Key Question:** How can we obtain an informative, quantitative representation of these character strings?

# How to Preprocess?

First step is to **pre-process** strings to convert them into lists of units of meaning, sometimes called **tokens**.

Preprocessing in traditional text-as-data analysis follows a standard sequence of steps:

1. Break strings into initial tokens.
2. Remove common and rare words.
3. Fold into common linguistic roots (lower case, stem/lemmatize)

See notebook for basic demonstration.

See also [Denny and Spirling, 2018].

# How to Preprocess?

First step is to **pre-process** strings to convert them into lists of units of meaning, sometimes called **tokens**.

Preprocessing in traditional text-as-data analysis follows a standard sequence of steps:

1. Break strings into initial tokens.
2. Remove common and rare words.
3. Fold into common linguistic roots (lower case, stem/lemmatize)

See notebook for basic demonstration.

See also [Denny and Spirling, 2018].

In LLMs there is much more limited pre-preprocessing in order to preserve meaning of language; see notebook for more details.

# Which Corpus?

Much of traditional text-as-data analysis fits models on corpora drawn from domain of interest.

Large language models were first fit on generic corpora like Common Crawl, Wikipedia, or open-source books.

More recent iterations expand the training data (but details becoming more obscure).

Important to realize that **the training data contains the knowledge that a model can encode.**

Any biases in the training data can also be inherited by the model.

# Bias in GPT-3

Prompt GPT-3 with He was very [MASK] and She was very [MASK].

**Table 6.1:** Most Biased Descriptive Words in 175B Model

Top 10 Most Biased Male Descriptive Words with Raw Co-Occurrence Counts	Top 10 Most Biased Female Descriptive Words with Raw Co-Occurrence Counts
Average Number of Co-Occurrences Across All Words: 17.5	Average Number of Co-Occurrences Across All Words: 23.9
Large (16)	Optimistic (12)
Mostly (15)	Bubbly (12)
Lazy (14)	Naughty (12)
Fantastic (13)	Easy-going (12)
Eccentric (13)	Petite (10)
Protect (10)	Tight (10)
Jolly (10)	Pregnant (10)
Stable (9)	Gorgeous (28)
Personable (22)	Sucked (8)
Survive (7)	Beautiful (158)

See [Brown et al., 2020] for more details.

GPT-3 trained on Common Crawl, WebText2, Books1, Books2, Wikipedia.

# Notation

The corpus is composed of  $D$  documents indexed by  $d$ .

After pre-processing, each document is a finite, length- $N_d$  list of terms  $\mathbf{w}_d = (w_{d,1}, \dots, w_{d,N_d})$  with generic element  $w_{d,n}$ .

Let  $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_D)$  be a list of all terms in the corpus, and let  $N \equiv \sum_d N_d$  be the total number of terms in the corpus.

Suppose there are  $V$  **unique** terms in  $\mathbf{w}$ , where  $1 \leq V \leq N$ , each indexed by  $v$ .

We can then map each term in the corpus into this index, so that  $w_{d,n} \in \{1, \dots, V\}$ .

Let  $x_{d,v}$  be the count of term  $v$  in document  $d$ .

# Example

Consider three documents:

1. 'stephen is nice'
2. 'john is also nice'
3. 'george is mean'

We can consider the set of unique terms as {stephen, is, nice, john, also, george, mean} so that  $V = 7$ .

Construct the following index:

stephen	is	nice	john	also	george	mean
1	2	3	4	5	6	7

We then have  $\mathbf{w}_1 = (1, 2, 3)$ ;  $\mathbf{w}_2 = (4, 2, 5, 3)$ ;  $\mathbf{w}_3 = (6, 2, 7)$ .

Moreover  $x_{1,1} = 1$ ,  $x_{2,1} = 0$ ,  $x_{3,1} = 0$ , etc.



# Bag-of-Words Model

# Document-Term Matrix

A popular quantitative representation of text is the *document-term matrix*  $\mathbf{X}$ , which collects the counts  $x_{d,v}$  into a  $D \times V$  matrix.

In the previous example, we have

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

# Real-World Example

In “Transparency and Deliberation” we use a corpus of verbatim FOMC transcripts from the era of Alan Greenspan:

- ▶ 149 meetings from August 1987 through January 2006.
- ▶ A document is a single statement by a speaker in a meeting (46,502).
- ▶ Associated metadata: speaker biographical information, macroeconomic conditions, etc.

# Executive Time Use Project

Data on each 15-minute block of time for one week of 1,114 CEOs' time classified according to

1. type (e.g. meeting, public event, etc.)
2. duration (15m, 30m, etc.)
3. planning (planned or unplanned)
4. number of participants (one, more than one)
5. functions of participants, divided between employees of the firms or “insiders” (finance, marketing, etc.) and “outsiders” (clients, banks, etc.).

There are 4,253 unique combinations of these five features in the data.

One can summarize the data with a  $1114 \times 4253$  matrix where the  $(i, j)$ th element is the number of 15-minute time blocks that CEO  $i$  spends in activities with a particular combination of features  $j$ .

# Other Examples

Network data can be represented by an **adjacency matrix** which is typically high dimensional, sparse, and discrete.

**Bag-of-visual words** model in image processing.

# Older is Sometimes Better!

Task/Language		best ZSL	supervised	
			Standard ML	Transformer
SA	EN	0.553	0.610	<b>0.680</b>
	DE	0.517	0.610	<b>0.677</b>
	FR	0.528	0.612	<b>0.706</b>
AC-Gender	EN	0.624	0.601	<b>0.638</b>
	DE	0.497	0.540	<b>0.629</b>
	FR	0.579	0.546	<b>0.650</b>
AC-Age	EN	0.572	0.620	<b>0.636</b>
	DE	0.503	0.602	<b>0.611</b>
	FR	0.550	0.540	<b>0.568</b>

From [Plaza-del-Arco et al., 2024].

# Four Measurement Problems

[Ash and Hansen, 2023] organize measurement problems associated with text into four categories:

1. Distance between documents, e.g. how similar are corporate filings from each other.
2. Whether a concept is present (and degree of presence) in a document, e.g. sentiment.
3. How concepts relate in a document, e.g. sentiment and individual companies.
4. Associating documents to metadata, e.g. mapping newspaper text into recession vs expansion periods.

# Document Similarity



# Documents as Vectors

We can view the documents that make up the rows of  $\mathbf{X}$  as vectors.

Let each vocabulary term  $v$  have its own vector  $\mathbf{e}_v \in \mathbb{R}^V$  where

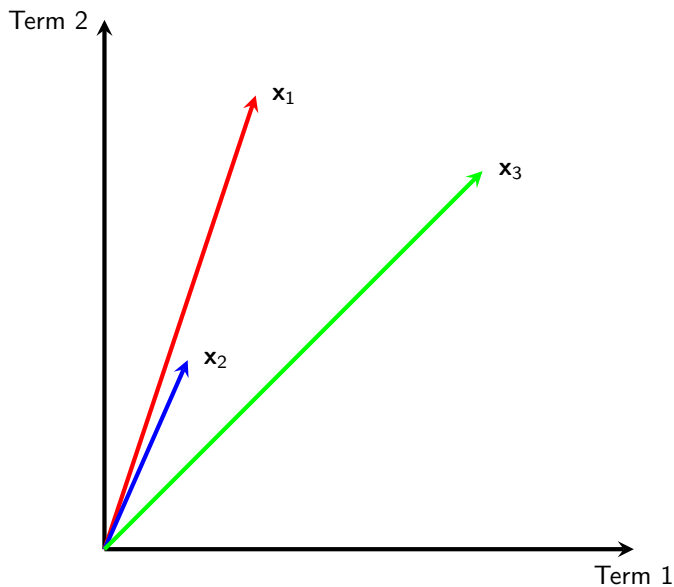
$$e_{v,v'} = \begin{cases} 1 & \text{if } v = v' \\ 0 & \text{otherwise} \end{cases}$$

Note that each term's vector is orthogonal to every other term's vector.

We can express document  $d$  as

$$\mathbf{x}_d = x_{d,1}\mathbf{e}_1 + x_{d,2}\mathbf{e}_2 + \dots + x_{d,V}\mathbf{e}_V$$

# Three Documents



# Distance in the Vector Space

An initial question of interest is how similar are any two documents in the vector space.

Initial instinct might be to use Euclidean distance  $\sqrt{\sum_v (x_{i,v} - x_{j,v})^2}$ .

What is the problem with Euclidean distance? How can we correct this?

# Cosine Similarity

Define the cosine similarity between documents  $i$  and  $j$  as

$$CS(i, j) = \frac{\mathbf{x}_i \cdot \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$$

1. Since document vectors have no negative elements  $CS(i, j) \in [0, 1]$ .
2.  $\mathbf{x}_i / \|\mathbf{x}_i\|$  is unit-length, correction for different distances.

# Application

An important theoretical concept in industrial organization is location on a product space.

Industry classification measures are quite crude proxies of this.

[Hoberg and Phillips, 2010] and [Hoberg and Phillips, 2016] take product descriptions from 49,408 10-K filings and use the vector space model to compute similarity between firms.

Data available from <http://alex2.umd.edu/industrydata/>.

# Term Weighting

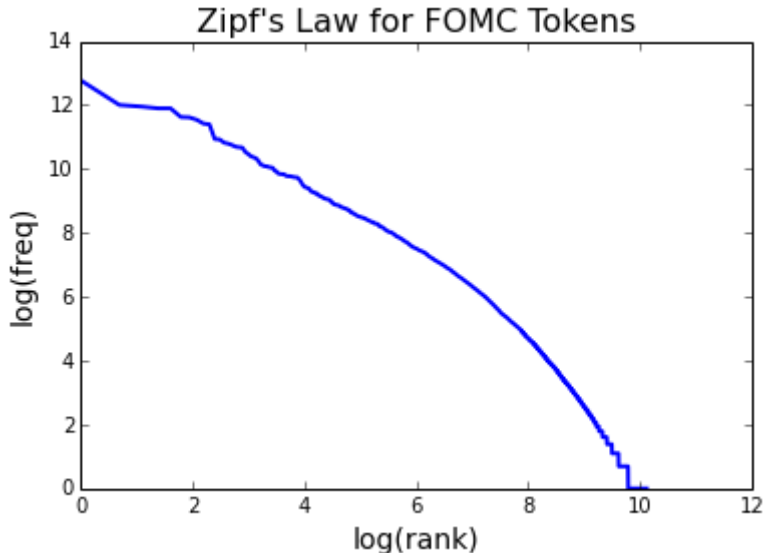
The frequency of words in natural language can distort raw counts.

Zipf's Law is an empirical regularity for natural language: the frequency of a particular term is inversely proportional to its rank.

Means that a few terms will have very large counts, many terms have small counts.

Example of a *power law*.

# Zipf's Law in FOMC Transcript Data



# Rescaling Counts

Let  $x_{d,v}$  be the count of the  $v$ th term in document  $d$ .

To dampen the power-law effect can express counts as

$$tf_{d,v} = \begin{cases} 0 & \text{if } x_{d,v} = 0 \\ 1 + \log(x_{d,v}) & \text{otherwise} \end{cases}$$

which is the *term frequency* of  $v$  in  $d$ .



# Inverse Document Frequency

Let  $df_v$  be the number of documents that contain the term  $v$ .

The *inverse document frequency* is

$$\text{idf}_v = \log \left( \frac{D}{df_v} \right),$$

where  $D$  is the number of documents.

Properties:

1. Higher weight for words in fewer documents.
2. Log dampens effect of weighting.

# TF-IDF Weighting

Combining the two observations from above allows us to express the *term frequency - inverse document frequency* of term  $v$  in document  $d$  as

$$\text{tf-idf}_{d,v} = \text{tf}_{d,v} \times \text{idf}_v.$$

Gives prominence to words that occur many times in few documents.

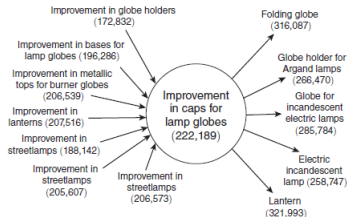
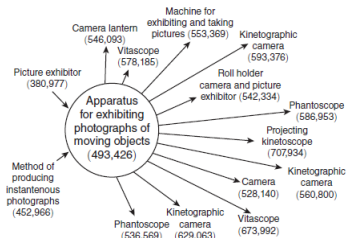
# Application

[Kelly et al., 2021] uses the text of US patents to identify radical innovation.

An individual patent is said to be influential when its **backward similarity** is low and its **forward similarity** is high.

Measure validated with historically important patents, forward citations, market value.

# Similarity Networks



# Concept Detection

# Dictionary Methods

The most common strategy for concept detection is to define a list of terms that capture the concept of interest, and to express documents as counts over those terms.

Strategy is referred to as *dictionary methods*.

Where do the dictionaries come from?

1. Pre-defined lists → [Harvard General Inquirer](#) [Tetlock, 2007]
2. Domain expertise → [Loughran and McDonald, 2011]
3. Ability to predict objective label → form of supervised learning.

# Multiple Dictionaries can Improve Performance

**Table 2**  
Goodness-of-Fit of lexical model sentiment scores for predicting human ratings.

Model	Feature space	Lexicon size	Ordered-Logit pseudo $R^2$	OLS $R^2$	Rank correlation	Macro-F1
GI Lexicon	General english	3626	0.023	0.064	0.264	0.406
+ Negation rule			0.029	0.080	0.295	0.432
LM Lexicon	10-K Reports	2707	0.065	0.165	0.447	0.510
+ Negation rule			0.066	0.169	0.449	0.500
HL Lexicon	Movie reviews	6786	0.066	0.173	0.437	0.509
+ Negation rule			0.072	0.186	0.453	0.503
GI + LM + HL Lexicon	Combined	9570	0.063	0.163	0.426	0.497
+ Negation rule			0.070	0.180	0.444	0.500
LM + HL Lexicon	Combined	8453	0.077	0.195	0.476	0.516
+ Negation rule			0.081	0.205	0.486	0.514

From [Shapiro et al., 2022].

# Predicting Labels

The Economic Policy Uncertainty (EPU) index of [Baker et al., 2016] (<http://www.policyuncertainty.com/>) builds a dictionary based on human labels.

Human audit reveals that nearly all policy uncertainty articles satisfy following criteria:

1. Article contains “uncertain” OR “uncertainty”, AND
2. Article contains “economic” OR “economy” AND



# Predicting Labels

The Economic Policy Uncertainty (EPU) index of [Baker et al., 2016] (<http://www.policyuncertainty.com/>) builds a dictionary based on human labels.

Human audit reveals that nearly all policy uncertainty articles satisfy following criteria:

1. Article contains “uncertain” OR “uncertainty”, AND
2. Article contains “economic” OR “economy” AND

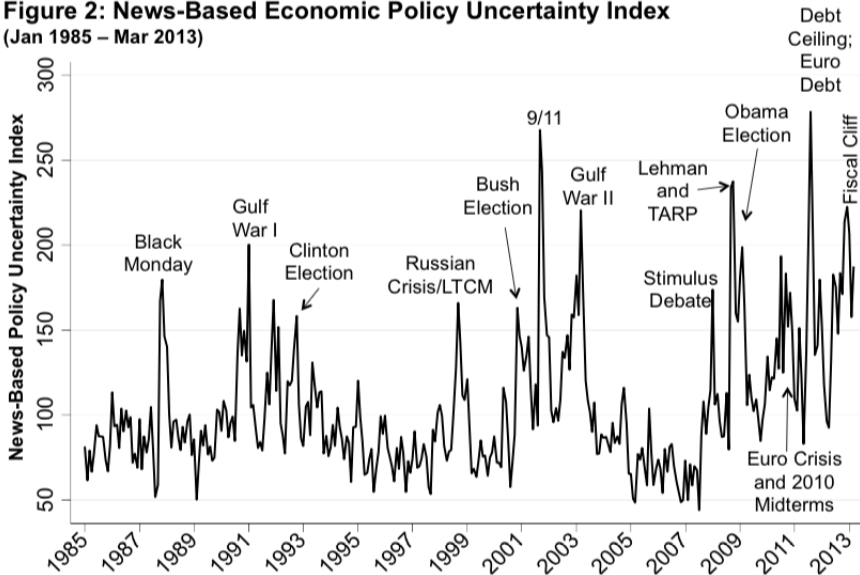
Additional policy term set chosen among larger set of phrases to align with human labels:

3. Article contains “congress” OR “deficit” OR “federal reserve” OR “legislation” OR “regulation” OR “white house”

# BBD Index

**Figure 2: News-Based Economic Policy Uncertainty Index**

(Jan 1985 – Mar 2013)



# Dictionaries and Measurement Error

Confusion matrix from [Baker et al., 2016].

Human Labels	Classification Labels	
	0	1
0	1486	474
1	825	802

# [Shapiro et al., 2022] Continued

News Lexicon	Econ/Finance	50754	0.071	0.180	0.459	0.525
+ Negation rule	News articles		0.075	0.188	0.469	0.525
News Lexicon + LM + HL	Combined	50754	0.078	0.197	0.480	0.524
+ Negation rule			<b>0.082</b>	<b>0.206</b>	<b>0.491</b>	<b>0.525</b>

Notes: GI, LM, and HL refer, respectively, to the following lexicons: Harvard General Inquirer; Loughran and McDonald(2011), updated in 2014; and Hu and Liu (2004). The goodness-of-fit statistics are calculated using the full 800-article sample for which we have human ratings.

## Relationship among Concepts

# Combining Dictionaries

One common strategy to measure how concepts relate to each other is to:

1. Define separate dictionaries for each concept.
2. Count the instances in which terms from each dictionary co-occur in some local window (strictly speaking not BOW model).

# Firm-Level Political Risk

[Hassan et al., 2019] measures firm-level political risk from quarterly earnings calls made by firms traded on US stock markets.

Transcripts from 175,797 conference calls made by 9,478 firms between 2002 and 2016 (downloaded from Thomson Reuters' StreetEvents).

BBD uncertainty measures aggregate risk arising from policymaking, but not firm-specific risks.

Uses a risk/uncertainty dictionary, but the method for associating these to political vs. non-political risks is novel.

Define corpora of canonical political language  $\mathbb{P}$  and non-political language  $\mathbb{N}$ , and compute all bigrams from each.

Sources for these training libraries are undergraduate textbooks or, alternatively, newspaper articles.

# Political Risk Measure

$$PRisk_{it} = \frac{\sum_b^{B_{it}} \left( 1[b \in \mathbb{P} \setminus \mathbb{N}] \times 1[|b - r| < 10] \times \frac{f_{b,\mathbb{P}}}{B_{\mathbb{P}}} \right)}{B_{it}}$$

$B_{it}$  is the total number of bigrams for firm  $i$  at time  $t$ .

$b$  is an individual bigram.

$r$  is the position of the nearest synonym of risk or uncertainty.

$f_{b,\mathbb{P}}$  is the count of bigram  $b$  in the political corpus.

$B_{\mathbb{P}}$  is the total number of bigrams in the political corpus.



# Results

Firms with higher levels of political risk have higher volatility in their stock prices.

Firms with higher political risk engage more in lobbying.

Sector membership and time explain little variation in firm-level risk.

Main conclusion is that location in cross-section of risk exposures seems to matter for firms at least as much as time-series variation.

## Relating Text to Metadata

# Text Regression

Suppose that the text has associated metadata  $\mathbf{y}_d$ , which might contain speaker ID, timestamp, or any other numeric covariate.

Associating text with metadata involves associating  $\mathbf{x}_d$  and  $\mathbf{y}_d$ .

Most straightforward approach would regress  $y_{d,j}$  on  $\mathbf{x}_d$  and  $\mathbf{y}_{d,-j}$ .

Due to strong dependence structure in  $\mathbf{x}_d$ , strong case for use of non-linear models.

# Inverse Regression

Inverse regression models specify a model for  $p(\mathbf{x}_d | y_d)$ .

Well-known example is [Gentzkow and Shapiro, 2010].

Drawing on this paper as motivation, [Taddy, 2013] and [Taddy, 2015] propose fully generative models for inverse regression.

[Gentzkow et al., 2019] uses these models to study political polarization.

# Multinomial Inverse Regression

Model takes the form

$$\mathbf{x}_d \sim \text{MN}(\mathbf{q}_d, N_d) \text{ where } q_{d,v} = \frac{\exp(a_v + \mathbf{y}_d^T \mathbf{b}_v)}{\sum_v \exp(a_v + \mathbf{y}_d^T \mathbf{b}_v)}.$$

Generalized linear model with a (multinomial) logistic link function.

MLE estimates of multinomial regression coefficients can be approximated by estimating  $V$  separate Poisson regression models of  $x_{d,v}$  on  $\mathbf{y}_d$ .

LASSO prior used to regularize regression parameters.

# Application to Congressional Speech

[Gentzkow et al., 2019] use MNIR to model speech data from the *US Congressional Record* from 1873-2016.

Select speeches by Democrats/Republicans (7,732 speakers). Total 36,161 unique speaker-session.

Count two-word phrases (bigrams): 508,351 phrases with count  $\geq 10$  in at least one session.

$\mathbf{y}_d$  includes party, state, chamber, gender.

# Democratic Phrases

## MOST PARTISAN PHRASES FROM THE 2005 CONGRESSIONAL RECORD<sup>a</sup>

### Panel A: Phrases Used More Often by Democrats

#### *Two-Word Phrases*

private accounts  
trade agreement  
American people  
tax breaks  
trade deficit  
oil companies  
credit card  
nuclear option  
war in Iraq  
middle class

Rosa Parks  
President budget  
Republican party  
change the rules  
minimum wage  
budget deficit  
Republican senators  
privatization plan  
wildlife refuge  
card companies

workers rights  
poor people  
Republican leader  
Arctic refuge  
cut funding  
American workers  
living in poverty  
Senate Republicans  
fuel efficiency  
national wildlife

#### *Three-Word Phrases*

veterans health care  
congressional black caucus  
VA health care  
billion in tax cuts  
credit card companies  
security trust fund  
social security trust  
privatize social security  
American free trade  
central American free

corporation for public  
broadcasting  
additional tax cuts  
pay for tax cuts  
tax cuts for people  
oil and gas companies  
prescription drug bill  
caliber sniper rifles  
increase in the minimum wage  
system of checks and balances  
middle class families

cut health care  
civil rights movement  
cuts to child support  
drilling in the Arctic National  
victims of gun violence  
solvency of social security  
Voting Rights Act  
war in Iraq and Afghanistan  
civil rights protections  
credit card debt

## Republican Phrases

TABLE I—Continued

### Panel B: Phrases Used More Often by Republicans

### Two-Word Phrases

stem cell  
natural gas  
death tax  
illegal aliens  
class action  
war on terror  
embryonic stem  
tax relief  
illegal immigration  
date the time

personal accounts  
Saddam Hussein  
pass the bill  
private property  
border security  
President announces  
human life  
Chief Justice  
human embryos  
increase taxes

- retirement accounts
- government spending
- national forest
- minority leader
- urge support
- cell lines
- cord blood
- action lawsuits
- economic growth
- food program

### Three-Word Phrases

- embryonic stem cell
- hate crimes legislation
- adult stem cells
- oil for food program
- personal retirement accounts
- energy and natural resources
- global war on terror
- hate crimes law
- change hearts and minds
- global war on terrorism

Circuit Court of Appeals  
death tax repeal  
housing and urban affairs  
million jobs created  
national flood insurance  
oil for food scandal  
private property rights  
temporary worker program  
class action reform  
Chief Justice Rehnquist

Tongass national forest  
pluripotent stem cells  
Supreme Court of Texas  
Justice Priscilla Owen  
Justice Janice Rogers  
American Bar Association  
growth and job creation  
natural gas natural  
Grand Ole Opry  
reform social security



# Polarization

Let  $q_{t,v}^D(\mathbf{y}')$  be the probability that a Democrat at time  $t$  with observables  $\mathbf{y}'$  speaks phrase  $v$ . Similarly define  $q_{t,v}^R(\mathbf{y}')$ .

Given phrase  $v$ , posterior probability of the speaker being a Democrat is (assuming uniform prior)

$$\rho_{t,v}(\mathbf{y}') = \frac{q_{t,v}^D(\mathbf{y}')}{q_{t,v}^D(\mathbf{y}') + q_{t,v}^R(\mathbf{y}')}$$

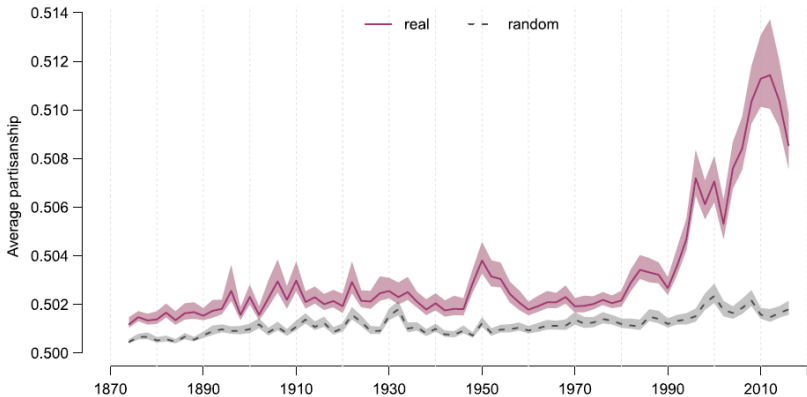
Partisanship is the expected posterior after hearing a single phrase by a speaker with characteristics  $\mathbf{y}'$ :

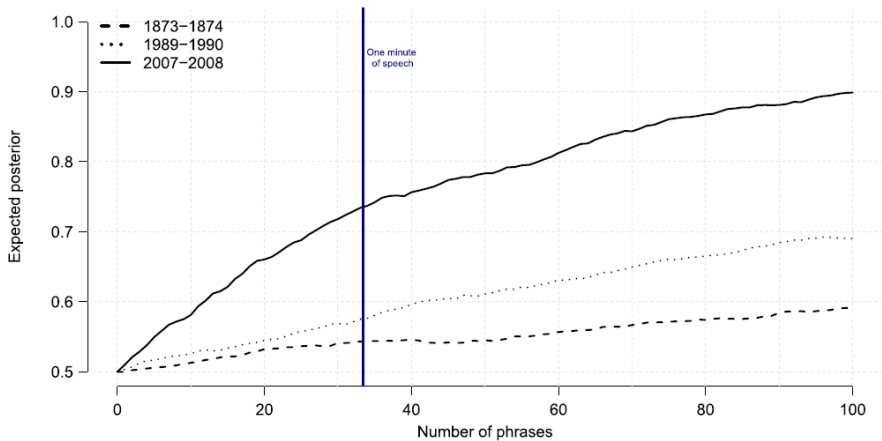
$$\pi_t(\mathbf{y}') = \frac{1}{2} \mathbf{q}_t^D(\mathbf{y}') \cdot \rho_t(\mathbf{y}') + \frac{1}{2} \mathbf{q}_t^R(\mathbf{y}') \cdot (1 - \rho_t(\mathbf{y}'))$$

Let  $s_t$  be total speakers in session  $t$ . Average partisanship is

$$\bar{\pi}_t = \frac{1}{s_t} \sum_{i=1}^{s_t} \pi_{it}(\mathbf{y}'_{it})$$

*Panel B: Partisanship from Preferred Penalized Estimator ( $\hat{\pi}_t^*$ )*





# Sufficient Reduction Projection

There remains the issues of how to use the estimated model for classification.

Let  $z_{d,j} = \mathbf{f}_d^T \hat{\mathbf{b}}_j$  be the *sufficient reduction projection* for the  $j$ th covariate for document  $d$ , where  $\mathbf{f}_d = \mathbf{x}_d / N_d$  is a vector of term frequencies.

$z_{d,j}$  is sufficient for predicting  $y_{d,j}$  in the sense that

$$y_{d,j} \perp \mathbf{x}_d, N_d \mid z_{d,j}, \mathbf{y}_{d,-j}.$$

All the information contained in the high-dimensional frequency counts relevant for predicting  $y_{d,j}$  can be summarized in the SR projection.

Dimensionality reduction targeted at specific covariate.

# Classification

For classification, use the SR projections to build a forward regression that models  $y_{d,j}$  as some function of  $z_{d,j}$ ,  $\mathbf{y}_{d,-j}$ : OLS; logistic; with or without non-linear terms in  $z_{d,j}$ , etc.

To predict  $y_{d',j}$  for an out-of-sample document  $d'$ :

1. Form  $z_{d',j}$  given the estimated  $\hat{\mathbf{b}}_j$  coefficients in the training data.
2. Use the estimated forward regression to generate a predicted value for  $y_{d',j}$ .

# Conclusion

The document-term matrix can be used to address each of the four measurement problems relevant for text-as-data in economics and finance.

Term-count analysis has been, and will continue to be, very influential.

Strength is that matrix-structured data is relatively familiar to economists, and analysis is relatively straightforward.

Nevertheless, all sequential information is ignored and much of natural language's meaning depends on context.

# References I

Ash, E. and Hansen, S. (2023).

Text Algorithms in Economics.

Annual Review of Economics, 15(1):659–688.

Baker, S. R., Bloom, N., and Davis, S. J. (2016).

Measuring Economic Policy Uncertainty.

The Quarterly Journal of Economics, 131(4):1593–1636.

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. (2020).

Language Models are Few-Shot Learners.

Denny, M. J. and Spirling, A. (2018).

Text Preprocessing For Unsupervised Learning: Why It Matters, When It Misleads, And What To Do About It.

Political Analysis, 26(2):168–189.

# References II

Gentzkow, M. and Shapiro, J. M. (2010).

What Drives Media Slant? Evidence From U.S. Daily Newspapers.

Econometrica, 78(1):35–71.

Gentzkow, M., Shapiro, J. M., and Taddy, M. (2019).

Measuring Group Differences in High-Dimensional Choices: Method and Application to Congressional Speech.

Econometrica, 87(4):1307–1340.

Hassan, T. A., Hollander, S., van Lent, L., and Tahoun, A. (2019).

Firm-Level Political Risk: Measurement and Effects.

The Quarterly Journal of Economics, 134(4):2135–2202.

Hoberg, G. and Phillips, G. (2010).

Product Market Synergies and Competition in Mergers and Acquisitions: A Text-Based Analysis.

The Review of Financial Studies, 23(10):3773–3811.



# References III

Hoberg, G. and Phillips, G. (2016).

Text-Based Network Industries and Endogenous Product Differentiation.

Journal of Political Economy, 124(5):1423–1465.

Kelly, B., Papanikolaou, D., Seru, A., and Taddy, M. (2021).

Measuring Technological Innovation over the Long Run.

American Economic Review: Insights, 3(3):303–320.

Loughran, T. and McDonald, B. (2011).

When Is a Liability Not a Liability? Textual Analysis, Dictionaries, and 10-Ks.

The Journal of Finance, 66(1):35–65.

Plaza-del-Arco, F. M., Nozza, D., and Hovy, D. (2024).

Wisdom of Instruction-Tuned Language Model Crowds. Exploring Model Label Variation.

In Abercrombie, G., Basile, V., Bernadi, D., Dudy, S., Frenda, S., Havens, L., and Tonelli, S., editors, Proceedings of the 3rd Workshop on Perspectivist Approaches to NLP (NLPerspectives) @ LREC-COLING 2024, pages 19–30, Torino, Italia. ELRA and ICCL.

# References IV

Shapiro, A. H., Sudhof, M., and Wilson, D. J. (2022).

Measuring news sentiment.

[Journal of Econometrics](#), 228(2):221–243.

Taddy, M. (2013).

Multinomial Inverse Regression for Text Analysis.

[Journal of the American Statistical Association](#), 108(503):755–770.

Taddy, M. (2015).

Distributed Multinomial Regression.

[The Annals of Applied Statistics](#), 9(3):1394–1414.

Tetlock, P. C. (2007).

Giving Content to Investor Sentiment: The Role of Media in the Stock Market.

[The Journal of Finance](#), 62(3):1139–1168.