Machine Learning, Artificial Intelligence, and Natural Language Processing

Marcelo C. Medeiros

Department of Economics The University of Illinois at Urbana-Champaign





Natural Language Processing (NLP)

- NLP is a discipline that intersects computer science, linguistics, and artificial intelligence.
- NLP focuses on how computers comprehend, interpret, and handle human languages.
- It can involve things such as interpreting the semantic meaning of a language, translating between human languages, or recognizing patterns in human languages.
- It uses statistical methods, machine learning, and text mining.



"The ability to take data—to be able to understand it, to process it, to extract value from it, to visualize it, to communicate it—that's going to be a hugely important skill in the next decades [...]. Because now we really do have essentially free and ubiquitous data. So the complimentary scarce factor is the ability to understand that data and extract value from it.

-Hal Varian, The McKinsey Quarterly, January 2009



Motivation What is the dataset of interest?



- Sews from newspapers
- 🥃 Reports, general documents, ...
- 🥃 Social media posts
- 🥃 Central Banks's communications
- 🛢 Books, papers, etc
- ┋....





- Index proposed by Baker, Bloom, and Davies (Quarterly Journal of Economics, 2016)
- A measure of policy-related economic uncertainty
- Quantifies newspaper coverage of policy-related economic uncertainty



Vol. 131 November 2016 Issue 4

MEASURING ECONOMIC POLICY UNCERTAINTY*

SCOTT R. BAKER NICHOLAS BLOOM STEVEN J. DAVIS

We develop a new index of ensume index pinet matrixer (1970) based as meaning well (2000 mempione ariseline and the model of the second secon

⁴⁰ Units Alem Jorring, Keli King, Abdulla Al-Korent, Ispini Biller, Jerr Boltak, Valkari Dakeber, Veng Derry, Biller Bort, Ven Borry, Baber, Ori-Bort, Sanger, S

© The Asthor(s) 2016. Published by Oxford University Press, on helaff of President and Follows of Harvard College. All rights reserved. For Permissions, please email: journals.permission@buy.com The Districts.Aurent of Researcher (2016). 1920. IEEE, doi:10.1003/sightsch04.

The Quarterly Journal of Economics (2016), 1503–1636. doi:10.1003/qa/qw024 Advance Access publication on July 11, 2016.

1593



For the US, the index is mostly based on a normalized index of search results from 10 large newspapers

- * USA Today, the Miami Herald, the Chicago Tribune, the Washington Post, the Los Angeles Times, the Boston Globe, the San Francisco Chronicle, the Dallas Morning News, the New York Times, and the Wall Street Journal
- Extensions to many countries

Used as input to several econometric/economic models



Why should we care about text? Example: Economic Policy Uncertainty (EPU)



Source: "Measuring Economic Policy Uncertainty" by Scott R. Baker, Nicholas Bloom and Steven J. Davis, as updated at www.policyuncertainty.com. Monthly data normalized to 100 prior to 2010.



Source: http://www.policyuncertainty.com and Baker, Bloom, and Davis (QJE, 2016).

Why should we care about text?

Example: Economic Policy Uncertainty (EPU



Notes: Global EPU calculated as the GDP-weighted average of monthly EPU index values for US, Canada, Brazil, Chile, UK, Germany, Italy, Spain, France, Netherlands, Russia, India, China, South Korea, Japan, Ireland, Sweden, and Australia, using GDP data from the IMF's World Economic Outlook Database. National EPU index values are from www.PolicyUncertainty.com and Baker, Bloom and Davis (2016). Each national EPU Index is renormalized to a mean of 100 from 1997 to 2015 before calculating the Global EPU Index.

Source: http://www.policyuncertainty.com and Baker, Bloom, and Davis (QJE, 2016).

Why should we care about text?



Notes: Index reflects scaled monthly counts of articles in Folha de São Paulo containing "incerto" or "incerta", "econômico" or "economia", and one or more policy-relevant terms that include regulação, déficit, orçamento, imposto, "banco centra", jeanaito, congresso, senado, legislação, and tarifa. Normalized to a mean of 100 from 1991 to 2011. Index methods follow "Measuring Economic Policy Uncertainty" by Baker, Bloom and Davis. Data avaitable at <u>www.PolicyUncertaintruc.com</u>.

Source: http://www.policyuncertainty.com and Baker, Bloom, and Davis (QJE, 2016).



Inflation Forecasting

Even structured (real-time) data can provide better forecasts... ... when carefully estimated ML models are considered



County, Marala C. Madedon





Inflation Forecasting

However, news and social-media-based indexes can help further The art of transforming text into numbers



	Focus (available)	Focus (benchmark)	Bias correction (OLS)	Bias correction (w/ indices)	adaLASSO (no indices)	adaLASSO (w/ indices)
inf12m	1.000	0.996	1.121	1.003	0.703	0.614
inf6m	1.000	0.984	0.957	0.753	0.810	0.673
inf3m	1.000	0.960	0.917	0.766	0.913	0.902
infim_30d	1.000	0.914	0.914	0.861	0.871	0.887
infim_5d	1.000	0.878	0.912	0.903	0.917	0.917



APPLIED ECONOMETRICS

Making text count: Economic forecasting using newspaper text

Deri Kalamara, Arthur Turrell 🌰 Chris Redi, George Kapetanico, Sujit Kapadia

Test published: 11 May 2022 | https://doi.org/10.1002/jae.2907 | Chatloric S

II ((10)

Summary

This gave maximes several away in end at large processing spin from endpoint of and detects that a share that the several several response for exact and the several several several several several several several several several from the response of the several several several several several several several from the response of the several several several several several several from the response of the several maximum devices of the several several several several several several several maximum devices of the several several several several several several several maximum devices the several several several several several several several maximum devices the several several several several several several several several several devices that several sev

APPLIED ECONOMETRICS

RESEARCH ANTICLE : @ Open Access : I © ③ ③ ⑤ News media versus FRED-MD for macroeconomic forecasting

jon Ellingsen 🖨 Vegard H. Larsen 🖨 Leif Anders Thorsrud 🖨

First published: 26.July 2021 | https://doi.org/10.1002/jae.2859 | Citations: 9

🗮 FOF 🔨 TOOLS 📢 SHARE

Summary

Using a unique distance of 22.5 million news articles from the dow jones hereaves whitely, we perform an in a specific performance of the performance of the most weight of the the PEDE Modulate. Focusing of the SGD constraints and the specific performance of the most of the most specific performance of the most specific performance of the performance of



Include adapting for the states includes all provides states and s



meaning an tables terms, we employ torout data and machine transmission that the source allow means of increments "distances that the source of the source of the source of the source "distances that the source of the source the source of the source sources. This is which extends the source transmission that the source of the source sources the source of the source of the source the source of the source sources. This is the source of the source of the source the source of the



1

Inflation Forecasting

The age of AI The role of large language models





Beyond Inflation Forecasting Nowcasting



Introduction to Natural Language Processing Let us start with a simple setup

Transforming text into numbers



Computers and ML models are able to process numbers

▶ We must transform text data into numbers



1. Pre-process and represent raw text ${\cal D}$ as a numerical array ${m C}$

* From text to numbers

2. Apply models to map $oldsymbol{C}$ to predicted outcomes $\widehat{oldsymbol{V}}$

3. Use \widehat{V} for descriptive or causal inference



- ▶ Web pages as HTML or XML files
- Noneditable formats as pdf
- Editable formats as .docx.
- Application Programming Interfaces (API)
- Optical Character Recognition for text extraction
- Manual extraction



Textual data are inherently high-dimensional

- ${\ensuremath{\,\ast\,}}$ Collection ${\ensuremath{\mathcal{D}}}$ of D documents, each with N words
- * Each document has the same number of words for simplicity
- $\ast\,$ If each word is drawn from a vocabulary of p-possible words, then the unique representation of these documents has dimension p^N

▶ A simple tweet with N = 280 words can have a p^{280} -dimension

 ${\mbox{ * The length of a basic English speaker is about $10,000$ words}$



- ► How do we deal with such huge-dimensional data?
 - * Computational capacity quickly escalates without any prior cleaning
- Cleaning is fundamental to organize textual data
 - * Assumptions regarding useful or useless words need to be made
 - * The order of the cleaning process may differ through different applications



Introduction

Pre-processing Textual Data

Topic Models

Latent Dirichlet Allocation (LDA)





- Divide the raw corpus of text into individual documents
- Level of aggregation depends on the application
 - 1. If the objective is to nowcast the GDP, and documents are news, we should aggregate by a reference period, e.g. day, week, month
 - 2. If we are interested in Central Bank Communication, it may be interesting to treat each communication isolated
- Aggregation level may not be straightforward
 - * If we analyze news from several newspapers, do we aggregate all articles by day or treat them individually?



Tokenization: Breaking text into words or phrases

Lowercasing, removing punctuation

Stopword removal: Remove common words (e.g., "the", "and")

► Stemming/Lemmatization: Reduce words to base form (e.g., "running" → "run")

Feature selection: Keep only informative



Original text:

"Good night, good night! Parting is such sweet sorrow."

After Preprocessing:

Remove punctuation, stop words

Apply stemming

Result: [good, night, good, night, part, sweet, sorrow]



- After tokenizing, cleaning and extracting the linguistic roots, we are able to structure the dara
 - * The usual representation of text is made in terms of **counts** of sentences, e.g. **bag-of-words** or **dictionaries**
 - * Any representation of text will throw away information
 - * The goal is to keep the relevant part for the application

Keep in mind: without simplifications, the problem becomes unfeasible



- The decision whether to aggregate or not articles may depend on the application:
 - * If we are interested in forecasting using news-analysis, it may be interesting to aggregate articles in daily-terms.
 - * If we are interested in Central Bank communication, it may be interesting to treat each document individually.
- Suppose we are interested in forecasting: If we have a set of different newspapers, should we treat them individually, or aggregate the articles as a single representative newspaper?



The document-term matrix - cross-section of documents

After pre-processing, we are left with a finite list of (potentially non-unique) terms

- * Index each unique term in the corpus by $n \in \{1, \cdots, N\}$ where N is the total number of unique terms
- * For each document $d \in \{1, \cdots, D\}$, compute the **count** $x_{d,n}$ of occurrences of term n in document d
- * The $D \times N$ matrix $\boldsymbol{X}_{D,N}$ of all such counts is called the document-term matrix
- Each article treated separately



- Consider the following extension:
 - * Index each unique term of the corpus (overall articles) by $n \in \{1, \dots, N\}$ where N is the number of unique terms
 - * For each article $d\in\{1,\cdots,D\}$ compute the count $x_{d,n}$ as the number of occurrences of term n in document D
 - * Aggregate the number of occurrences over all documents $x_{d,n}$ for a given period $t\in(1,\cdots,T)$
 - * Generate period-counts $x_{n,t}$
 - * The $T\times N$ matrix $\boldsymbol{X}_{T,N}$ of all such counts is called the document-term matrix
- This formulation aggregates articles by a period of reference (e.g., days, weeks, months, ...)

Bag-of-words: the order of words is ignored altogether, and a phrase of length n is called a n-gram:

 \blacktriangleright Higher order *n*-grams encode more complex structures

* Usually, we consider at most 3-grams





- Dictionary: focus on variation across observations along a limited number of words
- Define a dictionary of terms $\mathcal{H} = \{1, \cdots, H\}$:
 - * a) Boolean: The dictionary-based boolean $b_{d,n}$ of term n in document d is given by: $b_{d,n} = \mathbb{I}(x_{d,n} \in \mathcal{H})$, where \mathbb{I} is the indicator function
 - * b) Counts: The dictionary-based count $\bar{x}_{d,n}$ of term n in document d is given by: $\bar{x}_{d,n} = \mathbb{I}(x_{d,n} \in \mathcal{H})x_{d,n} = b_{d,n}x_{d,n}$, where \mathbb{I} is an indicator function.
- Dictionaries are arbitrary

Introduction

Pre-processing Textual Data

Topic Models

Latent Dirichlet Allocation (LDA)





Topic models are statistical models for discovering abstract topics in a collection of documents

They help uncover hidden thematic structures in large text corpora

Useful for document classification, recommendation systems, and information retrieval



Colection of documents



Documents may consist of one or more topics

▶ We want the computer to discover the topics of each document



Introduction

Pre-processing Textual Data

Topic Models

Latent Dirichlet Allocation (LDA)





Most popular topic modeling technique

Assumes each document is a mixture of topics, and each topic is a mixture of words

Probabilistic generative model



Introduction to Latent Dirichlet Allocation (LDA)

A simple machine to create text



https://www.youtube.com/watch?v=9mNV4AwA9QI



Introduction to Latent Dirichlet Allocation (LDA)

A simple machine to create text



https://www.youtube.com/watch?v=9mNV4AwA9QI





35

Introduction to Latent Dirichlet Allocation (LDA)

A simple machine to create text



https://www.youtube.com/watch?v=9mNV4AwA9QI


A simple machine to create text – Version 1





A simple machine to create text – Version 1





A simple machine to create text – Version 1





A simple machine to create text – Version 1





A simple machine to create text - Version 1





A simple machine to create text – Version 2

	Topics	pear	bird	planet	burger	sun	dog	chair	notebook	phone	table
	Τ1	0.07	0.30	0.03	0.01	0.05	0.40	0.04	0.05	0.04	0.01
	T2	0.25	0.10	0.01	0.31	0.01	0.15	0.05	0.01	0.02	0.10
<u>ш</u>	тз	0.05	0.04	0.35	0.01	0.25	0.02	0.05	0.10	0.03	0.10
					Text G	eneratio	n Machin	e			
						W _n	N	, ,			



A simple machine to create text – Version 2





A simple machine to create text – Version 2





A simple machine to create text – Version 2





A simple machine to create text – Version 2





A simple machine to create text – Version 3





A simple machine to create text – Version 3





A simple machine to create text – Version 3

	Topics	pear	bird	planet	burger	sun	dog	chair	notebook	phone	tabl
	Τ1	0.07	0.30	0.03	0.01	0.05	0.40	0.04	0.05	0.04	0.01
	T2	0.25	0.10	0.01	0.31	0.01	0.15	0.05	0.01	0.02	0.1
ш	тз	0.05	0.04	0.35	0.01	0.25	0.02	0.05	0.10	0.03	0.1
For each word, we w	ill sample	e a topic	from θ		Text G	eneration	n Machin	e			
			Ud								





A simple machine to create text – Version 3





A simple machine to create text – Version 3





A simple machine to create text – Version 3





A simple machine to create text – Version 3





A simple machine to create text – Version 3







https://www.youtube.com/watch?v=9mNV4AwA9QI



Generating multinomial distributions



https://www.youtube.com/watch?v=9mNV4AwA9QI



56

Generating multinomial distributions



https://www.youtube.com/watch?v=9mNV4AwA9QI



57



https://www.youtube.com/watch?v=9mNV4AwA9QI





https://www.youtube.com/watch?v=9mNV4AwA9QI



Generating multinomial distributions







https://www.youtube.com/watch?v=9mNV4AwA9QI



Generating multinomial distributions



https://www.youtube.com/watch?v=9mNV4AwA9QI



62

Generating multinomial distributions



https://www.youtube.com/watch?v=9mNV4AwA9QI



63







Generating multinomial distributions







65



https://www.youtube.com/watch?v=9mNV4AwA9QI





https://www.youtube.com/watch?v=9mNV4AwA9QI



Generating multinomial distributions







https://www.youtube.com/watch?v=9mNV4AwA9QI





https://www.youtube.com/watch?v=9mNV4AwA9QI





https://www.youtube.com/watch?v=9mNV4AwA9QI





https://www.youtube.com/watch?v=9mNV4AwA9QI


Generating multinomial distributions

This is the Dirichlet distribution





The Dirichlet distribution







A simple machine to create text – Version 4





A simple machine to create text – Version 4





Another Dirichlet distribution



https://www.youtube.com/watch?v=9mNV4AwA9QI



Another Dirichlet distribution







Another Dirichlet distribution





A simple machine to create text – Version 5





A simple machine to create text – Version 5





A simple machine to create text – Version 5







A simple machine to create text – Version 5

For every document d = 1,..., 0, draw a multinomial distribution θ_d from the Dirichlet distribution with parameter α . For every topic k = 1,..., K, draw a multinomial distribution ϕ_k from the Dirichlet distribution with parameter β . For every word position n = 1,..., N in the document d, assign a topic Z, from θ_d and draw a work from ϕ_k with k = Z,





► Algorithm:

- 1. Draw β_k independently for $k = \{1, \dots, K\}$ from $\text{Dir}(\boldsymbol{\nu})$.
- 2. Draw θ_d independently for $d = \{1, \dots, D\}$ from $Dir(\alpha)$.
- 3. Each word $w_{d,n}$ in document d is generated from a two-step process:
 - 3.1 Draw a topic assignment $z_{d,n}$ from θ_d .
 - 3.2 Draw $w_{d,n}$ from $\beta_{\boldsymbol{z}_{d,n}}$

Estimate the hyperparameters u and α from the Dirichlet distributions.



Scrape online news from major newspapers

* News from January, 2009 up to December, 2021 from Folha de São Paulo, Estadão, and Valor Econômico

Daily collected and weekly aggregated.





- Bag-of-words: for each individual article count the number of occurrences of a given term (or a composition of terms) on the cleaned text.
 - * Generate a matrix of all such counts
 - * Approach: 1-grams up to 3-grams counting
 - + e.g. 'Covid-19', 'Presidente.Bolsonaro' and 'Supremo.Tribunal.Federal'
 - * Aggregate the term-counts to generate the Corpus.
- Generate the Corpus: a matrix X_{T,V} where rows denote months and columns denotes terms

	2009	2010	2011	2012	2013	2014	2015
Folha	598	5329	6012	5235	15476	30628	24048
Estadão	55298	94182	90182	92639	61469	54209	31085
Valor	-	-	-	65545	51057	45509	45048
Sum	55896	99511	96194	163419	128002	130346	100181

	2016	2017	2018	2019	2020	2021	Total
Folha	22073	18479	19853	16375	15278	13988	194944
Estadão	35758	40258	34957	32405	36337	31459	692866
Valor	47010	46156	47001	46059	60773	64969	525309
Sum	104841	104893	101811	94839	112388	110416	1413119



Textual Data

grafton, wisconsin - guinhentas doses da vacina contra a covid-19 tiveram que ser descartadas porque não foram devidamente refrigeradas em wisconsin . nos es tados unidos, de acordo com a rede de hospitais aurora medical center, os imu nizantes foram aparentemente estragadas de forma deliberada por um funcionário , no sábado , o hospital havia revelado que as doses foram acidentalmente deixa das em temperatura ambiente durante a noite por um funcionário da unidade de gr afton , no entanto , nesta guarta-feira , 30 , o aurora divulgou uma nota afirm ando que o funcionário envolvido `` reconheceu que as vacinas foram retiradas i ntencionalmente da geladeira ", o comunicado diz ainda que o funcionário foi demitido e o assunto foi entregue às autoridades para uma investigação mais apr ofundada . o depoimento não menciona o possível motivo dessa ação , e os execut ivos do sistema de saúde não responderam no momento às mensagens que lhes foram enviadas em busca de mais informações . " continuamos a acreditar que a vacinac ão é a nossa saída para a pandemia , estamos mais do que frustrados com o fato de o comportamento desse indivíduo atrasar a vacinação de mais de 500 pessoas " . disse a nota . o aurora medical center se recusou a fornecer informações adic ionais, mas disse que daria mais detalhes na quinta-feira./ap

Cleaning and pre-processing of textual data from a particular article from Estadão. We mark the punctuation extraction in blue together with the number and single-letter removal. We highlight the stopwords, rare words, and temporal markers removed in red. In green, we point out the words that should be lemma-reduced. In black, we highlight words that should not be ex-ante modified a-ante modified areante modified.



Textual Data





Obrigado

