

Artificial Intelligence for Business Research @Antai

## Traditional NLP

Renyu (Philip) Zhang

1

## Agenda

- Natural Language Processing Framework
- N-Gram and Naïve Bayes
- Traditional NLP Applications in Business/Econ Research

2

2

## Classic NLP Framework

- Reference: <https://www.coursera.org/specializations/natural-language-processing>  
<https://web.stanford.edu/~jurafsky/slp3/>
- A classic NLP framework usually contains 2 parts:
  - Pre-processing: Text → Numeric representations
  - Classification: Numeric representations → outcome
    - Sentiment Classifier: Text → Sentiment Score
    - Review Classification: Text → Review Problem
    - Machine Translation: Text → Other language
- Typical NLP:
  - Sentiment Classification
  - Machine Translation
  - Document Similarity
  - Topic Modelling
  - Etc.

I am happy because I am learning NLP

```

graph LR
    A["I am happy  
because I am  
learning NLP"] --> B[X]
    B --> C[Train  
LR]
    C --> D[Classify]
    D --> E["Positive: 1"]
  
```

3

## NLP History

Probability & Vector Representations of NLP

Before 2010

Word2vec & Seq2seq

2013-2017

Transformer-based Models

2018-2022

Large Language Models (LLM)

2023-now

Efficient estimation of word representations in vector space

T Mikolov, K Chen, G Corrado, J Dean - arXiv preprint arXiv:1301.3781, 2013 - arxiv.org

... vector  $X = \text{vector}(\text{"biggest"}) - \text{vector}(\text{"big"}) + \text{vector}(\text{"small"})$ . Then, we search in the vector space for the word ... question (we discard the input question words during this search). When the ...

☆ Save 99 Cite Cited by 27227 Related articles All 48 versions 86

[PDF] Glove: Global vectors for word representation

J Pennington, R Socher, - Proceedings of the 2014 ... 2014 - aclanthology.org

... We use our insights to construct a new model for word representation which we call GloVe, for Global Vectors, because the global corpus statistics are captured directly by the model. ...

☆ Save 99 Cite Cited by 26265 Related articles All 34 versions 86

[PDF] Language models are unsupervised multitask learners

A Radford, J Wu, R Chen, D Luan, D Amodei, - OpenAI blog, 2019 - persagen.com

Natural language processing tasks, such as question answering, machine translation, reading comprehension, and generation, have traditionally been approached when supervised learning on task-specific datasets. We demonstrate that language models begin to learn these tasks without any explicit supervision when trained on a new dataset of millions of webpages called WebText. When conditioned on a document plus questions, the answers generated by the language model reach 55 F1 on the CoQA dataset-matching or exceeding

☆ Save 99 Cite Cited by 2834 Related articles All 16 versions 86

[PDF] Training language models to follow instructions with human feedback

LQuyang, J Wu, X Jiang, D Almeida, - Advances in ... 2022 - proceedings.neurips.cc

Making language models bigger does not inherently make them better at following a user's intent. For example, large language models can generate outputs that are untruthful, toxic, or simply not helpful to the user. In other words, these models are not aligned with their users. In this paper, we show an avenue for aligning language models with user intent on a wide range of tasks by fine-tuning with human feedback. Starting with a set of labeler-written prompts and prompts submitted through a language model API, we collect a dataset of ...

☆ Save 99 Cite Cited by 4049 Related articles All 10 versions 86

4

## NLP Roadmap

- A Typical NLP Task:



- The fundamental ways of thinking about **words and pre-processing** are the same across eras.

- **Old Architecture:**

- **Explicit probabilistic model of linguistics (and possibly utilities).** This **transparency** in modeling makes these kinds of model useful in **economics-driven** research with ML.

- **New DL Architecture:**

- Better performance + enables new tasks such as machine language generation, chatbot, etc.
- The framework is particularly useful **when you use your own methods** (to do tasks that are not accomplished in CS before).

5

## Text as Data in Biz/Econ Research

*Journal of Economic Literature* 2019, 57(3), 535–574  
<https://doi.org/10.1257/jel.20181020>

### Text as Data<sup>†</sup>

MATTHEW GENTZKOW, BRYAN KELLY, AND MATT TADDY\*

An ever-increasing share of human interaction, communication, and culture is recorded as digital text. We provide an introduction to the use of text as an input to economic research. We discuss the features that make text different from other forms of data, offer a practical overview of relevant statistical methods, and survey a variety of applications. (JEL C38, C55, L82, Z13)

#### 0. Pre-processing:

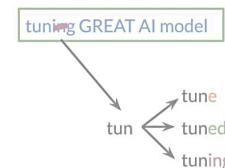
1. Represent raw text  $\mathcal{D}$  as a numerical array  $\mathbf{C}$ ;
2. Map  $\mathbf{C}$  to predicted values  $\hat{\mathbf{V}}$  of unknown outcomes  $\mathbf{V}$ ; and
3. Use  $\hat{\mathbf{V}}$  in subsequent descriptive or causal analysis.

6

6

## Pre-processing

- References: <https://nlp.stanford.edu/IR-book/pdf/02voc.pdf>  
[https://web.stanford.edu/~jurafsky/slp3/slides/2\\_TextProc\\_Mar\\_25\\_2021.pdf](https://web.stanford.edu/~jurafsky/slp3/slides/2_TextProc_Mar_25_2021.pdf)
- Text normalization: Transforming sentences into words.
- Text normalization includes 2 tasks: Word segmentation (i.e., tokenization) and word normalization:
  - Elimination of non-words: URL, HTML, handles, punctuations etc.
  - Tokenization: Parse strings into words.
  - Stop-word removal: Get rid of stop-words which are extremely common, such as "a, an, is, the, of..."
  - Stemming: Convert every word to its stem.
  - Normalization: Normalize accents and diacritics; change all letters into lower-cases.



7

7

## Word Representation

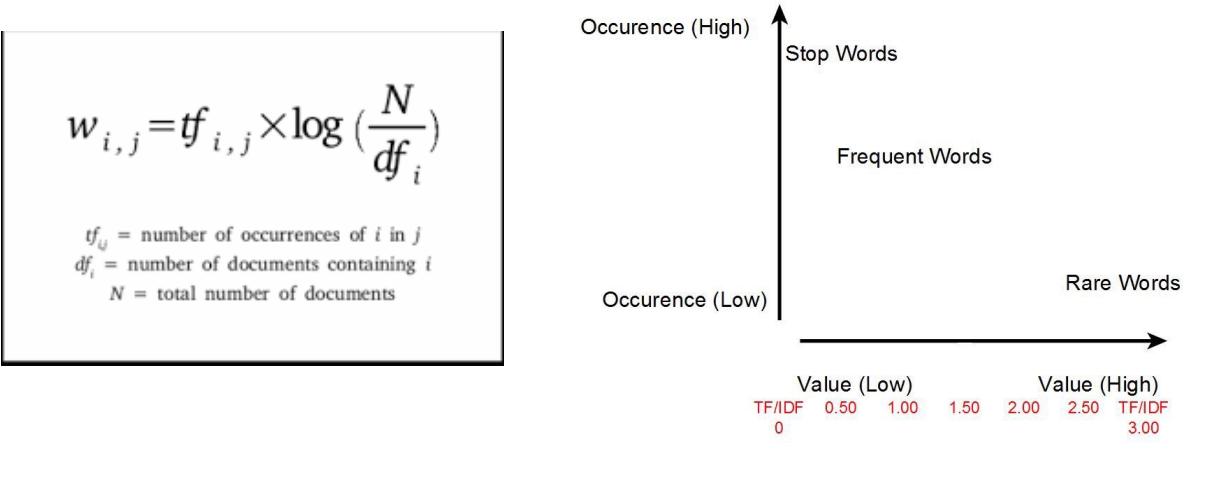
- **Traditional NLP**: Representation based on vocabulary of words and word count.
- **Frequentist** view: Represent words as vectors, which are low-dimensional projection of one-hot encoding of the words depending on its neighbors.
- **Bayesian** view: Represent words as probabilities; each word has a prior to be used and each sentence then has a conditional probability of words.

8

8

## Term Frequency-Inverse Document Frequency

- Each word has different importance for a document/sentence.
- TF-IDF: A word appearing in **fewer documents** and appearing **more times** may be more important.

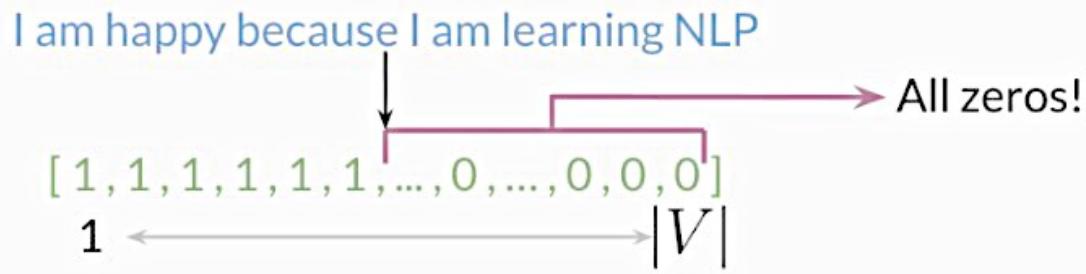


9

9

## One-hot Encoding

- One-hot encoding: **Sparse representation**; think about the dummy variable in econometrics.
- You need  $k$  variables to represent a document if you have vocabulary length equal to  $k$ .



10

10

## Low-Dimensional Document Representation of Words

- Reference: <https://web.stanford.edu/~jurafsky/slp3/>
- Use a word's occurrence in documents to represent a word's meaning.
- Basic idea: Similar words have similar vectors because they tend to occur in similar documents.

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

**Figure 6.5** The term-document matrix for four words in four Shakespeare plays. The red boxes show that each word is represented as a row vector of length four.

11

11

## Sentence/Document Representation

- You can also use word occurrence to represent sentence and document.

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

**Figure 6.3** The term-document matrix for four words in four Shakespeare plays. The red boxes show that each document is represented as a column vector of length four.

- This is called term-document matrix, allowing us to find similar documents.

12

12

## Downstream Task: Sentiment Classification

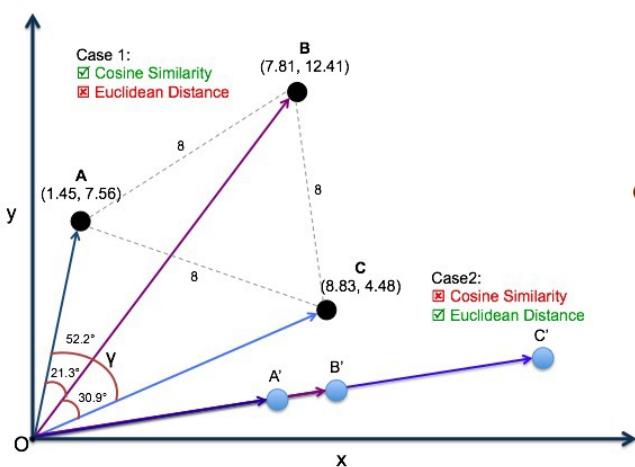
- Now your data is pre-processed from: Sentence<sub>i</sub> → Label<sub>i</sub> = positive or negative
- To: (X<sub>1\_i</sub>, X<sub>2\_i</sub>, X<sub>3\_i</sub>) → y<sub>i</sub> = 1 or 0, where X is the word/sentence/document-representation vector.
- We can use different ML methods to find the mapping between X and y. For example,
  - Linear probability model;
  - Logistic regression or other generalized linear models;
  - Deep learning.

13

13

## Downstream Task: Semantic Similarities

- Cosine Similarity: distance(OA, OC) > distance (OA, OB)
- Euclidean Similarity: distance(OA, OC) < distance(OA, OB)

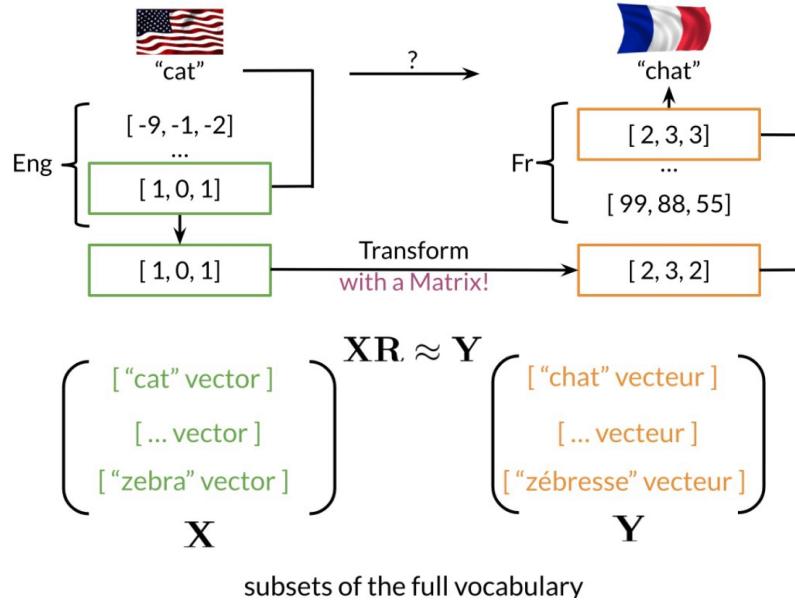


$$\cosine(\mathbf{v}, \mathbf{w}) = \frac{\mathbf{v} \cdot \mathbf{w}}{|\mathbf{v}| |\mathbf{w}|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

14

14

## Downstream Task: Machine Translation



15

15

## Putting Things Together

1. Pre-processing: Take a dataset and do text normalization.
2. Word-representation: Use your favorite way to do word representation.
3. Sentence/Document-representation: Use your favorite way to do sentence/document representation.
4. Downstream-task: Use your favorite model to do downstream tasks.

16

16

## Agenda

- Natural Language Processing Framework
- N-Gram and Naïve Bayes
- Traditional NLP Applications in Business/Econ Research

17

17

## Bayesian Perspective: N-Gram

- Reference: <https://web.stanford.edu/~jurafsky/slp3/>, Chapter 2.
- Language Model: To find the probability of a word given the entire history of the sentence so far.

$$\Pr(W_n | W_1, W_2, W_3, \dots, W_{n-1})$$

- Suppose the sentence is "its water is so transparent that the..."

$$P(\text{the} | \text{its water is so transparent that}) = \\ \frac{\text{Count}(\text{its water is so transparent that the})}{\text{Count}(\text{its water is so transparent that})}$$

- Too complex model + curse of dimensionality.
  - The number of combinations of word history grow exponentially with text length, requiring prohibitively large datasets to compute.

18

18

## Bayesian Perspective: N-Gram

- Chain rule of probability:

$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i | w_1 w_2 \dots w_{i-1})$$

$$\begin{aligned} P(\text{"its water is so transparent"}) &= \\ P(\text{its}) \times P(\text{water} | \text{its}) \times P(\text{is} | \text{its water}) \\ &\quad \times P(\text{so} | \text{its water is}) \times P(\text{transparent} | \text{its water is so}) \end{aligned}$$

19

19

## Bayesian Perspective: N-gram

- N-gram model: Limiting the dependencies on history, a.k.a. **Markov Chain**.

- Unigram:

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i)$$

- Bi-gram:

$$P(w_i | w_1 w_2 \dots w_{i-1}) \approx P(w_i | w_{i-1})$$

- We can extend to trigrams, 4-grams, 5-grams, etc., but there is a trade-off:

- A larger N: (a) longer-distance dependency; (b) much more data to estimate the conditional probabilities.

- Rule-of-thumb: N is no more than 5.

20

20

## Estimating Bigram

Maximum Likelihood Estimation (MLE)   $P(w_i | w_{i-1}) = \frac{count(w_{i-1}, w_i)}{count(w_{i-1})}$

An Example 

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

`<s> I am Sam </s>`  
`<s> Sam I am </s>`  
`<s> I do not like green eggs and ham </s>`

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

$$P(I | <s>) = \frac{2}{3} = .67 \quad P(Sam | <s>) = \frac{1}{3} = .33 \quad P(am | I) = \frac{2}{3} = .67$$

$$P(</s> | Sam) = \frac{1}{2} = 0.5 \quad P(Sam | am) = \frac{1}{2} = .5 \quad P(do | I) = \frac{1}{3} = .33$$

21

21

## Laplace/Add-One Smoothing

Pretend that we saw each word one more time than we did.

$$P_{MLE}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

 See  $(w_{i-1}, w_i)$  one more time.

$$P_{Add-1}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + V}$$

 See each of  $(w_{i-1}, w_k)$  one more time for all  $w_k$  in the vocabulary.

22

22

## More about N-gram

- If we have 10,000 unique words, we will have  $10,000 \times 10,000 = 10e8$  possible combinations of bigram.
- How about we use tri-gram or even 4-gram?
- Some toolkits for n-gram models:
  - SRILM: <http://www.speech.sri.com/projects/srilm/>
  - KenLM: <https://kheafield.com/code/kenlm/>
- Google n-gram viewer: <https://books.google.com/ngrams/>
  - Dataset: <https://storage.googleapis.com/books/ngrams/books/datasetsv2.html>

What drives media slant? Evidence from US daily newspapers  
M Gentzkow, JM Shapiro  
Econometrica, 2010 · Wiley Online Library

We construct a new index of media slant that measures the similarity of a news outlet's language to that of a congressional Republican or Democrat. We estimate a model of newspaper demand that incorporates slant explicitly, estimate the slant that would be chosen if newspapers independently maximized their own profits, and compare these profit-maximizing points with firms' actual choices. We find that readers have an economically significant preference for like-minded news. Firms respond strongly to

SHOW MORE ▾

[☆ Save](#) [99 Cite](#) [Cited by 2330](#) [Related articles](#) [All 33 versions](#) [Web of Science: 710](#) [⊗](#)

Measuring group differences in high-dimensional choices: method and application to congressional speech  
M Gentzkow, JM Shapiro, M Taddy  
Econometrica, 2019 · Wiley Online Library

We study the problem of measuring group differences in choices when the dimensionality of the choice set is large. We show that standard approaches suffer from a severe finite-sample bias, and we propose an estimator that applies recent advances in machine learning to address this bias. We apply this method to measure trends in the partisanship of congressional speech from 1873 to 2016, defining partisanship to be the ease with which an observer could infer a congressperson's party from a single utterance. Our

SHOW MORE ▾

[☆ Save](#) [99 Cite](#) [Cited by 372](#) [Related articles](#) [All 14 versions](#) [Web of Science: 95](#) [⊗](#)

23

23

## Bayesian Perspective: Naïve Bayes

- Reference: <https://web.stanford.edu/~jurafsky/slp3/>, Chapter 3.
- Text classification is a fundamental application in NLP.
  - Sentiment analysis
  - Spam detection
  - Authorship identification
  - Language identification
  - Assigning subject categories, topics, or genres
  - Many more.....
- Input:
  - A document  $d$
  - A fixed set of classes  $C = \{c_1, c_2, \dots, c_J\}$
- Output:
  - A predicted class of document  $d$  in  $C$ .

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c | d)$$

MAP is "maximum a posteriori" = most likely class

$$= \operatorname{argmax}_{c \in C} \frac{P(d | c)P(c)}{P(d)}$$

Bayes Rule

$$= \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

Dropping the denominator

24

## Bayesian Perspective: Naïve Bayes

- Reference: <https://web.stanford.edu/~jurafsky/slp3/>, Chapter 3.

$$\begin{aligned}
 c_{MAP} &= \operatorname{argmax}_{c \in C} P(d | c)P(c) \\
 &= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c)
 \end{aligned}$$

Document d represented as features x<sub>1..xn</sub>

"Likelihood"

O(|X|<sup>n</sup> • |C|) parameters

Could only be estimated if a very, very large number of training examples was available.

"Prior"

How often does this class occur?

We can just count the relative frequencies in a corpus

25

25

## Bayesian Perspective: Naïve Bayes

- Reference: <https://web.stanford.edu/~jurafsky/slp3/>, Chapter 3.
- Bag of words assumption:** Position does not matter.
- Conditional independence:** The feature probabilities  $\Pr(x_i | c)$  are independent conditioned on class  $c$ .

$$P(x_1, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \cdot P(x_3 | c) \cdot \dots \cdot P(x_n | c)$$

positions  $\leftarrow$  all word positions in test document

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c)$$

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$

26

26

## Naïve Bayes Estimation

- Reference: <https://web.stanford.edu/~jurafsky/slp3/>, Chapter 3.
- Maximum likelihood estimates:

$$\hat{P}(c_j) = \frac{N_{c_j}}{N_{total}}$$

The fraction of documents in class  $c_j$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

The fraction of times word  $w_i$  appears among all words in documents of class  $c_j$

$$c_{MAP} = \operatorname{argmax}_c \hat{P}(c) \prod_i \hat{P}(x_i | c)$$

Most likely class of the document

27

27

## Naïve Bayes Example

- Reference: <https://web.stanford.edu/~jurafsky/slp3/>, Chapter 3.

Positive tweets		
I am happy because I am learning NLP		
I am happy, not sad.		
Negative tweets		
I am sad, I am not learning NLP		
I am sad, not happy		

word	Pos	Neg
I	3	3
am	3	3
happy	2	1
because	1	0
learning	1	1
NLP	1	1
sad	1	2
not	1	2
$N_{\text{class}}$	13	12

word	Pos	Neg
I	0.24	0.25
am	0.24	0.25
happy	0.15	0.08
because	0.08	0
learning	0.08	0.08
NLP	0.08	0.08
sad	0.08	0.17
not	0.08	0.17

Let's classify the following tweets as **positive or negative**:

1. I am not sad.
2. I am learning NLP.

28

28

## Naïve Bayes (Laplace) Smoothing

- Reference: <https://web.stanford.edu/~jurafsky/slp3/>, Chapter 3.
- Like n-gram, normalize the estimated probabilities and bound them away from 0.

$$\hat{P}(w_i | c) = \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} (\text{count}(w, c) + 1)}$$

See the word  $w_i$  in class  $c$   
one more time.

$$= \frac{\text{count}(w_i, c) + 1}{\left( \sum_{w \in V} \text{count}(w, c) \right) + |V|}$$

See each word  $w$  of the  
vocabulary in class  $c$  one  
more time.

29

29

## Putting Things Together for Naïve Bayes

- Labeling:** Annotate a dataset with class labels (positive, negative, etc.).
- Pre-processing:** Pre-process the text to words.
- Frequency computation:** Compute the Freq(word, class).
- Probability computation:** Compute  $P(\text{word}|\text{class})$  and  $P(\text{document}|\text{class})$ .
- Classification.**

30

30

## Agenda

- Natural Language Processing Framework
- N-Gram and Naïve Bayes
- Traditional NLP Applications in Business/Econ Research

31

31

## Authorship Identification

JOURNAL OF THE AMERICAN  
STATISTICAL ASSOCIATION

Number 302

JUNE, 1963

Volume 58

INFERENCE IN AN AUTHORSHIP PROBLEM<sup>1-4</sup>  
A comparative study of discrimination methods applied  
to the authorship of the disputed *Federalist* papers

FREDERICK MOSTELLER  
*Harvard University*  
and  
*Center for Advanced Study in the Behavioral Sciences*  
AND  
DAVID L. WALLACE  
*University of Chicago*

This study has four purposes: to provide a comparison of discrimination methods; to explore the problems presented by techniques based strongly on Bayes' theorem when they are used in a data analysis of large scale; to solve the authorship question of *The Federalist* papers; and to demonstrate the value of solving an authorship problem.

With counts are the variables used in discrimination. Since the topic written about heavily influences the rate with which a word is used, care in selection of words is necessary. The filler words of the language such as *an*, *of*, and *upon*, and, more generally, articles, prepositions, and conjunctions provide fairly stable rates, whereas more meaningful words like *war*, *executive*, and *legislature* do not.

After an investigation of the distribution of these counts, the authors execute an analysis employing the usual discriminant function and an analysis based on Bayesian methods. The conclusions about the authorship problem are that Madison rather than Hamilton wrote all 12 of the disputed papers.

The findings about methods are presented in the closing section on conclusions.

This report, summarizing and abbreviating a forthcoming monograph [8], gives some of the results but very little of their empirical and theoretical foundation. It treats two of the four main studies presented in the monograph, and none of the side studies.

- Who wrote the *Federalist Papers*, Alexander Hamilton or James Madison?
- Applying **Naïve Bayes**, where the class c is either Hamilton or Madison, Mosteller and Wallace (1963) find overwhelming evidence that the disputed papers were authored by Madison.
- A similar method was applied to identify who invented instrumental variables estimator.

**Inference in an Authorship Problem: A Comparative Study of Discrimination Methods Applied to the Authorship of the Disputed *Federalist* Papers**  
F. Mosteller, DL Wallace - Journal of the American Statistical ..., 1963 - Taylor & Francis

... Problems of discrimination are widespread, and we wished a case study that would give us ... problem in authorship as a case study. (The classical approach to discrimination problems, ...

☆ Save 99 Cite Cited by 779 Related articles All 4 versions Web of Science: 253 ▷

[PDF] **Retrospectives: Who invented instrumental variable regression?**

JH Stock, F Trebbi - Journal of Economic Perspectives, 2003 - pubs.aeaweb.org

... derivations of the **instrumental variables** estimators of the ... B was showing that **instrumental variables** regression can be ... do, which makes **instrumental variables** regression a central ...

☆ Save 99 Cite Cited by 242 Related articles All 13 versions Web of Science: 82 ▷

32

32

16

# Sentiment and Stock Price

THE JOURNAL OF FINANCE • VOL. LXII, NO. 3 • JUNE 2007

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

PAUL C. TETLOCK\*

### ABSTRACT

I quantitatively measure the interactions between the media and the stock market using daily content from a popular *Wall Street Journal* column. I find that high media pessimism predicts downward pressure on market prices followed by a reversion to fundamentals, and unusually high or low pessimism predicts high market trading volume. These and similar results are consistent with theoretical models of noise and liquidity traders, and are inconsistent with theories of media content as a proxy for new information about fundamental asset values, as a proxy for market volatility, or as a sideshow with no relationship to asset markets.

- Based on the dictionary approach, convert the word counts in WSJ's "Abstract of the Market" into sentiment scores and condense them into a single principal component, called "pessimism factor".
- This pessimism score is used to forecast stock market activity.

**When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks**  
[T Loughran, B McDonald - The Journal of finance, 2011 - Wiley Online Library](#)

... this paper is whether a ... is inherently imprecise, we provide evidence based on 50,115 firm-year 10-Ks between 1994 and 2008 that the H4N list substantially misclassifies words when ...  
☆ Save 99 Cite Cited by 5231 Related articles All 8 versions Web of Science: 1875 »»

## Giving content to investor sentiment: The role of media in the stock market

[PC Tetlock - The Journal of finance, 2007 - Wiley Online Library](#)

... investor sentiment or ... investor sentiment, resulting in downward pressure on prices. It is unclear whether media pessimism forecasts investor sentiment or reflects past investor sentiment...

☆ Save 99 Cite Cited by 5098 Related articles All 18 versions Web of Science: 1858 »»

33

# Measuring Policy Uncertainty

## THE QUARTERLY JOURNAL OF ECONOMICS

Vol. 131 November 2016 Issue 4

### MEASURING ECONOMIC POLICY UNCERTAINTY\*

SCOTT R. BAKER  
NICHOLAS BLOOM  
STEVEN J. DAVIS

We develop a new index of economic policy uncertainty (EPU) based on newspaper coverage frequency. Several types of evidence—including human readings of 12,000 newspaper articles—indicate that our index proxies for movements in policy-related economic uncertainty. Our U.S. index spikes near tight presidential elections, Gulf Wars I and II, the 9/11 attacks, the failure of Lehman Brothers, the 2011 debt ceiling dispute, and other major battles over fiscal policy. Using firm-level data, we find that policy uncertainty is associated with greater stock price volatility and reduced investment and employment in policy-sensitive sectors like defense, health care, finance, and infrastructure construction. At the macro level, innovations in policy uncertainty foreshadow declines in investment, output, and employment in the United States and, in a panel vector autoregressive setting, for 12 major economies. Extending our U.S. index back to 1900, EPU rose dramatically in the 1930s (from late 1931) and has drifted upward since the 1960s. *JEL Codes:* D80, E22, E66, G18, L50.

- Based on the dictionary approach, count the number of news articles containing at least one key word from the three categories: economy, policy, and uncertainty. Use these counts to predict the level of economic policy uncertainty, defined as the simple average of the counts across different newspapers
- The created index is validated by a human audit, i.e., it is highly correlated with the human-coded index.

### Measuring economic policy uncertainty

[SR Baker, N Bloom, SJ Davis - ... quarterly journal of economics, 2016 - academic.oup.com](#)

... battles over fiscal policy. Using firm-level data, we find that policy uncertainty is associated with greater stock price volatility and reduced investment and employment in policy-sensitive ...

☆ Save 99 Cite Cited by 10872 Related articles All 53 versions Web of Science: 4167 »»

34

34

## Media Slant

**ECONOMETRICA**  
JOURNAL OF THE ECONOMETRIC SOCIETY

Full Access

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

Matthew Gentzkow, Jesse M. Shapiro

First published: 08 February 2010 | <https://doi.org/10.3982/ECTA7195> | Citations: 861  
Get it @ NYU

PDF TOOLS SHARE

**Abstract**

We construct a new index of media slant that measures the similarity of a news outlet's language to that of a congressional Republican or Democrat. We estimate a model of newspaper demand that incorporates slant explicitly, estimate the slant that would be chosen if newspapers independently maximized their own profits, and compare these profit-maximizing points with firms' actual choices. We find that readers have an economically significant preference for like-minded news. Firms respond strongly to consumer preferences, which account for roughly 20 percent of the variation in measured slant in our sample. By contrast, the identity of a newspaper's owner explains far less of the variation in slant.

**A measure of media bias**

T Groseclose, J Milyo - The quarterly journal of economics, 2005 - academic.oup.com

... We measure media bias by estimating ideological scores for several major media outlets. ...  
Our results show a strong liberal bias: all of the news outlets we examine, except Fox News' ...  
☆ Save 59 Cite Cited by 1554 Related articles All 25 versions Web of Science: 467

RE ↴  
Cite Cited by 2330 Related articles All 33 versions Web of Science: 710

What drives media slant? Evidence from US daily newspapers  
M Gentzkow, JM Shapiro  
Econometrica, 2010 · Wiley Online Library

We construct a new index of media slant that measures the similarity of a news outlet's language to that of a congressional Republican or Democrat. We estimate a model of newspaper demand that incorporates slant explicitly, estimate the slant that would be chosen if newspapers independently maximized their own profits, and compare these profit-maximizing points with firms' actual choices. We find that readers have an economically significant preference for like-minded news. Firms respond strongly to consumer preferences, which account for roughly 20 percent of the variation in measured slant in our sample. By contrast, the identity of a newspaper's owner explains far less of the variation in slant.

35

35

## Industry Segmentation

Text-Based Network Industries and Endogenous Product Differentiation

Gerard Hoberg and Gordon Phillips

PDF PDF PLUS Abstract Full Text Supplemental Material

**Abstract**

We study how firms differ from their competitors using new time-varying measures of product similarity based on text-based analysis of firm 10-K product descriptions. This year-by-year set of product similarity measures allows us to generate a new set of industries in which firms can have their own distinct set of competitors. Our new sets of competitors explain specific discussion of high competition, rivals identified by managers as peer firms, and changes to industry competitors following exogenous industry shocks. We also find evidence that firm R&D and advertising are associated with subsequent differentiation from competitors, consistent with theories of endogenous product differentiation.

**Text-based network industries and endogenous product differentiation**

G Hoberg, G Phillips - Journal of Political Economy, 2016 - journals.uchicago.edu

... industries as time-varying intransitive networks. We name these new industries text-based network industry ... Relative to existing industry classifications, our text-based classifications ...  
☆ Save 59 Cite Cited by 1884 Related articles All 20 versions Web of Science: 597

- Produce the 2-grams and 3-grams by speaker and select the top one thousand phrases through a chi-square test to identify the frequently and asymmetrically used phrases by Democrats and Republicans.
- Predict newspaper slant from the counts of the selected phrases by a two-stage supervised generative method.

- Classify industries based on product descriptions of company disclosure text, the 10-K report.
- The cosine-similarities between the token counts of different product offerings are computed.
- Industries are defined by clustering firms according to their cosine similarities.

36

36

## Transparency and deliberation within the FOMC: A computational linguistics approach

S Hansen, M McMahon, A Prat

The Quarterly Journal of Economics, 2018 • academic.oup.com

### Abstract

How does transparency, a key feature of central bank design, affect monetary policy makers' deliberations? Theory predicts a positive discipline effect and negative conformity effect. We empirically explore these effects using a natural experiment in the Federal Open Market Committee in 1993 and computational linguistics algorithms. We first find large changes in communication patterns after transparency. We then propose a difference-in-differences approach inspired by the career concerns literature, and find evidence for both effects. Finally, we construct an influence measure that suggests the discipline effect dominates.

 Oxford University Press

SHOW LESS ^

 Save  Cite Cited by 726 Related articles All 37 versions Web of Science: 163 

## FOMC Transparency

- Use LDA to study 149 FOMC meeting transcripts during Alan Greenspan's tenure.
- Measure the proportion of language devoted to K topics (estimated from the fitted topic model) for each FOMC member at each meeting.
- DiD suggests that moving to a more transparent system prompts inexperienced FOMC members to discuss a wider range of topics.

37

37