

Reading Between (the Party) Lines

by

Sophie Beiying Chou

Submitted to the MIT Media Lab,
School of Architecture and Planning
in partial fulfillment of the requirements for the degree of

MS in Media Arts and Sciences

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2016

© Massachusetts Institute of Technology 2016. All rights reserved.

Author
MIT Media Lab
May 5, 2016

Certified by
Deb Roy
Associate Professor
Thesis Supervisor

Accepted by
WHO IS THE CHAIR(WO)MAN?
Chairman, Department Committee on Graduate Theses

Reading Between (the Party) Lines

by

Sophie Beiying Chou

Submitted to the MIT Media Lab
on May 5, 2016, in partial fulfillment of the
requirements for the degree of
MS in Media Arts and Sciences

Abstract

TO-DO

Thesis Supervisor: Deb Roy
Title: Associate Professor

The following people served as readers for this thesis:

Sepandar Kamvar.....
Associate Professor of Media Arts and Sciences
MIT Media Lab

Iyad Rahwan
Associate Professor of Media Arts and Sciences
MIT Media Lab

Acknowledgments

Thank you !!

Contents

1	Introduction	7
2	The Power of (Percieved) Media Bias	9
2.1	The Effects of Media Bias	9
2.2	The Role of the Reader in Perceptions of Bias	9
2.3	The Role of Media Brands in Perceptions of Bias	10
2.4	The Role of Language [Policial Persuasion]	11
2.4.1	Language and Politics	11
2.4.2	The Seductive Allure [... of Simple] Language	11
2.5	Importance in Political Outcomes	13
2.6	The 2016 Elections	13
2.6.1	Criticism of Media Bias	13
3	Data Collection	15
3.1	The Electome	15
3.2	Story Collection	16
3.3	Election Classification	18
3.4	Topic Classification	19
3.5	Flesch-Kincaid Readability Tests	20
4	Experimental Design	23
4.1	Data Selection	23
4.1.1	Topic Selection	24

4.1.2	Flesch Kincaid Cutoffs	24
4.1.3	Redaction of Stories	24
4.2	CrowdFlower	24
4.3	Demographic Survey	24
4.4	Political Affiliation Survey	24
4.5	Quality Assurance	24
5	Pre-Survey Analysis	25
5.1	Topic Analysis	25
5.2	Flesch-Kincaid Analysis	25
5.2.1	Comparisons to other Reading Level Tests	25
6	Study	27
6.1	Demographics of Readers	27
6.2	Overall Bias Reportings	27
7	Analysis	29
7.1	Media Brand Effect	29
7.2	Reading Level Effect	29
7.3	Other Linguistic Cues	29

Chapter 1

Introduction

Most Americans say that they want to read news that's unbiased. A survey from Pew Research in 2012 showed that more than two-thirds (68%) of readers want to read political articles with a neutral stance, compared to just a little less than a quarter (23%) of those who want to read those stories that share their point of view.¹ But what exactly does that mean?

To begin with, whether or not we perceive news as biased is biased in itself. Conservative readers tend to view media as more biased than both Democrats and Independents (49% to 32% and 35%, respectively)[?].

The Hostile Media Effect, first studied by Vallone, Ross, and Lepper in 1985, gives one possible explanation for discrepancies: it describes a phenomenon where people with strong stances on an issue tend to perceive media covered as biased against their opinions, even on the same article.²

Clearly, finding bias in news depends on who the reader is as much as what they are reading.

In my thesis, I seek to examine the effects of context versus content in perceptions of media bias. In particular, when the context of a story is removed, how do linguistic

features, in particular reading level and vocabulary, in the content affect the reader? Although studies have been conducted to both examine the psychological effect of wording on believability (see “Seductive Allure”) and the impact of media brands and bias (see Baum, 2008), I seek to combine and contrast the two.

To do so, I will perform an A/B study for a broad range of readers to read and annotate political news stories (collected daily and sorted using a machine learning classifier). Each story is determined to be primarily about one political candidate and one topic computationally. In the control group, readers are given the full text of the article with no additional content. In the experimental group, readers are given a link to the original article complete with the byline, publication, and images. Stories are classified as either “high reading level,” “average reading level,” or “low reading level” by the Flesch-Kincaid test.

For each reader, I will collect their demographic information, and self-reported political stances. I will then analyze the effects of reading level versus media brand in the reader’s perception of the article.

I want to measure just how strong the effect of the media brand and the reader’s beliefs are.

Chapter 2

The Power of (Perceived) Media Bias

2.1 The Effects of Media Bias

Why is media bias IMPORTANT? Why is the problem IMPORTANT?

- Fox News Effect - Does the media matter

2.2 The Role of the Reader in Perceptions of Bias

It comes as no surprise that our own political stances have a significant effect in our perceptions of bias in the media.

In even seemingly neutral stories, partisans tend to view reporting as biased against their own views. This phenomenon—deemed the “hostile media effect”—was first studied at Stanford University by Robert P. Vallone, Lee Ross, and Mark R. Lepper in 1985 [8]. Although “true” neutrality of a story is nearly impossible to quantify due to the subjective nature of the concept, Vallone et. al were able to successfully demonstrate that partisans of *both* sides (pro-Israeli and pro-Arab) viewed the same

news segments as hostile towards their beliefs and favorable to the other side.

? Perceptions of media bias, then, have as much to do as self-serving motivations to secure preferential treatment as they do with the media itself.

? The political leanings of the reader are essential considerations when attempting to measure other factors that contribute to bias. In

2.3 The Role of Media Brands in Perceptions of Bias

The media, of course, is not just one unified mass, and in an increasingly fragmented ecosystem, the role of media brands is a crucial factor in the perception of bias. Although most research [2]

For instance, most research on the hostile media phenomenon conceptualizes the news media as an undifferentiated mass of information sources that individuals can (and do) reasonably characterize as having a uniform political orientation (Giner-Sorolla and Chaiken 1994, Peffley et al. 2001, Eveland and Shah 2003). Yet, the past two decades have seen a dramatic increase in the number and variety of news sources. One consequence is that Democrats and Republicans are increasingly likely to differ systematically in their assessments of specific media outlets.

With the decline of print newspapers, a diverse number of new platforms and web-centric publications have risen.

How do you control for the above things in your study?

2.4 The Role of Language [Policial Persuasion]

2.4.1 Language and Politics

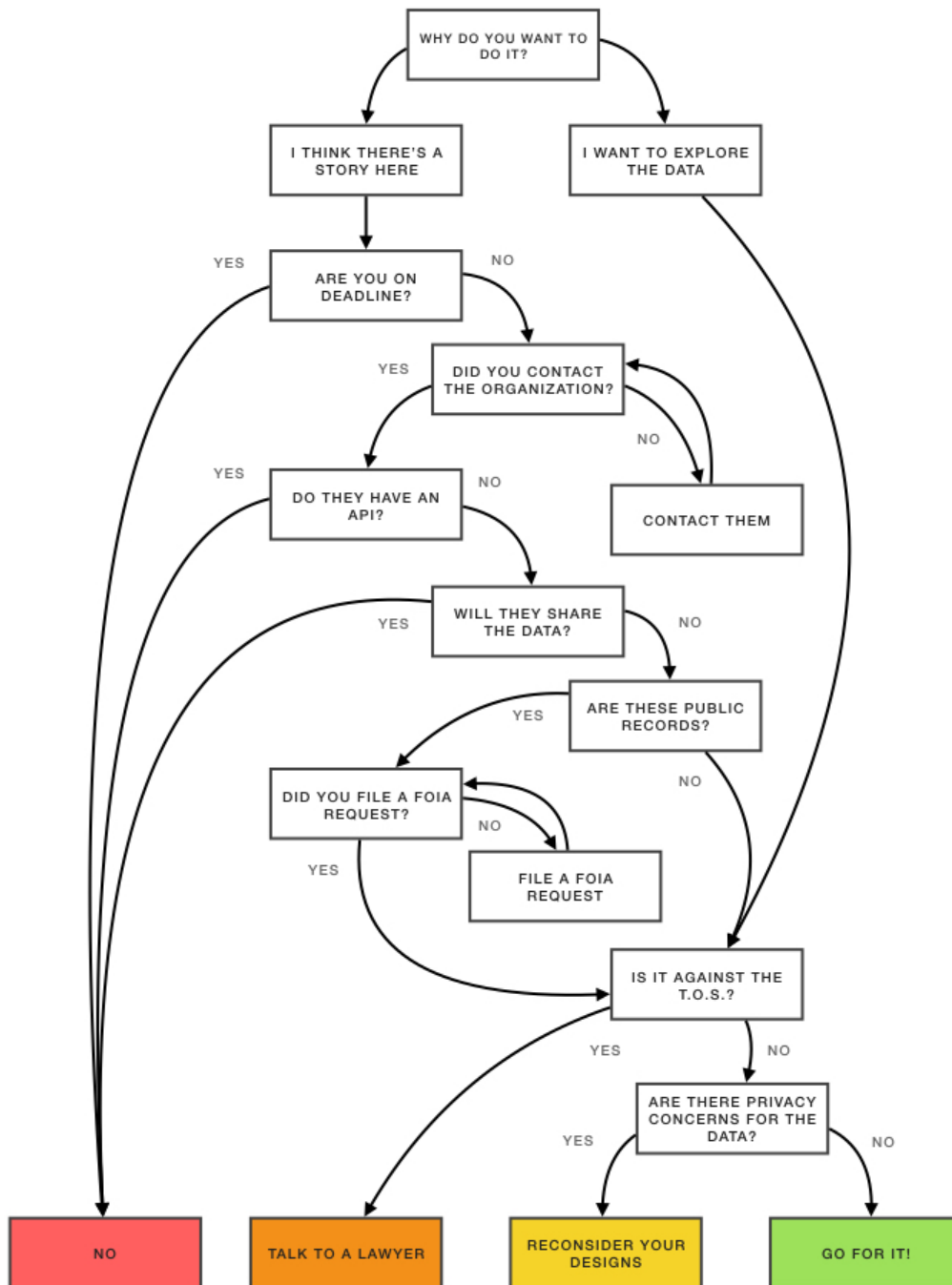
Presidential speeches degrading over time– ie simple language appeals to the masses in politics

2.4.2 The Seductive Allure [... of Simple] Language

But we trust complex language for explaining technical facts

Test image

Should You Build a Scraper?



2.5 Importance in Political Outcomes

Fox news effect

2.6 The 2016 Elections

2.6.1 Criticism of Media Bias

(Obama Speech)

So.... are you what you cover?

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 3

Data Collection

3.1 The Electome

The Electome is a large, collaborative, and ongoing effort in the Laboratory for Social Machines that seeks to analyze the “competition of ideas” in the upcoming 2016 elections. It does so by using techniques in natural language processing, machine learning, and network analysis to make sense of “big data” collected from two main sources: traditional media (online versions of news publications) and social media (Twitter) [9].

This thesis, which emerged from the Electome, examines a narrowed portion of the first dataset centered around specific topics and candidates. The following section will describe the methods used to gather this dataset as well as shared machine learning tools for article classification.

3.2 Story Collection

News articles from 14 different news publications were systematically collected every hour from RSS feeds beginning from January 2015. The outlets tracked are:

- CNN
- Fox News
- The Wall Street Journal
- ProPublica
- Politico
- McClatchy
- The Washington Post
- BuzzFeed (News only)
- National Public Radio (NPR)
- The Huffington Post
- The Associated Press
- Reuters
- The New York Times
- The Los Angeles Times

The above outlets were chosen to form a diverse subset of the current U.S. news ecosystem, including a combination of private and public, liberal and conservative, legacy and new media publications. Also included are wire services and a mix of media delivery formats for which the outlet is known (radio, television, print, or web).

Steps to collect the news stories were as follows:

1. For each news publication:
 - (a) Use regular expressions to extract all RSS feed urls for a news site.
 - (b) For each RSS feed:
 - i. Parse feed using open source xml reader library, Feedparser.
 - ii. For each link to a story in the feed:
 - A. Parse html using BeautifulSoup 3 (an open source python library)
 - B. Insert headline, authors, story text, publication date and retrieval date into an SQL database.

Data depulication (by story url and headline) is then performed to ensure only one

copy of each article is in the database. This step is necessary as articles from wire services often appear across many outlets and effect aggregate text analysis.

On average, 2,000 stories are collected per day across all outlets. However, volume follows a consistent pattern of fluctuation depending on weekday, ranging from approximately 1,000 to 3,000 stories.

[INSERT HERE GRAPH OF NEWS STORIES VOLUME BY WEEKDAY]

As of March 1st, 2016, there were 855,000 stories collected in the database and 43,000 journalists.

3.3 Election Classification¹

All stories collected from the sources above are passed through a machine learning classifier to determine if they are primarily about the 2016 U.S. elections. This thesis examines only those articles classified as election related.

The election classifier consists of a binary Maximum Entropy (MaxEnt) text classifier using Bag-of-Word (BoW) features selected from the news articles [6]. The features are ranked according to the chi-squared test (where high scores indicate that the null hypothesis of independence should be rejected and thus the occurrence of class and term are dependent) with a cutoff of 20,000. We use the open-source Python library scikit-learn for the implementation our MaxEnt classifier [7].

The classifier is trained on a balanced dataset of 1,000 manually labelled news articles and evaluated on a separate balanced test set of 300 articles. We achieved a precision of 90% and recall of 91% (F-score of 92%).

Between January 1, 2015 and March 1, 2016 there were 24,837 articles with over 90%

¹This section features shared machine learning tools within the Electome, with acknowledgements to Prashanth Vijayaraghavan.

confidence level of being election related. The number of stories classified as such has increased over time as election day nears.

[INSERT % ELECTION/ % TOTAL STORIES CHART HERE]

3.4 Topic Classification²

The final step of article processing within the Electome pipeline for this experimental dataset is the application of a 22-topic classifier. The following 22 topics were curated within the Laboratory for Social Machines as central issues of discussion within the election:

- Income Inequality
- Environment/Energy
- Jobs/Employment
- Guns
- Racial Issues
- Foreign Policy/National Security
- LGBT Issues
- Ethics
- Education
- Financial Regulation
- Budget/Taxation
- Veterans
- Campaign Finance
- Surveillance/Privacy
- Drugs
- Justice
- Abortion
- Immigration
- Trade
- Health Care
- Economy
- Other

3,000 articles classified as election related by the methods detailed in section 3.3 were manually labelled to form our training dataset. Articles were labelled as belonging to one or more topics. We then used a two-step model to create the classifier, due to the challenges of having a large number of classes and relatively small number of labeled stories. First, thousands of election related articles were inputted into a domain adaptive semi-supervised (stories were not all labeled) topic classification system. Then, a denoising autoencoder (DA) was used to learn salient features in an unsupervised fashion [10]. Then, these features were used to train a topic classifier using the labelled dataset.

²This section features shared machine learning tools within the Electome, with acknowledgements to Prashanth Vijayaraghavan.

The classifier was evaluated on an independent dataset of 400 manually annotated articles. We achieved a precision of 91% and a recall of 94% (weighted F-score of 92%).

3.5 Flesch-Kincaid Readability Tests

In this study, we focus primarily on the Flesch-Kincaid (F-K) tests for estimating text readability. Originally developed for the U.S. Navy in 1975 for assessing the difficulty of technical manuals, the F-K reading level corresponds roughly to U.S. grade level and the reading ease score is inversely proportional to the grade level on a scale from 0 to approximately 120 [3].

We chose the F-K tests over other comparable ones due to its popularity in educational assessment and other applications, including in legislation. For example, it is required by law in Florida that life insurance policies have a Flesch reading ease of 45 or greater (less than 12th grade in reading level) [4]. The F-K tests are also bundled in many common word processing services, including Microsoft Office Word. As a comparison, basic article analysis is also computed using the Gunning fog index (see Section 5.2.1).

The formula for Flesch reading ease is as follows:

$$206.835 - 1.015 \left(\frac{\text{total words}}{\text{total sentences}} \right) - 84.6 \left(\frac{\text{total syllables}}{\text{total words}} \right)$$

And for reading grade level:

$$0.39 \left(\frac{\text{total words}}{\text{total sentences}} \right) + 11.8 \left(\frac{\text{total syllables}}{\text{total words}} \right) - 15.59$$

The two formulas are not directly comparable due to the difference in weighting factors. For ease of metaphor, we use the grade level tests in our analysis. Syllable

length is highly weighted in this formula, so it is possible to generate a story of very high reading level that consists of a single word in a single sentence (the longest English word, *pneumonoultramicroscopicsilicovolcanoconiosi*, a type of lung disease, has a reading grade level of 197.2), which is a limitation of the method, since texts with polysyllabic words are not always necessarily more difficult to read.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 4

Experimental Design

4.1 Data Selection

For this study, we chose to analyze stories collected between January 1, 2016 (the start of the election year) and March 1, 2016 (Super Tuesday). Since a large number of states hold primary elections and caucuses on Super Tuesday, it is seen as an early indicator of candidate electability. All stories had been filtered through both the election (see section 3.3) and topic (see section 3.4) classifiers.

Based on the results of Super Tuesday, we selected four candidates for this study by delegate count: Hillary Clinton (1,279), Bernie Sanders (1,027), Donald Trump (743), and Ted Cruz (517) [1].

News articles were then separated into single-candidate stories (i.e. articles featuring primarily one candidate in the headline) to be able to measure more clearly the perceived bias per candidate. This was done programmatically using regular expressions to determine if a headline contained one candidate and one candidate only. A dictionary of related names was created to make sure that stories were correctly categorized (i.e. “Hillary”, “Clinton”, and “Hillary Clinton” were to be categorized as pertaining

to “Hillary Clinton” but not if preceded by “Bill”).

4.1.1 Publication Selection

For the purposes of this study, stories were examined from five outlets:

- CNN
- Fox News
- The New York Times
- The Wall Street Journal
- The Associated Press

The choices consist of two pairs of outlets in both print and television across the liberal-conservative divide, plus a wire service. Of the 14 outlets above, both Fox News and the Wall Street Journal have an audience that leans conservative compared to the overall population (27% mostly conservative viewers versus 17% in the overall population for Fox News and 22% mostly conservative viewers versus 17% in the overall population) measured by a 2014 Pew survey [5].

On the other hand, the New York Times and CNN both have audiences that lean mostly liberal (25% liberal versus 22% in all respondents for CNN and 25% for the New York Times). The Associated Press, which was not included in the survey, has members in outlets across the political divide and was chosen as an experimental control.

[MIGHT INCLUDE THOSE DISTRIBUTIONS HERE]

4.1.2 Topic Selection

4.1.3 Flesch Kincaid Cutoffs

4.1.4 Redaction of Stories

4.2 CrowdFlower

4.3 Demographic Survey

4.4 Political Affiliation Survey

4.5 Quality Assurance

-Filter by nationality - highest setting on crowdflower - Gold questions - time limits
- price

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 5

Pre-Survey Analysis

5.1 Topic Analysis

5.2 Flesch-Kincaid Analysis

5.2.1 Comparisons to other Reading Level Tests

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 6

Study

We ran this over n days blah blah

6.1 Demographics of Readers

6.2 Overall Bias Reportings

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 7

Analysis

7.1 Media Brand Effect

7.2 Reading Level Effect

7.3 Other Linguistic Cues

THIS PAGE INTENTIONALLY LEFT BLANK

Bibliography

- [1] March 1 election results 2016 updates - the washington post. <https://www.washingtonpost.com/2016-election-results/super-tuesday/>. (Accessed on 04/06/2016).
- [2] Matthew A Baum, Phil Gussin, et al. In the eye of the beholder: How information shortcuts shape individual perceptions of bias in the media. *Quarterly Journal of political science*, 3(1):1–31, 2008.
- [3] J Peter Kincaid, Robert P Fishburne Jr, Richard L Rogers, and Brad S Chissom. Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel. Technical report, DTIC Document, 1975.
- [4] The Florida Legislature. The 2015 florida statutes title xxxvii. <http://www.leg.state.fl.us/Statutes>. (Accessed on 04/05/2016).
- [5] Amy Mitchell, Jeffrey Gottfried, Jocelyn Kiley, and Katerina Eva Matsa. Political polarization & media habits. <http://www.journalism.org/2014/10/21/political-polarization-media-habits/>, Oct 2014. (Accessed on 04/02/2016).
- [6] Kamal Nigam, John Lafferty, and Andrew McCallum. Using maximum entropy for text classification. In *IJCAI-99 workshop on machine learning for information filtering*, volume 1, pages 61–67, 1999.
- [7] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *The Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [8] Robert P Vallone, Lee Ross, and Mark R Lepper. The hostile media phenomenon: biased perception and perceptions of media bias in coverage of the beirut massacre. *Journal of personality and social psychology*, 49(3):577, 1985.
- [9] Prashanth Vijayaraghavan, Soroush Vosoughi, and Deb Roy. Automatic detec-

tion and categorization of election-related tweets. In *10th International AAAI Conference on Web and Social Media*, 2016.

- [10] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103. ACM, 2008.