Analyzing mass media influence in public opinion using natural language processing and time series analysis





Pablo Balenzuela, Viktoriya Semeshenko, Sebastián Pinto & Federico Albanese SoPhy Lab, Physics Department, FCEyN, University of Buenos Aires

Agenda-Setting Theory

- How the media influence public opinion is a question usually faced and analyzed within the context of **Agenda Setting theory** (*Agenda-Setting, McCombs 1972*).
- Media outlets are represented by their agenda: <u>compilation of topics</u> <u>presented with different emphasis</u>.
- The theory, in its basic version, studies the correlation between the topics covered by the media and those that are prominent in the public.

Introduction: Field experiments

American Economic Journal: Applied Economics 2009, 1:2, 35–52 http://www.aeaweb.org/articles.php?doi=10.1257/app.1.2.35

> Does the Media Matter? A Field Experiment Measuring the Effect of Newspapers on Voting Behavior and Political Opinions[†]

By Alan S. Gerber, Dean Karlan, and Daniel Bergan*

We conducted a field experiment to measure the effect of exposure to newspapers on political behavior and opinion. Before the 2005 Virginia gubernatorial election, we randomly assigned individuals to a Washington Post free subscription treatment, a Washington Times free subscription treatment, or a control treatment. We find no effect of either paper on political knowledge, stated opinions, or turnout in post-election survey and voter data. However, receiving either paper led to more support for the Democratic candidate, suggesting that media slant mattered less in this case than media exposure. Some evidence from voting records also suggests that receiving either paper led to increased 2006 voter turnout. (JEL D72, L82)

Field Experiment: (Governor elections in Virginia, 2005)

Three groups in Prince William County, Virginia:

- Group 1: free suscription (15/10 8/11) to Washington Post.
- Grupo 2: free suscription (15/10 8/11) to Washington Times.
- Group 3: Control

<u>Subjects:</u> registered voters. People already suscribed to a given newspaper were excluded.

<u>Post-elections survey</u>: Current events, news, politics, political leanings, whether they voted, for whom.

Conclusions:

While reading the newspaper did not increase the number of people who voted, it did influence who voted. They saw: 8% were more likely to vote for Democrats regardless of the newspaper they read, suggesting that media bias mattered less in this case than media exposure.

Gerber, A. S., Karlan, D., & Bergan, D. (2009). Does the media matter? A field experiment measuring the effect of newspapers on voting behavior and political opinions. American Economic Journal: Applied Economics, 1(2), 35-52.

Introduction: Web data analysis



Results suggest that voters seek information only when they are considering changing their vote.

[1] Yasseri, T., & Bright, J. (2016). Wikipedia traffic data and electoral prediction: towards theoretically informed models. *EPJ Data Science*, *5*(1), 1-15. ISO 690

Analysis of information consumption on Wikipedia sites related to different European elections. Increases in the consumption of these pages correlates with:

- Increased participation in elections
- Increase of votes in given political parties.

2009 & 2014 European Elections



Our proposal

- To understand what features of the news media articles can influence public opinion in an electoral context.
- We focus on 2016 USA presidential elections and analyze:
 - News articles mentioning candidates Hillary Clinton and Donald Trump
 - Electoral polls on voting intentions for each candidate.
- We analyze different features of news media articles :
 - Amount of mentions to each candidate.
 - Sentiment bias toward each candidate.
 - Topics of discussion.
- We studied correlation and causality measures between these aspects and the time series of the surveys.

Surveys as proxy of public opinion





- 263 USA national surveys
- Different pollsters: NBC, New York Times, LA Times, CBS, Fox News, Gravis, ABC, IBD (among others)

Final difference:

+3.2 Clinton according to average surveys +2.1 Clinton according to final election

+2.1 Clinton according to final elections results

Corpus of news articles

We select online articles containing at least the name of one of the two main candidates: Hillary Clinton (Democrat) and Donald Trump (Republican).

- The New York Times (5672 news articles)
- Fox News (5750)
- USA Today (833)
- **CNN** (2920)



At least one of the candidates is mentioned in the notes.

1 – Mentions to candidates





1 – Mentions to candidates

Mentions in Fox News:

1

Difference Clinton – Trump: 🦊



2 – Sentiment Analysis

Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank

Richard Socher, Alex Perelygin, Jean Y. Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng and Christopher Potts Stanford University, Stanford, CA 94305, USA richard@socher.org, {aperelyg, jcchuang, ang}@cs.stanford.edu {jeaneis, manning, cgpotts}@stanford.edu

Abstract

Semantic word spaces have been very useful but cannot express the meaning of longer phrases in a principled way. Further progress towards understanding compositionality in tasks such as sentiment detection requires richer supervised training and evaluation resources and more powerful models of composition. To remedy this, we introduce a Sentiment Treebank. It includes fine grained sentiment labels for 215,154 phrases in the parse trees of 11,855 sentences and presents new challenges for sentiment compositionality. To address them, we introduce the Recursive Neural Tensor Network, When trained on the new treebank, this model outperforms all previous methods on several metrics. It pushes the state of the art in single sentence positive/negative classification from 80% up to 85.4%. The accuracy of predicting fine-grained sentiment labels for all phrases reaches 80.7%, an improvement of 9.7% over bag of features baselines. Lastly, it is the only model that can accurately capture the effects of negation and its scope at various tree levels for both positive and negative phrases.



Figure 1: Example of the Recursive Neural Tensor Network accurately predicting 5 sentiment classes, very negative to very positive (--, -, 0, +, ++), at every node of a parse tree and capturing the negation and its scope in this sentence.

models to accurately capture the underlying phenomena presented in such data. To address this need, we introduce the Stanford Sentiment Treebank and a powerful Recursive Neural Tensor Network that can accurately predict the compositional semantic effects present in this new corpus.

The Stanford Sentiment Treebank is the first cor-

Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A., & Potts, C. (2013). Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing* (pp. 1631-1642). **Standford Core NLP**: They build a binary tree from each phrase to determine whether the phrase has a positive, negative or neutral connotation from :

- A database with positive, negative or neutral words.
- The sentence structure.
- The presence of intensifiers, appeasers or deniers.



2 - Sentiment Analysis

Steps:

1) Phrases where Donald Trump and Hillary Clinton are mentioned are detected.

2) Sentiment analysis is run on the sentences (in the case where both are mentioned, the sentence is split and then the SA algorithm is applied).

3) The number of positive, neutral, negative and total mentions of each candidate are counted.

A Look at Clinton's Marriage Woes Raises a Question: Really?

The Public Editor By LIZ SPAYD OCT. 4, 2016



Hillary Clinton talking to the media on her campaign plane. Doug Mills/The New York Times

Late last Saturday evening, The New York Times delivered an eye-popping scoop to its readers — tax documents of Donald Trump's that showed he

<u>New:</u> https://www.nytimes.com/2016/10/05/public-editor/hillaryclinton-bill-clinton-marriage.html

2 - Sentiment Analysis

- #C₊ = Fraction of positive mentions to Hillary Clinton
- #C₋ = Fraction of negative mentions to Hillary Clinton
- #T₊ = Fraction of positive mentions to Donald Trump
- #T_ = Fraction of negative mentions to Donald Trump

Sentiment Bias (SB_i): measure de bias towards the candidates in the "i" media outlet. (if SB>0, the bias is towards Clinton compared with Trump and viceversa)

$$SB_i = (\#C_+ - \#C_-) - (\#T_+ - \#T_-)$$



SB_i = es nonnegative with p< 0.001 (boostrapping) for all media oulets

Obs: While SB is positive for all, it is higher for NYT and very close to zero for Fox News. (We will see that the time dependence of this statistic is more informative).

2 - Sentiment Analysis



3 – Topic decomposition:

We describe each news articles as a frequency vector in the semantic space (tf-idf)

Topic's keywords

U.S. Cyberforce Was Deployed to Estonia to Hunt for Russian Hackers

An operation ahead of the November election was part of stepped-up efforts by the military to stop Russian interference in American politics.

f ¥ = * 🗌



Voting machines being checked on Election Day last month in Houston. No fureign power was able to discupt the American vote. Tumir Ralifa for The New York Times

By Julian E. Barnes

WASHINGTON — The United States deployed operatives to Estonia in the weeks before the November election to learn more about defending against Russian hackers as part of a broader effort to hunt down foreign cyberattacks, American and Estonian officials said.

Estonian officials believe the growing cooperation with the United States will be an important deterrent to any attacks by neighboring



News corpus (tf-idf)



Topic decomposition: We look identify the group of news speaking of the same topic and described by a limited number of keywords (unsupervised clustering).

3 - Topic decomposition

We describe each news articles as a frequency vector in the semantic space (tf-idf)

V=[1 idf(t) ,...,F] # of time term "t" appears in the article Specificity factor of term "t" Topic's keywords News corpus (tf-idf) Similar 50 50 articles NMF 100 100 200 Documentos 250 Documentos 500 520 300 300 350 350 2000 0 1000 3000 4000 5000 6000 7000 8000 1000 2000 3000 4000 5000 6000 7000 Términos Términos

Topic decomposition: We look identify the group of news speaking of the same topic and described by a limited number of keywords (unsupervised clustering).

U.S. Cyberforce Was Deployed to Estonia to Hunt for Russian Hackers

An operation ahead of the November election was part of stepped-up efforts by the military to stop Russian interference in American politics.

f ¥ = *



Voting machines being checked on Election Day last month in Houston. No foreign power was able to discupt the American vote. Tamir Kalifa for The New York Times

By Julian E. Barnes

WASHINGTON — The United States deployed operatives to Estonia in the weeks before the November election to learn more about defending against Russian hackers as part of a broader effort to hunt down foreign cyberattacks, American and Estonian officials said.

Estonian officials believe the growing cooperation with the United States will be an important deterrent to any attacks by neighboring

Emergent Topics





3 – Topic decomposition and media agenda



Торіс	NYT (SRL)	Fox (SRL)	CNN (SRL)	USA (SRL)
Clinton email controversy	-0.46 (10-20)	-0.42 (11-20)	-0.45 (13-20)	-0.54 (10-20)
Economy	0.56 (4-20)	0.59 (8-20)	0.48 (5-18)	0.40 (10-15)
Clinton foundation affair	-0.53 (3-20)	-0.43 (15-20)	-0.53 (1-20)	-0.40 (5-12)
Immigration	-0.42(0-12)	-0.44 (5-20)		-0.47 (12-20)
Foreign affairs	0.44 (17–20)	1/	1	1



3 – Topic decomposition and media agenda



Торіс	NYT (SRL)	Fox (SRL)	CNN (SRL)	USA (SRL)
Clinton email controversy	-0.46 (10-20)	-0.42 (11-20)	-0.45 (13-20)	-0.54 (10-20)
Economy	0.56 (4-20)	0.59 (8-20)	0.48 (5-18)	0.40 (10-15)
Clinton foundation affair	-0.53 (3-20)	-0.43 (15-20)	-0.53 (1-20)	-0.40 (5-12)
Immigration	-0.42(0-12)	-0.44 (5-20)	20	-0.47 (12-20)
Foreign affairs	0.44 (17–20)		10	1



4 – Combined topic decomposition and sentiment analysis

Sentiment Bias (SB): relative number of positive and negative mentions of Clinton and Trump.

$$SB = (\#C_{+} - \#C_{-}) - (\#T_{+} - \#T_{-})$$

4 – Combined topic decomposition and sentiment analysis

Sentiment Bias (SB): relative number of positive and negative mentions of Clinton and Trump (by topic and media outlet).

$$SB = (\#C_{+} - \#C_{-}) - (\#T_{+} - \#T_{-})$$

Tópico	SB_{NYT}	SB_{FN}	SB_{CNN}	SB_{USA}
$Clinton\ email\ controversy$	-0.475	-0.429	-0.302	-0.315
Economy	0.332	0.070	0.152	0.168
$Clinton\ foundation\ affair$	-0.256	-0.257	-0.304	-
Immigration	0.501	0.347	0.306	0.382
Foreign affairs	0.146	0.053	0.115	0.166

Bias of each topic (SB < 0 unfavorable to Clinton).

4 – Combined topic decomposition and sentiment analysis

Sentiment Bias (SB): relative number of positive and negative mentions of Clinton and Trump.

$$SB = (\#C_{+} - \#C_{-}) - (\#T_{+} - \#T_{-})$$

Tópico	SB_{NYT}	SB_{FN}	SB_{CNN}	SB_{USA}	SB < 0
Clinton email controversy	-0.475	-0.429	-0.302	-0.315	\triangleright
Economy	0.332	0.070	0.152	0.168	
Clinton foundation affair \triangleleft	-0.256	-0.257	-0.304	-	\triangleright
Immigration	0.501	0.347	0.306	0.382	
Foreign affairs	0.146	0.053	0.115	0.166	

Bias of each topic (SB < 0 unfavorable to Clinton).









5 – Granger causality

It reveals if the dynamics of a temporal series have a causal relationship with another time series.

Basically, it consist in determine if additional information provided by a second lagged time series, $y(t - \tau)$, improves the forecast of x(t).

$$\Delta CT(t) = CT(t) - CT(t-1) = w_t$$
 (1)

$$\Delta CT(t+\tau) = \beta \cdot \Delta T_i(t) + w_{t+\tau}$$
 (2)

We say that the **coverage** of a given issue effectively **affects** the **difference** between voting intention between candidates when the parameter β in la equation (2) significatively **differs from zero.**

Τόριςο	NYT	FN	CNN	USA Today
Clinton email controversy	-	β < 0 (τ=19)		
Clinton Foundation affaire	β<0 (t=19)			
Inmigration	β < 0 (t=10)			
Economy		$\beta > 0$ (τ =11-16)	$\beta > 0$ (τ =11-16)	$\beta > 0$ (τ =11-16)

5 – Granger causality

It reveals if the dynamics of a temporal series have a causal relationship with another time series.

Basically, it consist in determine if additional information provided by a second lagged time series, $y(t - \tau)$, improves the forecast of x(t).

$$\Delta CT(t) = CT(t) - CT(t-1) = w_t$$
 (1)

$$\Delta CT(t+\tau) = \beta \cdot \Delta T_i(t) + w_{t+\tau}$$
 (2)

We say that the **coverage** of a given issue effectively **affects** the **difference** between voting intention between candidates when the parameter β in la equation (2) significatively **differs from zero.**



5 – Granger causality

It reveals if the dynamics of a temporal series have a causal relationship with another time series.

Basically, it consist in determine if additional information provided by a second lagged time series, $y(t - \tau)$, improves the forecast of x(t).

$$\Delta CT(t) = CT(t) - CT(t-1) = w_t$$
 (1)

$$\Delta CT(t+\tau) = \beta \cdot \Delta T_i(t) + w_{t+\tau}$$
 (2)

We say that the **coverage** of a given issue effectively **affects** the **difference** between voting intention between candidates when the parameter β in la equation (2) significatively **differs from zero.**



Conclusions

- A combined approach of topic decomposition and sentiment analysis to news stories mentioning candidates in an election allowed us to see which relevant aspects influence people's voting intentions.
- The analysis of correlations by topic with the difference in voting intentions observed in the polls allowed us to understand the fundamental role played by the topics that favor or disfavor the candidates.
- Sentiment analysis discriminated by topic allowed us to understand how a given topic (<u>Clinton email controversy</u>) played a key role in Trump's vote increase in the days leading up to the election.
- The Granger test is consistent with previous results and reinforces the interpretation of a causal relationship between the topic addressed by candidates in the media and voting behavior.

SoPhy Lab: The Social Physics Group



Thanks!!