



PRINCETON
UNIVERSITY

A

1. Motivation

Goal: Create a data visualizer that compares and
Top-rated and **Controversial-rated** Reddit he

- Users draw their own conclusions about how

Upvoted and downvoted textual analysis of “top” vs “ advice” posts

Jessica Zheng '19

Advi

contrasts
adlines.
Reddit treated

4a. Topic Modeling

NLTK Latent Dirichlet Allocation (LDA)



And Unbiased?

“controversial” Reddit headline

Author: Professor Brian Kernighan

IBM Watson

Topic Model

sexual harassment

harass female reporter

5. Findings

2 vulnerability

nes



Topic Model

r/technology on 10/16/2017



major events over the data collection period.
Las Vegas shooting, Net Neutrality debate etc.

- Existing Reddit visualizers do not account for post popularity, instead consider all posts.

Problem: Statistics are computed from content that is not meaningful or impactful on Reddit users – as few as 5% of submitted content actually becomes popular.



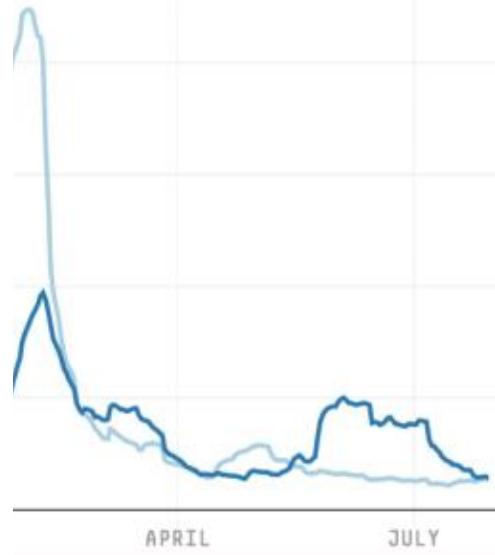
2. Approach

internet* Talks

I mostly male users of Reddit, anyway.

Isaac and Ritchie King

Enter terms



^ randomized model: problematic!

4b. Positive-Negative Sentiment

- NLTK's VADER Sentiment Analysis
- **Challenge:** VADER does not handle:
E.g. “Patients *unwillingly* treated
Norwegian mental hospitals” p
- **Solution:** Re-run neutral-scored h



returned more coherent topics

sentiment Analysis

was a good starting point
 use negative operators
 ended with electroshock at
 predicted a neutral (zero) score
 headlines through IBM Watson's

r/technology WPA:

Top

Other: 13.64%



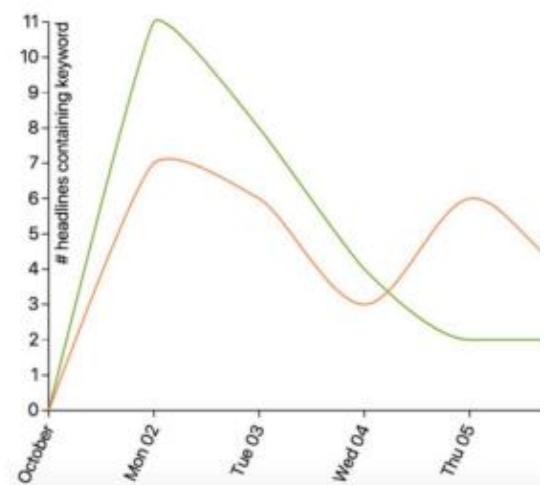
Tech News Site: 43.18%

r/news: desensitization

Las Vegas Shooting: 10/

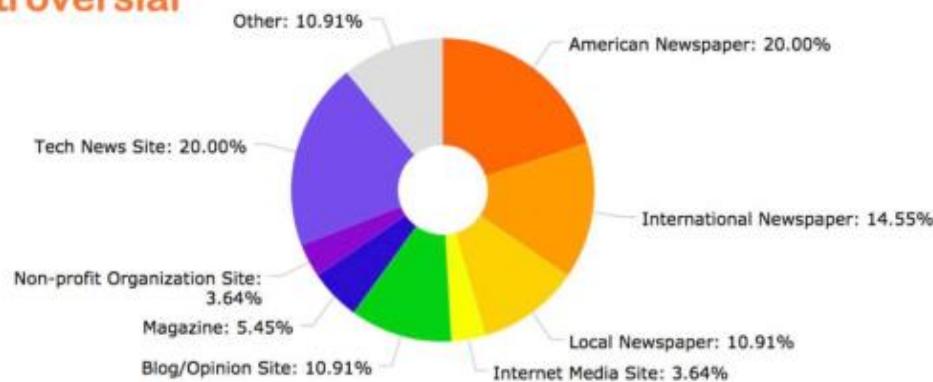
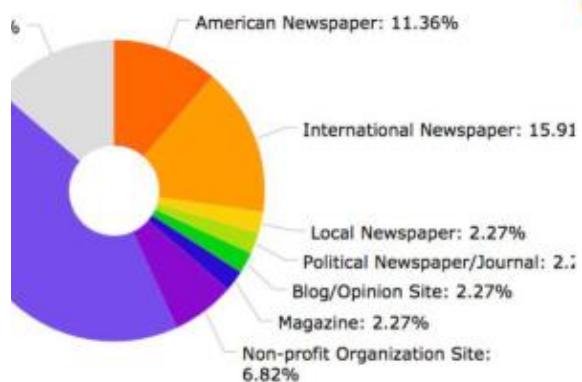
Keyword Frequency Over Time

vegas, shoot, gun



r/politics

Controversial



Organization to violence?

1-10/7

NYC Truck Attack: 10/31-11/4

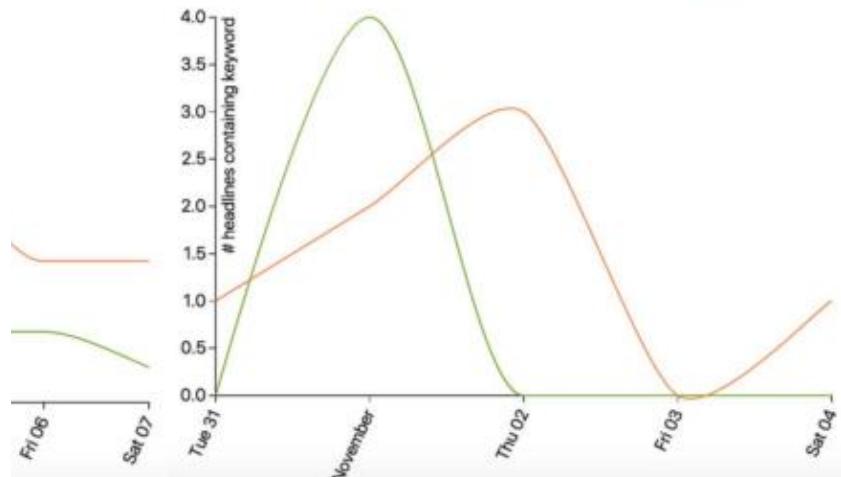
1e

Go

Keyword Frequency Over Time

nyc, new york, truck attack

Go

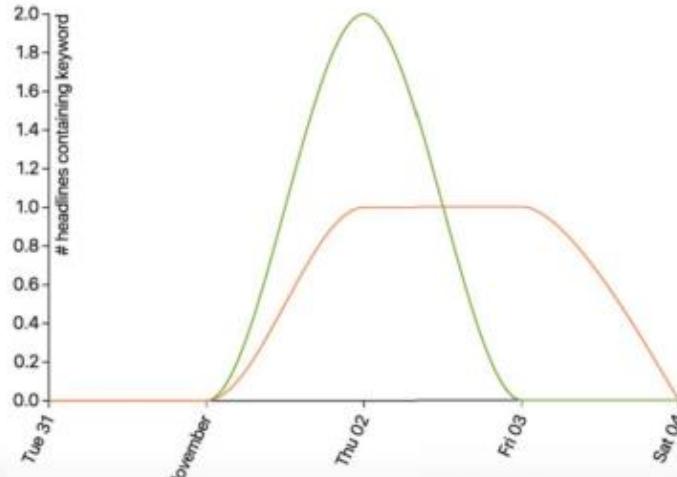


Thronton Walmart Shooting: 11/2-11/4

Keyword Frequency Over Time

Thornton, walmart shooting

Go



r/news

r/news

Analyze **Top-rated** and **Controversial-rated** Reddit posts

- topic modeling
- positive-negative sentiment analysis
- political sentiment analysis
- content source breakdown
- keyword frequency

Subreddits used: r/politics, r/news, r/worldnews, r/technology

Data collection period: October 1st 2017 to November 1st 2017

Observe patterns between analysis measures and content popularity (“top” or “controversial”, post count, number of comments)

3. Implementation

t headlines via:

/technology
ber 18th 2017.
indicators of
score, number

Sentiment Analyzer (which handles a second measurement).

4c. Political Sentiment Analysis

- Indico.io Text Analysis Service's AI
 - Negative operators: “I *don't* care”
 - Highly focused on conservative-leaning topics
- Somewhat better at detecting liberal topics, rather than liberal-conservative
- Machine learning has lots of room for improvement

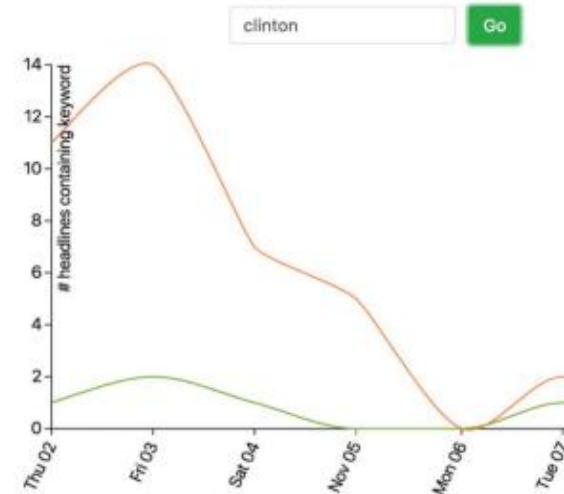
4d. Other Analysis Measures

ed negative operators) to obtain

ysis and Shortcomings

PI had limited accuracy
re for Trump.” scores more
than liberal leaning.
eral/conservative people and
native leanings
n for improvement!

controversial talks much
more about clinton

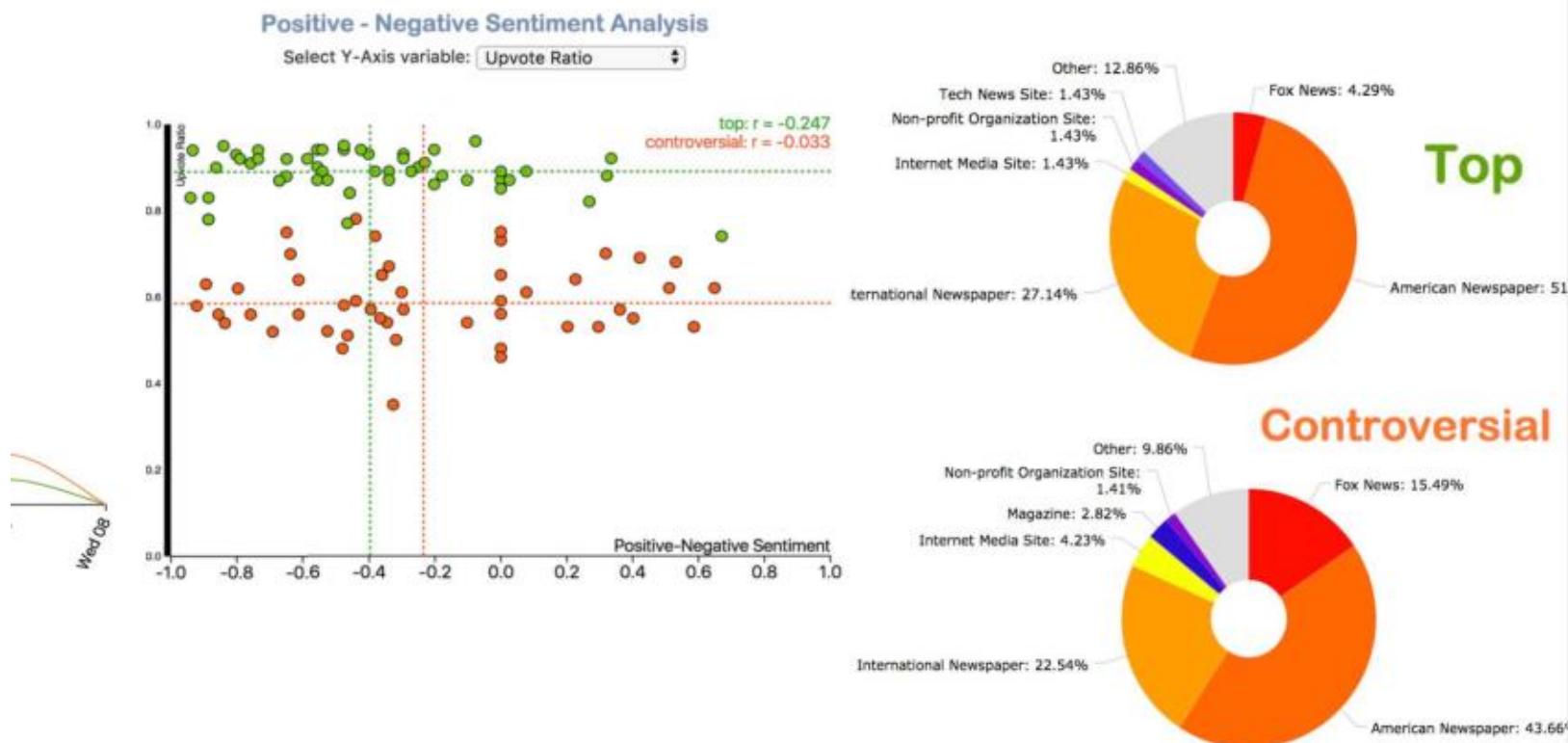


6. Usability Eva

- Effectiveness of a
- 4/5 users drew
- 5/5 users drew

top is on average more negative than controversial

controversial has more sources from Fox News



Evaluation

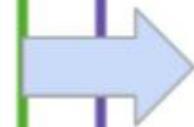
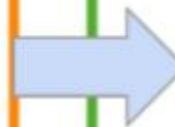
Analysis tool:

- ✓ the same conclusions I did for Las Vegas Event
- ✓ the same conclusions I did for the WPA2 event

Data Collection
Python Reddit
API Wrapper
(PRAW)

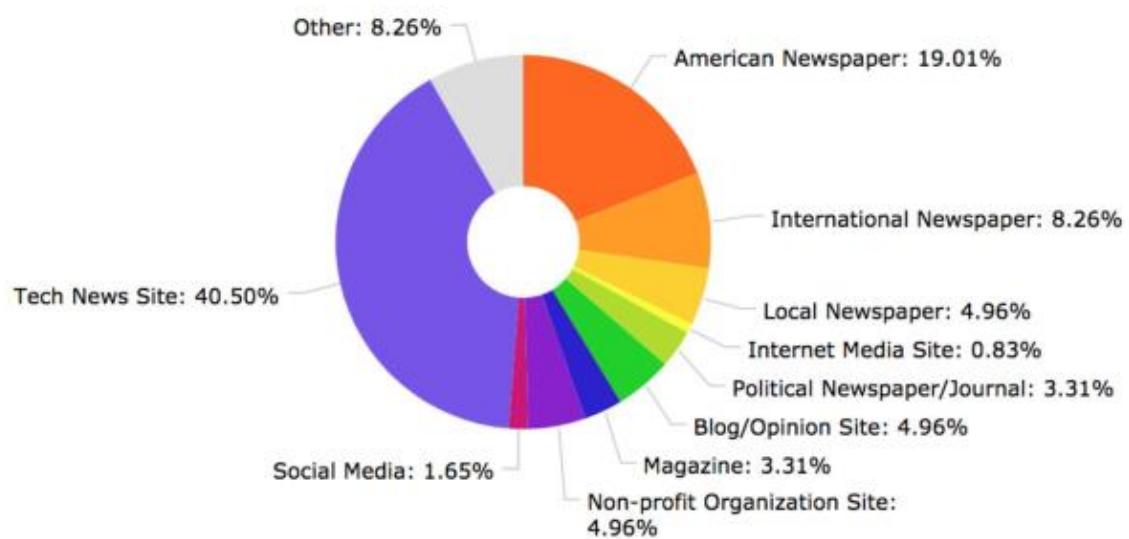
Analysis Tools
NLTK
IBM Watson
Indico

Data
D3.js
Pos



Content Source Breakdown

Content Type



Visualization

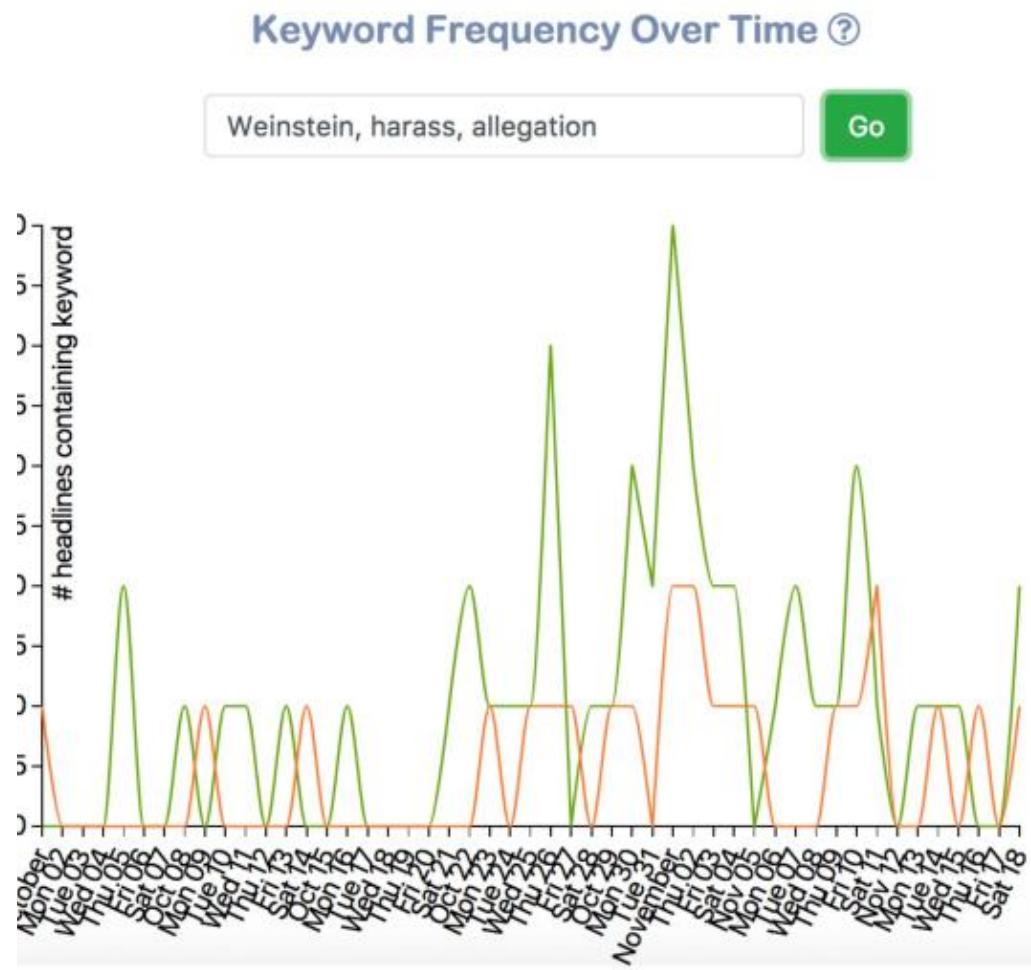
Visualizations

Flask

PostgreSQL DB

Heroku

Keyword Frequency in Headlines



- Usability issues w
resolved) and lim

7. Conclusions

- Topic model and
- Machine learning
- Implications:
 - For known cor
and promote h
- Future Work: Cor
create improved

were mostly concerned with learnability (easily
limitations of the sentiment analysis tools.

and Future Work

keyword frequency visuals were most informative
; still has room to go for sentiment analysis

controversial topics, Reddit moderators can intervene
healthy discussion

combine multiple analysis forms in the same visual,
sentiment classifiers, explore more subreddits.