# A Tool to Visualize Language Biases Over Time Using Reddit Comment Data

Filippo M. Libardi & Otto Mättas

# Related Work

- Very recent 2020 paper;

- NLP model to identify biased terms towards a list of concepts given a text corpus;

- Method validated on Google News;

- The model only returns a list of biased words, the usage of these is not studied at all;

- Their method allows to formalize language bias;

- We aimed at improving on the method by displaying the usage of each biased word trough time.

# Ferrer et al. (2020)



$$\text{Bias}(w, c_1, c_2) = \cos(\vec{w}, \vec{c_1}) - \cos(\vec{w}, \vec{c_2})$$

$$\text{Sent}(W) = \frac{1}{|W|} \sum_{w \in W} SA(w)$$

$$\text{Most Biased}(V, c_1, c_2) = \arg\max_{w \in V} \text{Bias}(w, c_1, c_2)$$

3

1/25/2021

# Our Additions



List of words representing the concept of "foreign" [1]

Extract N comments from subreddit

Model (A)

Train on data (A1)

Get most biased words towards L1 (A2)

[2] k most positively biased 🙂

[3] k most negatively biased 😠

End Tool

Time (B)

Extract UNIX timestamp out of comments (B1)

Bin comments based on time (N bins depends on what time range comments span). (B2)

Count occurrences of L2 and L3 in each time bin. (B3) Accumulate their sentiment value.

Normalising by total comments for that bin. We call this BIAS SCORE.

[4] Result in N bins containing normalised bias score per bin.

Visualize D1 over range from UNIX timestamp of bin one to UNIX timestamp bin N.

Legend
[1] L1
[2] L2
[3] L3
[4] D1

$$BiasScore = \frac{\sum\limits_{w \in C} Sent(w)}{|C|}$$

4

1/25/2021

# Issues

Ferrer et al. (2020)

APIs

Formatting

Subreddits

Ingestion

Processing

# Results and Validation



**Weekly avg.** bias score for most negatively biased words towards *"foreign"*.

20 November 2016
On November 20, 2016, Benjamin Marconi, a detective with the San Antonio Police Department, was shot to death in San Antonio, Texas. In the shooting, a motorist stopped his car, got out, and shot and wounded Marconi while the latter was sitting in his marked patrol car in front of the department's headquarters, writing a ticket for another driver during a routine traffic stop.

Source

3 - 9 September 2017
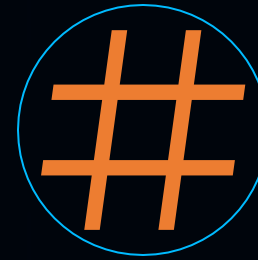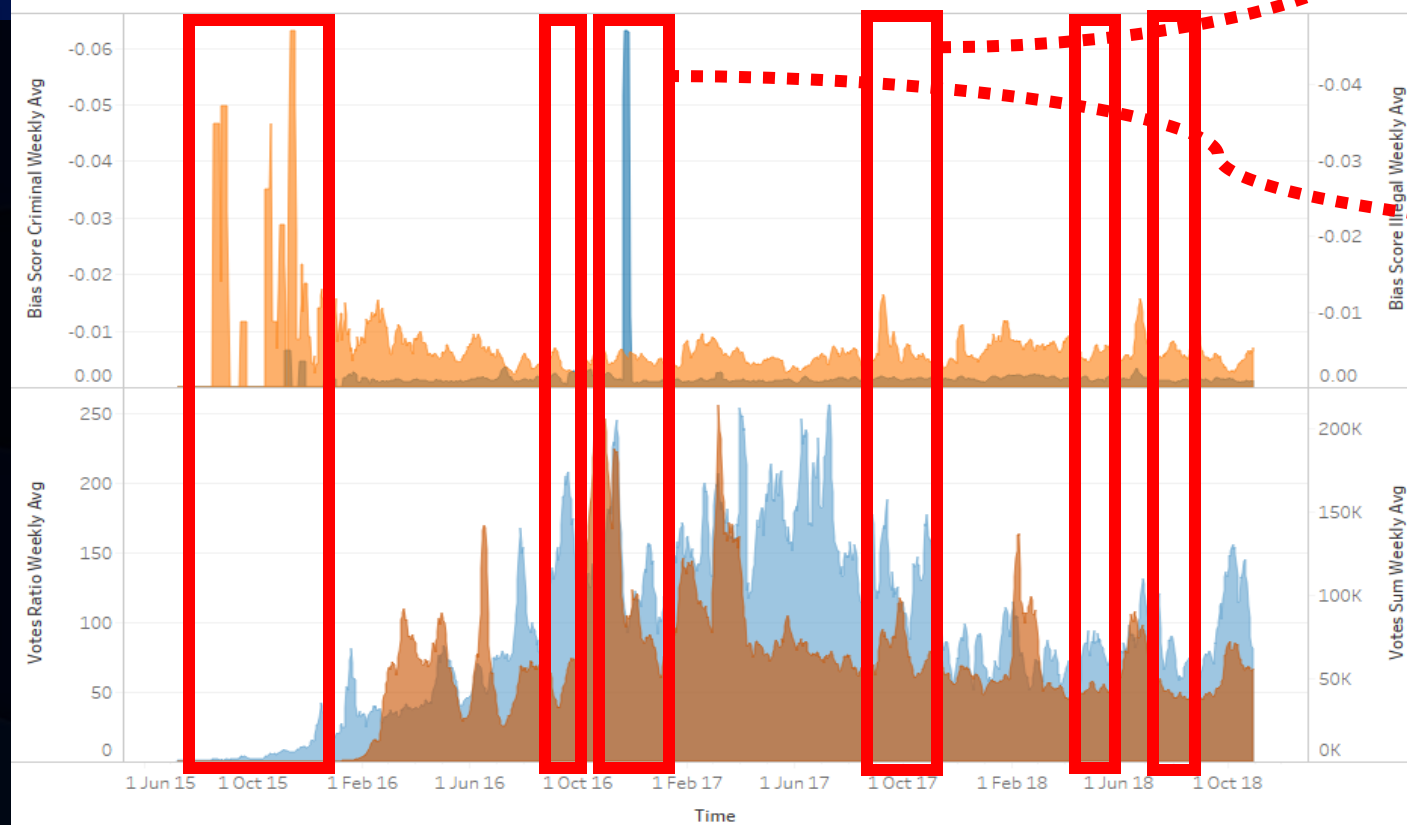U.S. Ambassador to the United Nations Nikki Haley on Monday urged the U.N. Security Council to enact "the strongest sanctions" against North Korea in response to its latest weapons test. "Enough is enough," Haley said at the emergency meeting. She said that Kim Jong Un is "begging for war" and the stakes "could not be higher." The meeting was called after North Korea conducted a test Sunday of what appears to have been a hydrogen bomb. Both President Trump and Treasury Secretary Steven Mnuchin have called for sanctioning countries that do not cut business ties with the defiant Hermit Kingdom.

Source

1/25/2021

6

# Progress

- We have further specified and improved the methodology proposed by Ferrer et al. (2020);
- We have given a detailed description of the technical implementation;
- We have validated the results and the methodology for our purpose;

- We invite you to inform our proposal for future research!
  *And to read our report, of course…*



INFOMMDI Assignment

VISUALIZING LANGUAGE BIASES OVER TIME USING REDDIT COMMENTS DATA.

September 2020

*Authors:*
Otto Mättas,
Filippo M. Libardi

*Supervisor:*
Alexandru Telea