

# COVID-19 Misinformation Detection

Caitlin Moroney

12/2/2020

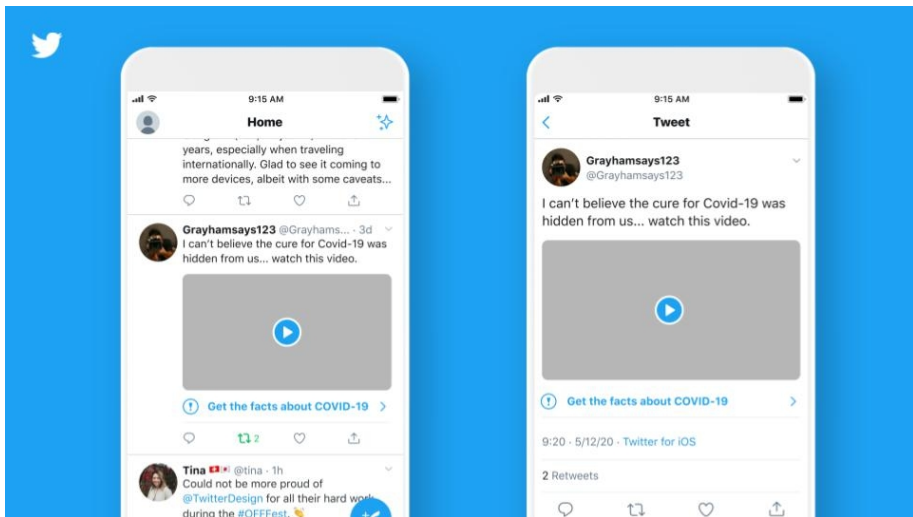


Figure 1: Twitter adds warning labels to tweets.

# Dataset from Boukouvalas et al. (2020) <sup>1</sup>

- 560 tweets, perfectly balanced classes
- sample of 282,201 users in Canada
- tweets posted between January 1 - March 13, 2020
- manually labeled as “reliable” or “unreliable”

---

<sup>1</sup>The data are available online at Dr. Boukouvalas' website.

Table 1: Misinformation rules from Boukouvalas et al. (2020)

Linguistic Feature	Example from Dataset
Hyperbolic, intensified, superlative, or emphatic language	e.g., 'blame', 'accuse', 'refuse', 'catastrophe', 'chaos', 'evil'
Greater use of punctuation and/or special characters	e.g., e.g., 'YA THINK!!?!?!', 'Can we PLEASE stop spreading the lie that Coronavirus is super super super contagious? It's not. It has a contagious rating of TWO'
Strongly emotional or subjective language	e.g., 'fight', 'danger', 'hysteria', 'panic', 'paranoia', 'laugh', 'stupidity' or other words indicating fear, surprise, alarm, anger, and so forth
Greater use of verbs of perception and/or opinion	e.g., 'hear', 'see', 'feel', 'suppose', 'perceive', 'look', 'appear', 'suggest', 'believe', 'pretend'

# Overview

- raw text
- word embeddings
  - word-word co-occurrence matrix
  - latent variable methods
- tweet embeddings
- classification
- evaluation

# Word-Context Matrix

1. I enjoy flying.
2. I like NLP.
3. I like deep learning.

The resulting counts matrix will then be:

$$X = \begin{array}{l} \begin{array}{c} I \\ like \\ enjoy \\ deep \\ learning \\ NLP \\ flying \\ . \end{array} \begin{bmatrix} 0 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 2 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{bmatrix} \end{array}$$

Figure 2: Example of a word-context matrix from Towards Data Science

# Latent Variable Methods

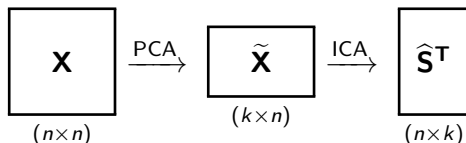


Figure 3: Truncated Singular Value Decomposition followed by Independent Component Analysis.

# Tweet Embeddings

A tweet embedding is the average of the word embeddings for the words that occur in that tweet.

$$\mathbf{v}_i = \frac{1}{T_i} \sum_{j=1}^{T_i} \mathbf{s}_j$$

Example tweet: "Covid is fake news."

$$\text{tweet} \begin{bmatrix} v_{i1} \\ \vdots \\ v_{ik} \end{bmatrix} = \frac{\text{covid} \begin{bmatrix} s_{11} \\ \vdots \\ s_{1k} \end{bmatrix} + \text{is} \begin{bmatrix} s_{21} \\ \vdots \\ s_{2k} \end{bmatrix} + \text{fake} \begin{bmatrix} s_{31} \\ \vdots \\ s_{3k} \end{bmatrix} + \text{news} \begin{bmatrix} s_{41} \\ \vdots \\ s_{4k} \end{bmatrix}}{4}$$



# LIME: Local Explainability

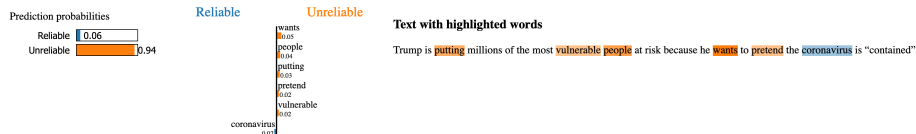


Figure 4: LIME output for unreliable tweet.

# ICA: Global Explainability

We define the importance of the  $i^{th}$  target word as follows:

$$g_i = \frac{1}{k} \sum_{j=1}^k |s_{ji}|$$

where  $k$  is the number of SVD features (and therefore the number of ICA features), and  $|s_{ji}|$  is the magnitude of the  $i^{th}$  word's importance in topic  $j$ .

## Example: Tweet 170

CNBC ADVICE NOW: Coronavirus is the flu. Wash your hands. Book a vacation. We'll look back on this and laugh.

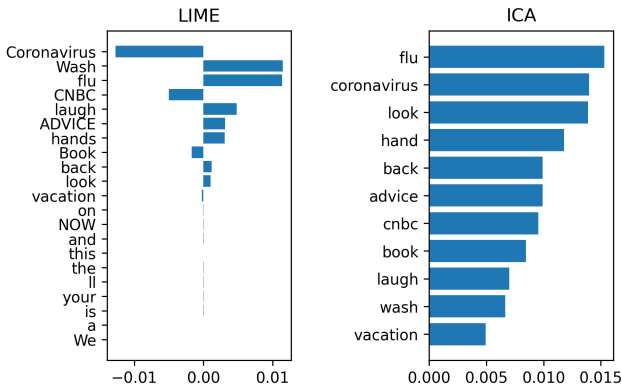


Figure 5: Comparing LIME and ICA explainability.

# Explainability metric

Goal: value that captures how “explainable” a classifier’s predictions are (with respect to human classification)

$$\frac{1}{N} \sum_{i=1}^N \frac{1}{T_i} \sum_{j=1}^{T_i} \mathbb{1}_{A_i}(w_j)$$

where there are  $T_i$  table 1 words in tweet  $i$ ,  $w_j$  is the  $j^{th}$  table 1 word in tweet  $i$ ,  $A_i$  is the set of words that the classifier associated with the unreliable class according to LIME for tweet  $i$ , and there are  $N$  tweets.

Table 2: Binary SVM performance

<b>Dimensions</b>	<b>AUC</b>	<b>Accuracy</b>	<b>F1</b>	<b>Precision</b>	<b>Recall</b>
50	0.903	0.804	0.801	0.818	0.804
100	0.911	0.796	0.793	0.817	0.796
150	0.906	0.795	0.792	0.810	0.795
200	0.901	0.800	0.798	0.815	0.800
250	0.904	0.807	0.804	0.827	0.807
500	0.908	0.789	0.785	0.814	0.789

## Results (continued)

Two experiments: (1) used strictly table 1 words, and (2) used table 1 words plus related terms. Both used stemming.

Table 3: Binary SVM explainability

Experiment	Explainability Score
1: Correctly predicted	0.593
1: Wrongly predicted	0.500
1: Aggregated	0.588
2: Correctly predicted	0.619
2: Wrongly predicted	0.407
2: Aggregated	0.608

# Future work

- local ICA explainability
- different word embeddings (e.g., BERT)
- different classifiers (e.g., neural net)
- improve explainability metric