# Interactive Dashboard for Text Label Exploration

Aisha Khatun
Mustapha Unubi Momoh

Team FAHM

UNIVERSITY OF WATERLOO

# Motivation

- Labeling is necessary for many supervised tasks. It is typically a **manual** yet **erroneous** labour.
- Automated labeling is limited by the black box nature of models. They are models specific, have static visualization, or not scalabile for complex tasks.

# Dataset

- Confusing or overlapping labels are hard to label, even for humans
- Labels in our dataset are:
  - Conspiracy
  - Controversy
  - Misconception
  - Fact
  - Fiction
  - Stereotype

# Examples

- Eskimos do have a disproportionate number of words representing snow in their languages.
    - Misconception of Controversy?
    - Requires deeper knowledge of topic to correctly identify if it is a debatable topic (controversy) or is widely believed but wrong (misconception)
- Water condensation trails ("contrails") from aircraft consist of chemical or biological agents under secret government policies
    - Conspiracy? Spread deliberately by <u>one group</u> against another?
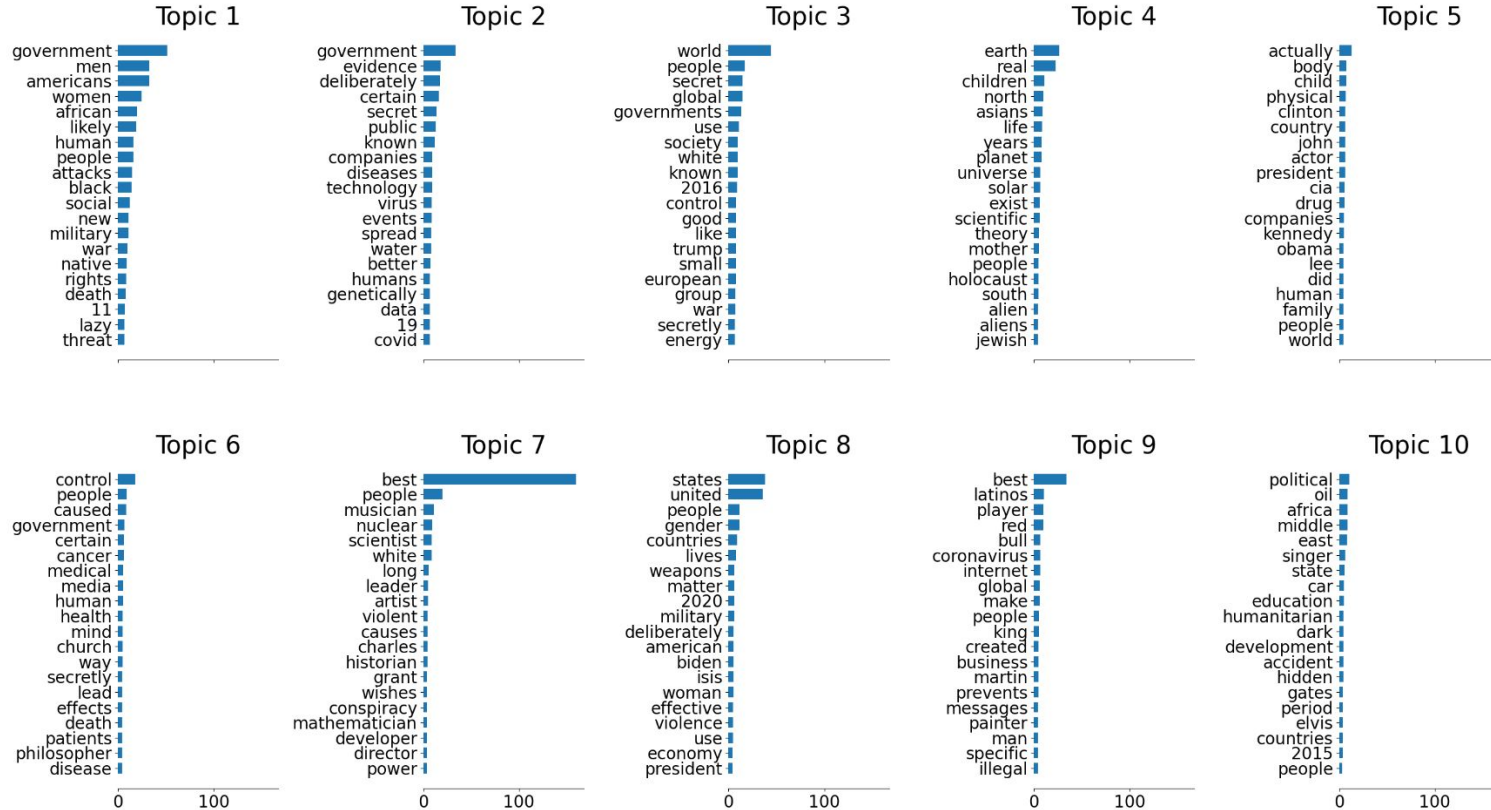    - Controversy? Just a belief held by <u>significant number</u> of people?

# Solution

- Visualize natural clusters in text labels
- Find topics in samples and see how they correlate with the existing labels
- Find semantically similar samples

# Methodology (Demo)

- Find embedding of each sample.
  - Doc2Vec, Universal Sentence Encoding, BERT encoding, …..
- Visualize natural clusters in data. Find suitable labels, or errors in existing labels.
  - PCA, T-SNE, ….
- Visualize data by topic by performing topic modeling.
  - NMF, LDA
  - Show topic words and topic word cloud for each sample
- Show top-n nearest samples for each data point
  - To help find outliers or perform semantic de-duplication

LDA

# LDA Topic 1

# LDA Topic 2