# Detecting Fake News

December 2019

Group 11

Rooney, Chu, Michael, Rachel

# Problem Definition and Motivation

## Pain Points

Algorithmic Filtering

Persistent False Rumours

Infectious Unverified Rumour

Subjective High Profile Users

We tackled **Fake news** by implementing a **machine learning** to discern fake news

01 | Exploration   02 | DIAGNOSTIC   03 | PREDICTIVE

# Relevant Data

**Common sources of fake news according to Canadians**

| | |
|---|---|
| Facebook | **68%** |
| Social media | **65%** |
| Websites | **62%** |
| YouTube | **49%** |
| TV | **45%** |

CBC NEWS

Source: Ipsos and CIGI

- **90%** Canadians fallen prey to fake news online.

- **80+%** Canadians agree or somewhat agree that search engines should be forced to remove inaccurate search results related to a person's name.

- **21,600 tweets** from troll accounts directly targeted Canadians in the 2019 Canada's federal election.

- Fake news has cost the US stock market **$39 billion** annually.

- At least **$200 million** will spend on fake news in 2020 US presidential election.

# Dataset Description

- Composed of 15,555 claims from 9 fact-checking websites
- Claims categorized by humans as false, partly true, or true
- Included claimant, date, and related articles
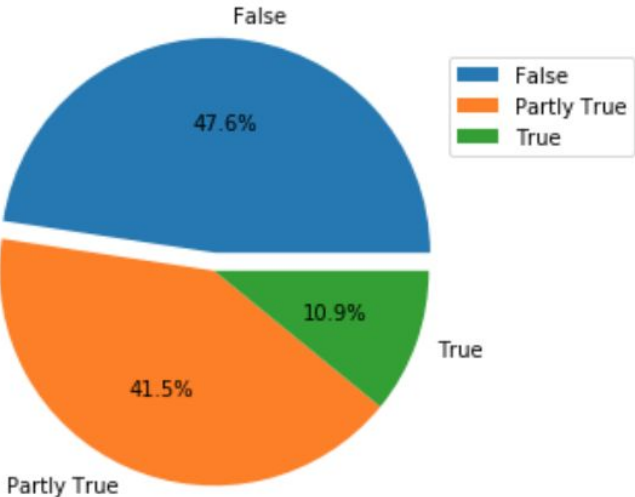- Total of 64,974 related articles composed of source and supporting articles

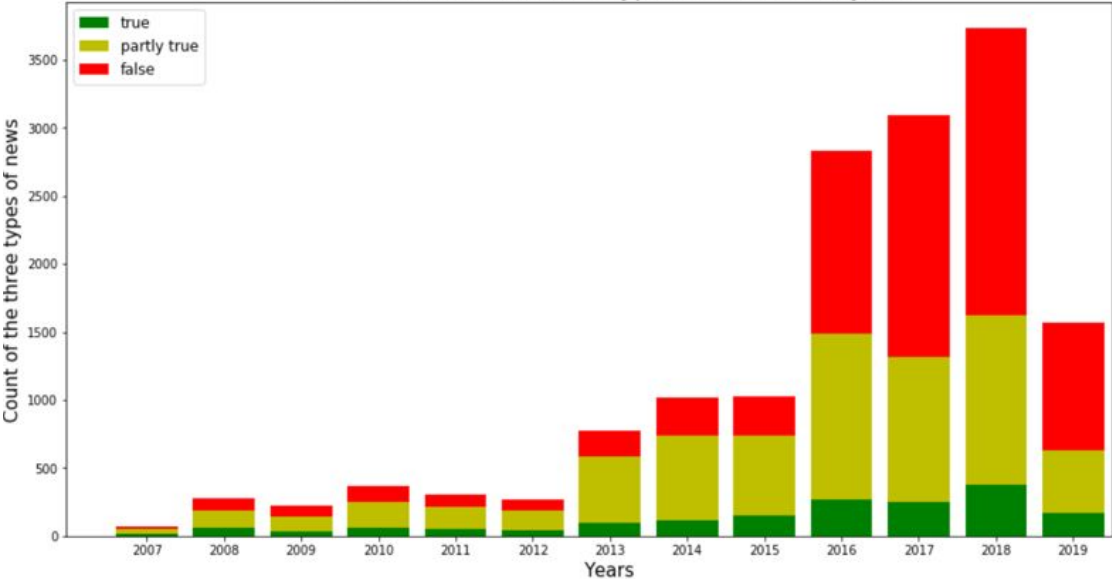| Claim | Claimant | Date | Related Articles | Label |
|---|---|---|---|---|
| When it comes to fighting terrorism, "Another thing we know that does not work, based on lots of empirical evidence, is torture." | Hillary Clinton | 2016-03-22 | Supporting Article<br>Source Article<br>Supporting Article<br>Supporting Article<br>Source Article | True |

DETECTING FAKE NEWS

# Exploration

- True news take a small portion
- Nearly half of the dataset is fake news
- News before 2013 contribute a small portion
- Ratio of fake news increased significantly compared to partly true and true news
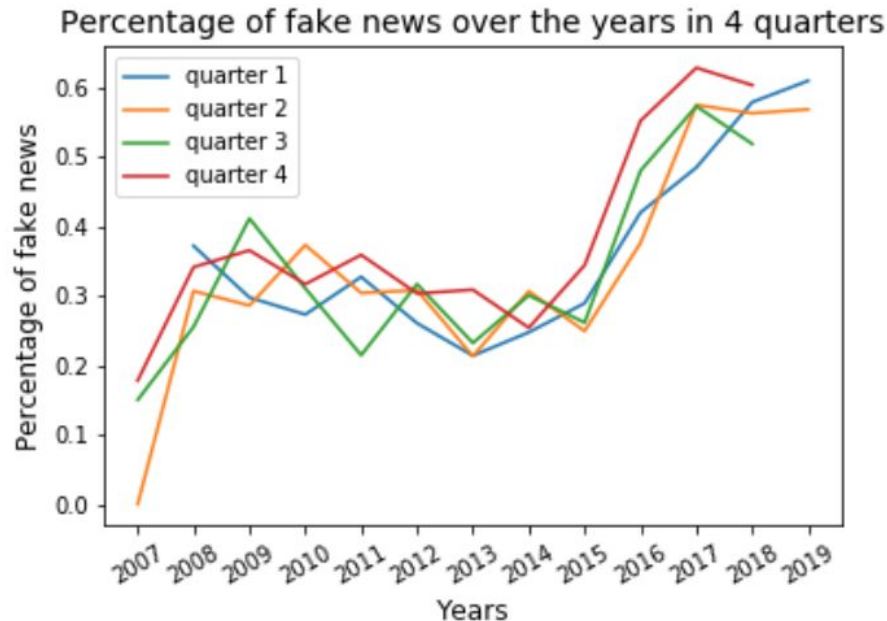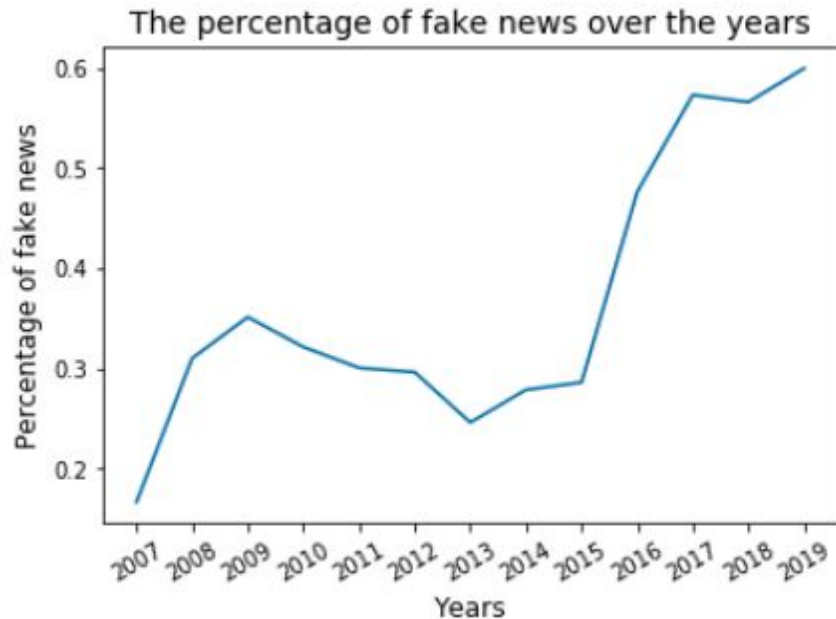
Distribution of the news



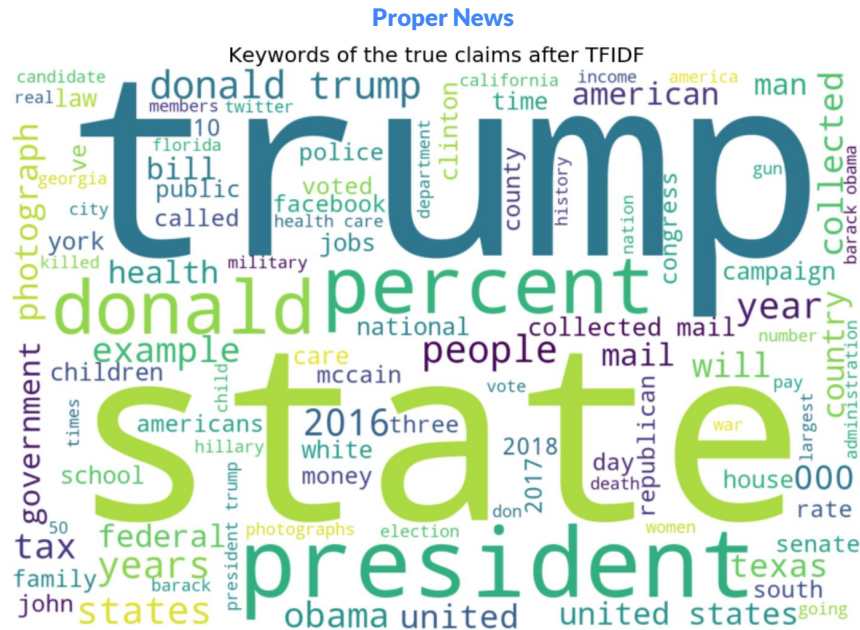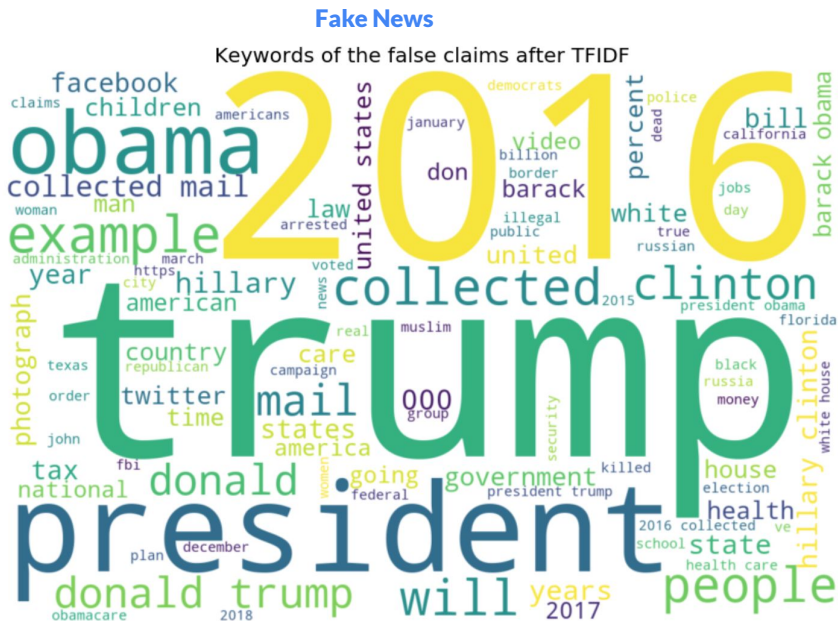The amount of the three types of news over years

# Exploration

- Overall increasing trend of fake news
- Ratio of fake news increased significantly since 2015
- Each quarter has a similar trend
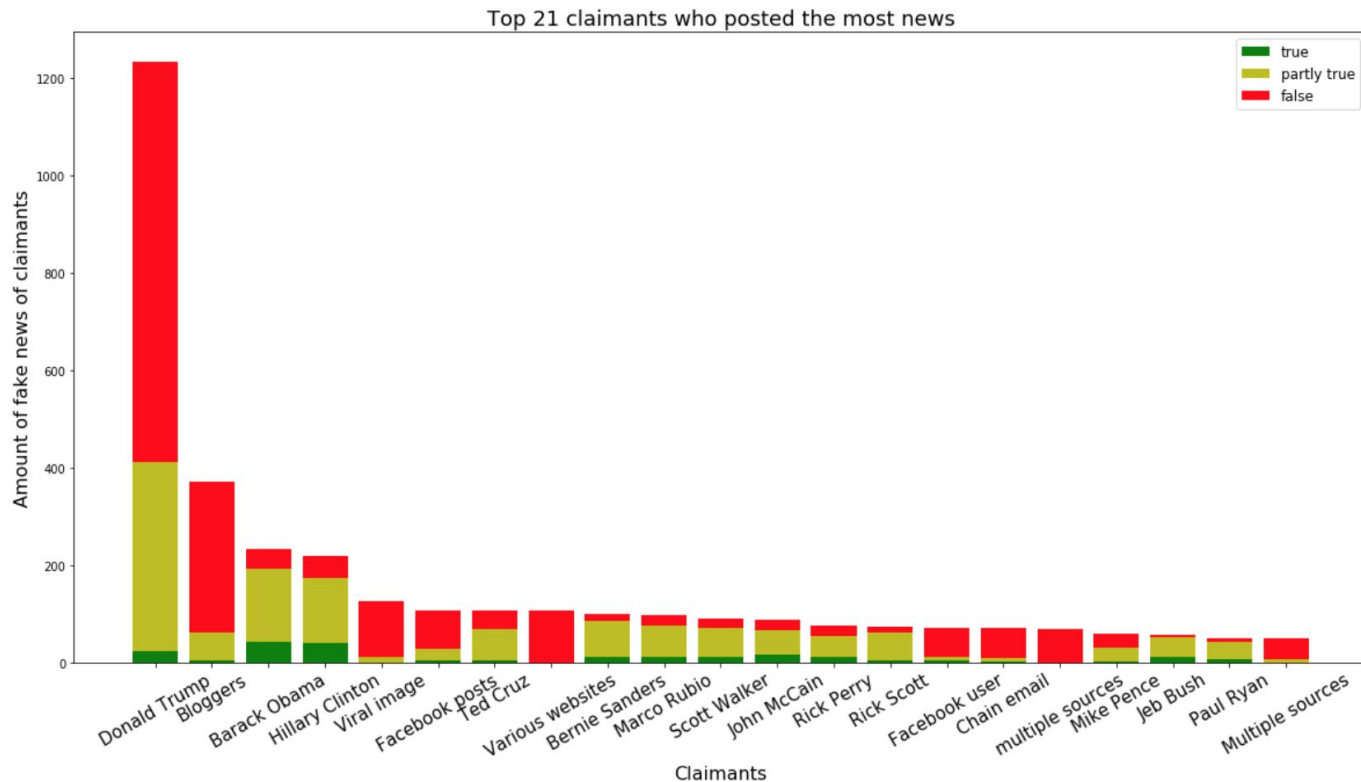- The fourth quarter has higher percentage of fake news



The percentage of fake news over the years



Percentage of fake news over the years in 4 quarters

# Exploration

- Frequent words are political
- 'Trump' frequently occurred in both fake and true claims



**Fake News**

Keywords of the false claims after TFIDF
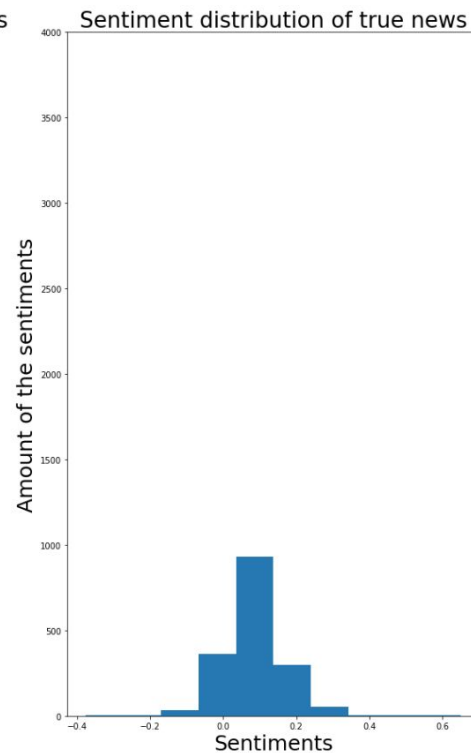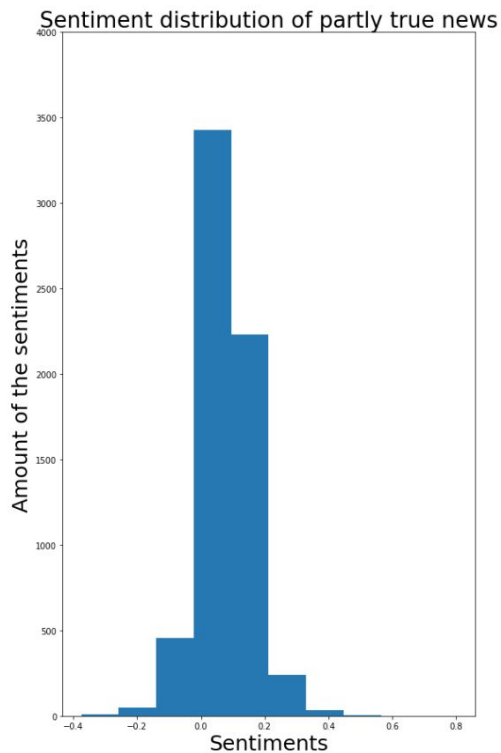
**Proper News**

Keywords of the true claims after TFIDF

# Exploration

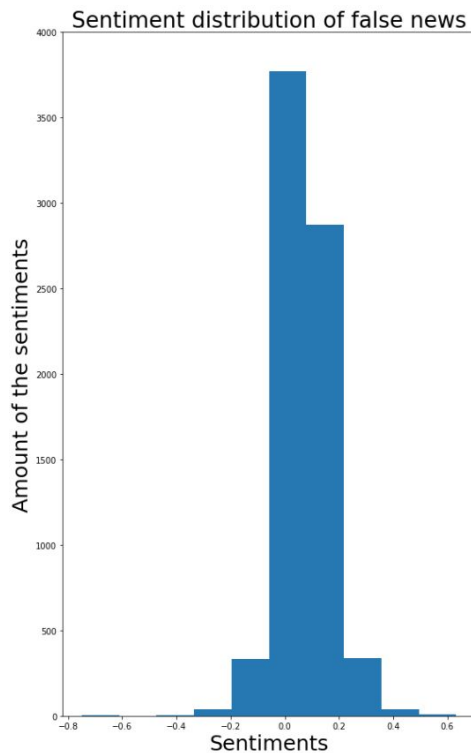- Social media is the main source of fake news
- Donald Trump posted more than 1200 news, more than half of them are false



Top 21 claimants who posted the most news
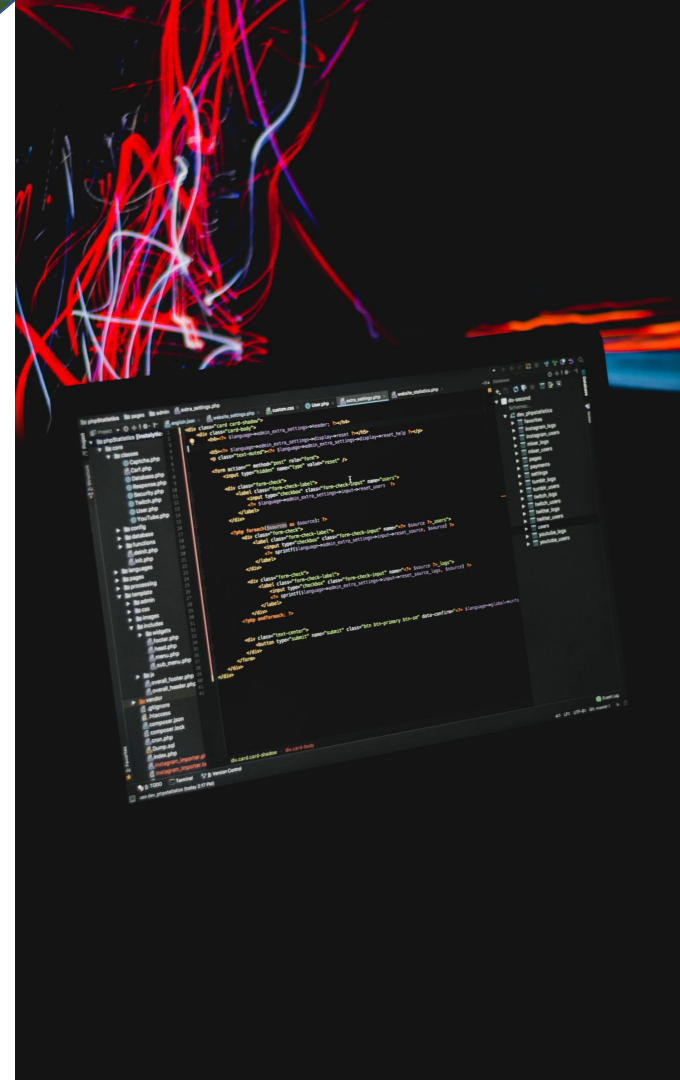
# Exploration

- Most of the news are positive
- True news are more positive than fake and partly true news



Sentiment distribution of false news · Sentiment distribution of partly true news · Sentiment distribution of true news

# Machine Learning Implementation

- Calculated how similar the claim is to each sentence in the supporting articles which indicates if the claim is well supported
- Analyzed the sentiment of the claims
- Also considered the claimant, month, and year
- Used the reliable XGBoost open source machine learning model that was used by several Kaggle competition winners

# Results

- Testing on a **30%** subset of the original data:
    - Achieved **59%** accuracy in classifying claims as false, partly true, or true
    - Detected **80%** of false news with **69%** accuracy
- Obtained F1 score of **36%** in the DataCup competition
- Possibility to flag misleading posts and news articles
- Automation of fact checking enables immediate results and processing huge volumes of news

01 | DESCRIPTIVE   02 | DIAGNOSTIC   03 | PREDICTIVE

# DECISION MAKING

## Text Mining

### Detect fake news
Through machine learning

### Flag Fake news

Flag likely Fake news as fake and draw attention to fact checkers

### Fact Check

Manual Fact check to validate and improve the machine learning results

### Continuous Improvement

Continuously improve the model by incorporating new labels from fact checkers