# An average citizen's guide to fight Fake News

By T. Bornet dit Vorgeat, T. Mortiniera & P. Wagner
Supervised by Léonore Guillain
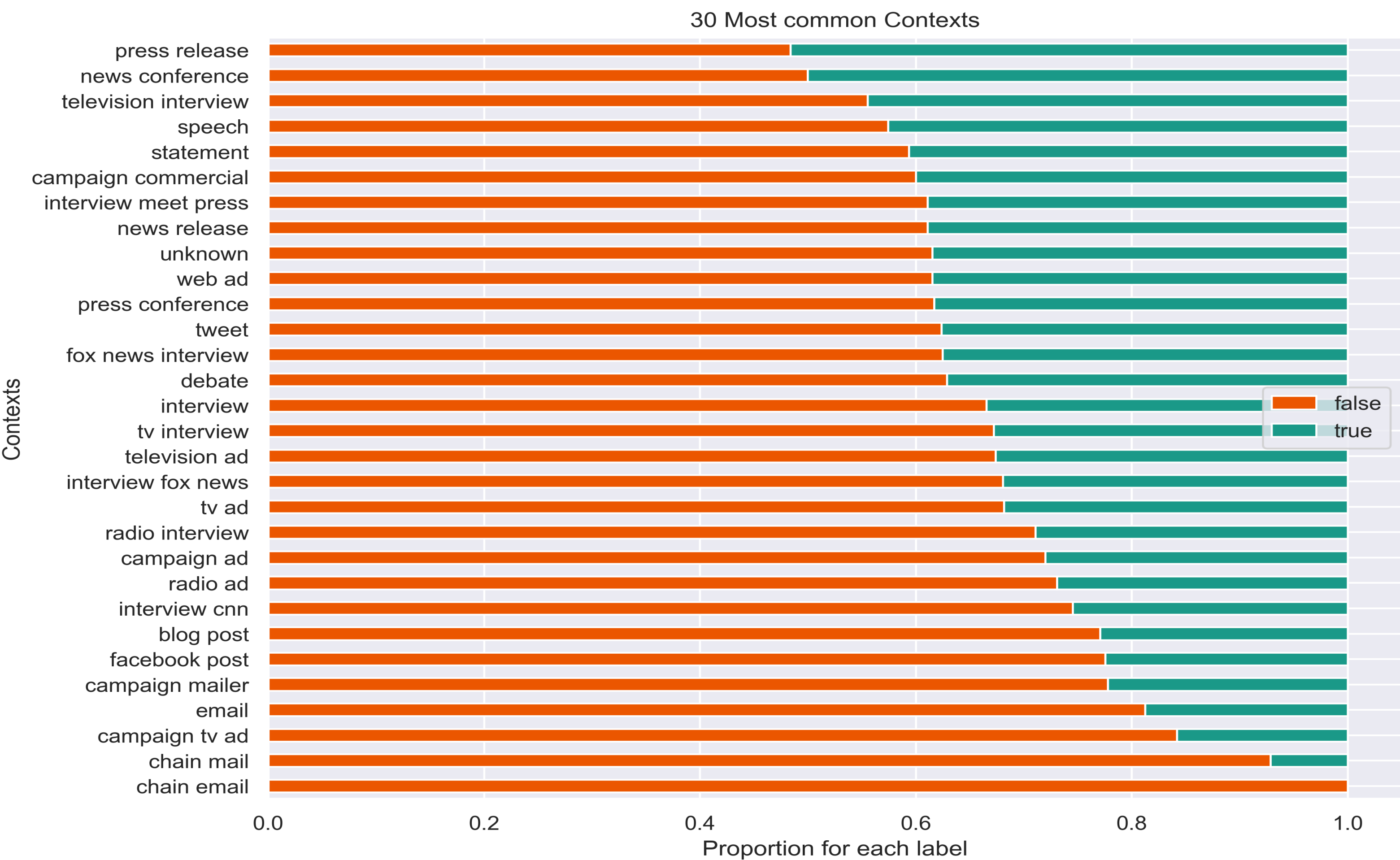Applied Data Analysis (CS-401)

## Introduction

With the exponential spread of data accessible to everyone, the proliferation of misleading information has made it a challenge to rightfully identify reliable news sources. This has led to an increase in the demand for automated tools to detect false news. In this work, we proposed **a guide for the average citizens** to make them more critical about information, and help them make the distinction between the safe and unreliable topics and sources of information.

To do this, we answer the following questions:
- **Can we classify the reliability of the context?**
- **Are there suspicious topics where extra caution is advised?**
- **Can we determine a vocabulary of fake news?**
- **Are there geographical relations to fake news?**



Word cloud of 10 most common topics in Fake News



Word cloud of 10 most common topics in Real News

Most important topics in Real news refer to topics delivered from "serious" press release. Science or economy newspapers, statements and figures on jobs, transportation and safety concerns, are platforms in which the experts making statements are quoted.

However regarding the Fake news, most of the information relates to conspiracy theory. It seems, that hype and fear are used to promote differences and intolerance.
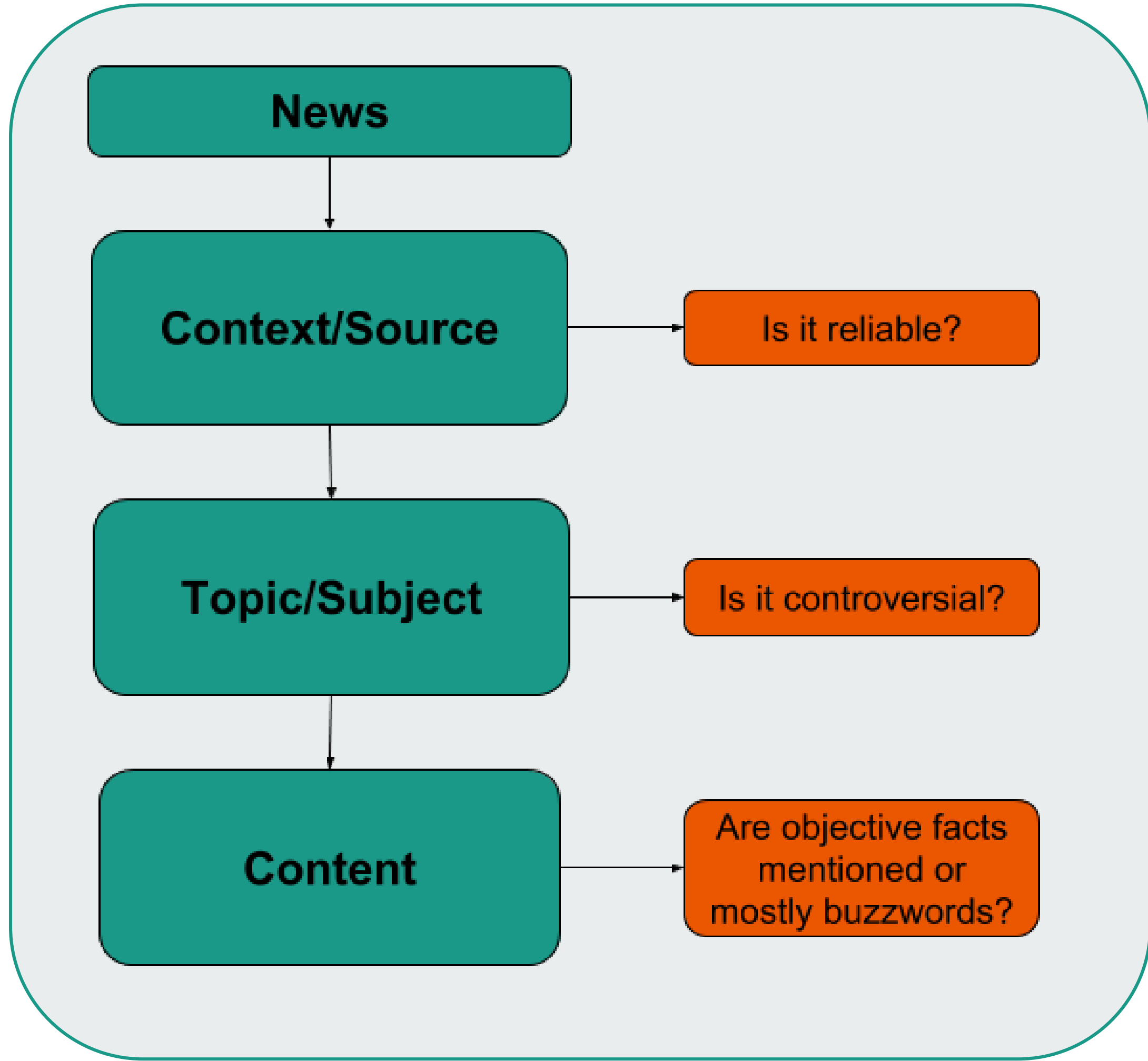


30 Most common Contexts

We observe that a baseline of truth, the situation where a normal level of distrust is needed, emerges as the central line where the proportion of true statements vary between 30% and 40%. The safest contexts, where 40% to 50% of true news are given, is a group needing the less amount of distrust of all. It contains information coming from official channels like press or news releases, television interviews (mostly from "pundits"). This is the complete opposite of the last case where the context is mostly online in e-mails like ads, chain mails, posts on Facebook or blogs. This falseness haven is made up of contexts where less than 30% of statements are true.

**Dataset** from Politifact.com made up of 12.8K manually labelled statements. We worked on a subset consisting of the most true and falsely labelled data. Hence we are left with 6415 statements with their metadata.

## Vocabulary

The following tables describe the most important features that help the model to predict if a statement is false or true. These features follow the intuition that controversial subjects are determinant in the prediction of fake news. Another intuition is confirmed, indeed the words which are determinant in the prediction of the true news are mostly objective terms used to describe facts. Surely they are made by objective statements that cannot be attacked on their truthfulness.

| Feature | Weight | Feature | Weight |
|---------|--------|---------|--------|
| obamacare | 2.138006 | georgia | -2.023354 |
| wisconsin | 1.678782 | months | -1.701455 |
| medicare | 1.632842 | child | -1.576920 |
| making | 1.472672 | day | -1.539615 |
| muslim | 1.331842 | half | -1.467554 |
| care law | 1.323680 | top | -1.453722 |
| raise | 1.295063 | terms | -1.415189 |
| illegal | 1.223156 | three | -1.410444 |
| work | 1.202895 | called | -1.408316 |
| obamas | 1.202583 | mccain | -1.375338 |

Word feature importance determining Fake and True news obtained using a Logit model



## Conclusion:

Multiple questions were answered during our investigations. We concluded that geographical origin does not affect the truthfulness of statements, while the contexts do. Indeed we found more or less three different kinds of contexts; with varying levels of truthfulness. Finally, it seems that controversy tends to bring misinformation. We also verified our assumptions with machine learning techniques, such as Logit Model which gave us a **F1-Score** of **71%** while its feature importance aligned with our assumptions.