

Initial Project Overview

Overview of project content and milestones

A recent phenomenon that has made its way into public perception is that of ‘fake news’ - misleading and deceptive information that can sway political opinion, manufacture outrage and split society. The aim of this project is to create and evaluate a method of determining if the headline of a news article appropriately reflects its content, which could offer a way to track and analyse the dissemination of misinformation.

The project will have several milestones. Firstly, a literature review will be undertaken that both analyses existing solutions, as well as exploring the problem space more generally. Then a prototype solution will be created, after which it will be tested and refined. Finally, the results and efficacy of the solution will be analysed and conclusions drawn.

The main deliverable(s)

The primary deliverable of this project will be an algorithm that is able to determine to what percentage a headline matches the content of an article.

Alongside this, a dataset of articles may be created that could be used by others to further develop incongruity detection.

The target audience for the deliverable(s)

This project could be of interest to a wide range of sectors - people in the field of sentiment analysis and natural language processing, those who produce and write articles, and people who want to determine the trustworthiness of the content they’re consuming could all benefit from this project.

The work to be undertaken

Before I start to work on the main bulk of the project, I’ll begin by creating a plan of action, in the form of a kanban board and a Gantt chart. Then, the literature will need to be reviewed and possible techniques explored.

Data will then need to be collected, and depending on the requirements of the technique I settle on it’ll need to be labelled. Alternatively, an existing dataset could be sourced, although it will be tricky to find one that matches the problem space perfectly.

Once a dataset is in place, the algorithm will need to be created that can classify each record and determine a percentage incongruity between the headline and body of an article.

The results of this classifier will need to be tested and analysed. Ideally the dataset will have information about the articles (e.g. date of publication, publisher etc.) and trends could be graphed to show correlations between this meta information and the incongruity.

Additional information/knowledge required

I’m currently unsure how to approach the algorithm design and implementation of this project. I have a very rudimentary understanding of neural networks and sentiment analysis, so research will need to be conducted into these areas, and a decision made whether to use one or the other, a combination of both or something else entirely.

Information sources that provide a context for the project

- *Detecting Incongruity Between News Headline and Body Text via a Deep Hierarchical Encoder, 2018* Uses two neural networks with hierarchical structures to determine how incongruent headlines and bodies in articles are
- *The effects of subtle misinformation in news headlines, 2014* Investigates how misinformation can be delivered via news headlines
- *Media-generated Shortcuts: Do Newspaper Headlines Present Another Roadblock for Low-information Rationality?, 2007* A manual content analysis that shows “considerable difference” between articles and headlines

The importance of the project

Since the idea of news was created, people in power have used a range of techniques to manipulate it to their advantage. As headlines are the most prominent aspect of any article, they are a prime target for alteration and miscommunication. By building a tool that can programmatically detect any incongruity, misleading and potentially manipulative headlines can be identified.

Additionally, the use-case of the project could be built upon by further research, and techniques used could apply to research papers or similar publications.

The key challenge(s) to be overcome

The first challenge is that of creating or obtaining a quality dataset. In order to have validity, the project will need a large collection of news articles and their headlines, with each being labelled with their incongruence. While it would be fairly trivial to automate the collection of news articles, labelling them would have to be done manually, which would be a very subjective and time-consuming process.

Another challenge is my lack of experience around the potential technologies I could use to approach the problem - building classifiers, machine learning and natural language processing are all outside of my current skillset. While I expect to do a lot of learning around this area during the literature review and background research, I anticipate that the development aspect of the project may take a while as I get to grips with some paradigms and techniques.