

Ben Kaufmann

Daniel Semenko

CptS 315

21 March 2022

Course Project Proposal

1. Data Mining Task

For our project, we are going to do the Toxic Comment Classification project, as recommended in the ideas text file. Originally the competition was designed by Jigsaw and Google to help improve online tools for detecting negative online behaviors. We were motivated to choose this project for a couple of reasons. We both play video games and as such are often exposed to these types of comments, whether they are verbal or through text chats. In addition, though the project may prove to be difficult, we believe that we will be able to successfully complete it.

2. Dataset

The dataset, provided through Kaggle.com, is comprised of comments from Wikipedia's talk page edits, which have been labeled by human raters for toxic behavior.

3. Methodology

Our problem consists of reading through text, and labeling the text based on certain keywords or phrases within the text. This is solved through sentiment classification, though ours may be adjusted slightly to work with the multiple types of toxicity we look for. Since the only data provided is comment strings, we will use the string lengths and common words to determine ratings for each of the types of

toxicity. Then after analyzing the dataset, we will create a model which then can be used to predict the toxicity ratings of the test data.

4. Final Product

The physical output of our project will be a text file which will contain the ID's of comments found within the test data, and our predicted probabilities of each of the six possible types of comment toxicity.

One of the benefits of working with data used in a competition, is that there already is a scoring/evaluation system set in place to measure the effectiveness of any given solution to the problem. Therefore, we can use the evaluation leaderboard and determine a score (from 0.0 to 1.0), and see how effective our algorithm turned out to be. If we can get above a 0.9 (based on the evaluations used in the competition), then that will be a success for us.

In addition to learning about Sentiment Analysis in general, I believe the most impactful thing we will learn to do is how to create a model out of data, and use and apply the model to make predictions.