Ezana N. Beyenne

MSDS 453, Section 57 2020

**Week 6: Proposal for Fourth Research/Programming Assignment**

Initially, I was building a corpus on the criticisms of Facebook, and wanted to use that model to

see if other companies followed a similar path and how they ended up rebuilding their public

image. As I did more research, I found two interesting issues that kept coming up. First, the

psychological effects of social media, especially when one is the target of toxic comments.

Secondly, Fake news (especially being posted and disseminated by bots) kept coming up in

association with Facebook. The two subtopics of my research were Toxic comments and Fake

news, which are sometimes related to one another because of the vitriol that currently exists on

social media when the topic is related to politics. Luckily, I could find the datasets on Kaggle,

and I am having a hard time choosing between the two. The two data sets are:

1. The Getting Real about Fake News**:** This is an interesting topic because it would be

great to find a way to detect fake news and stop it from going viral. There are several posts on

how to detect fake news, but there is no standard way to detect them, since they seem to lack any

consistency. The Kaggle dataset contains the data from 244 websites and has 12,999 posts that

were tagged by the BS Detector Chrome extension. There are a lot of websites that provide

guidelines, but as the fight against Fake news intensifies, the offenders keep evolving in more

sophisticated cat and mouse game.

2. The Toxic Comment Classification Challenge: The 1$^{st}$ topic that I was going to choose

was on the criticisms of Facebook, and one them was the psychological effect that online

comments had on people. So, I saw this Kaggle dataset, which attempted to identify and classify

toxic online comments. In this challenge, we will attempt to detect different types of toxicity like

threats, obscenity, insults and identity-based hate. It would be interesting to develop a model that can attempt to eliminate toxic comments from being posted.

There is a lot of literature and research being conducted on either topic, and I am hoping to let the reading help me determine which topic to pick. Any input or guidance on which way I should go would be greatly appreciated.