# Capstone Proposal Template

## Fake News Detector

### Business Understanding
The problem I'm aiming to solve is to limit the spread of misinformation. This topic is somewhat personal to me, as this was a heated discussion topic last time I went home to visit my family. While this problem isn't necessarily specific to any industry, the general public would find this beneficial, as it would allow them to identity misinformation in news articles quickly.

Pre-existing projects/research papers I've explored are below:
https://data-flair.training/blogs/advanced-python-project-detecting-fake-news/
https://journals.sagepub.com/doi/full/10.1177/2053951719843310

### Data Understanding
Data Sources:
https://www.uvic.ca/engineering/ece/isot/datasets/fake-news/index.php
https://data-flair.training/blogs/advanced-python-project-detecting-fake-news/
https://www.kaggle.com/c/fake-news/overview

Noone has worked on this specific dataset, as I combined the dataset from 3 separate sources. The features that will be used are described clearly (title of article, text of article, Real/Fake label).

### Data Preparation
The data is stored as a .csv, where the variables are all strings. Pre-processing steps/cleaning challenges will include removing stop-words and punctuation, and making bigrams. My current dataframe is ~72,000 rows. I am planning on creating a word cloud ot show the important words/phrases

### Modeling/Tools/Methodologies
The target variable is whether a news article is real or fake (classification). The baseline model would just be a model that predicts the majority class every single time. Some techniques I'll be using are the TFIDF Vectorizer, PassiveAggressiveClassifier, and looking at Confusion Matrices,

### Evaluation
The metric I'll be used to determine success will primarily be accuracy. The MVP involves data cleaning and providing data insights on which words/phrases strongly indicate if an article is real or fake, and basic modeling. My level-up stretch goals are to build a model with high accuracy.

### Deployment
Creating a plug-in that users can use that will scan to see if an article is real or fake.