

GROUP 4 FINAL PROJECT PROPOSAL

Group Members:

- Alexis Kaldany
- Ricardo Diaz
- Sagar Tripathi

Problem:

Our project examines COVID articles and tries to classify them, there are large quantities of covid articles which can be overwhelming if a person is trying to look for a specific article in a category (business, tech, etc). We are looking to simplify the search for articles by categories plus comparing accuracy from two NLP methods to evaluate which works better to classify articles.

Dataset:

[COVID News Articles \(2020 - 2022\) | Kaggle](#)

The dataset encapsulates approximately half a million news articles collected over a period of 2 years during the Coronavirus pandemic onset and surge. It contains several articles with the headline, content, and their respective category.

NLP Method:

There are two main NLP methods we are using.

1. Classical Classification Model
2. Seq2seq Model

Packages:

Pandas, Torch, Transformers (DistilBert), HuggingFace.

We need pandas to read, and manipulate the dataset, torch to use neural networks for our models, and HuggingFace to use summarization for the Seq2seq Model

Performance:

We will evaluate and compare the results based on their accuracy and their confusion matrix.

Time Frame for the Project:

NLP Final Project

Project Schedule

[illegible]