

## News Recommendation System

## 1. Research --

### Timeline followed -

- 14 Sep - Group Discussion and jotting down the main points for the dataset.
- 15 Sep - Dataset formation- Preprocessing.
- 16 Sep - EDA on the data collected
- 17 Sep - Final Documentation

### Articles read -

<https://www.datacamp.com/community/tutorials/recommender-systems-python>,  
<https://blog.insightdatascience.com/news4u-recommend-stories-based-on-collaborative-reader-behavior-9b049b6724c4>,  
<https://www.offerzen.com/blog/how-to-build-a-content-based-recommender-system-for-your-product>,  
<https://github.com/vivekpabani/News-Recommendation-System>

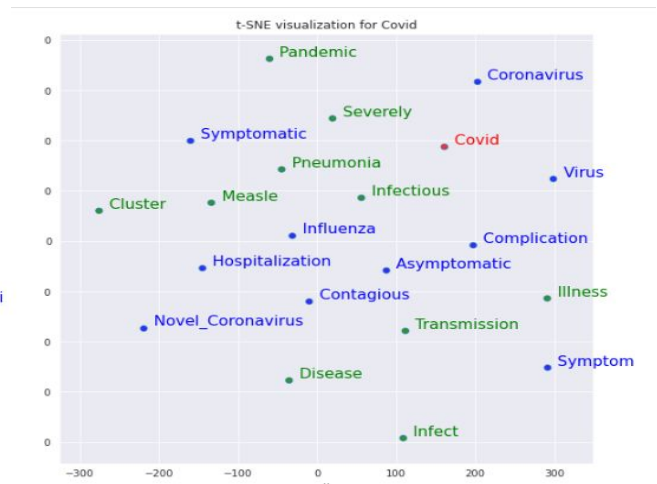
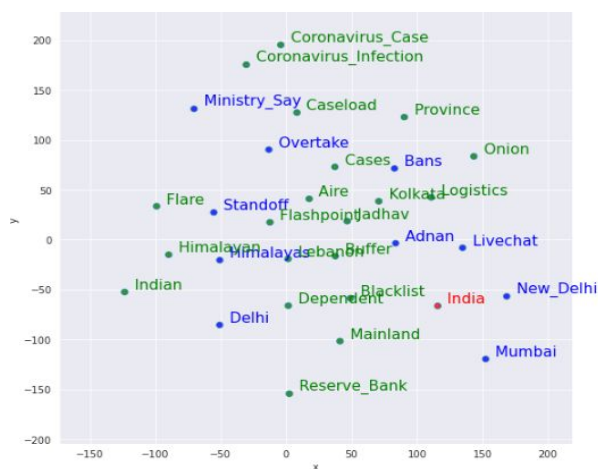
**Data sources** - We created our own dataset. As there was no dataset available which was up to date. We used <https://newsapi.org/> , to extract links of various articles and used “**newspaper3k**” for extracting **title,dates,author,summary,keywords etc.** We are also trying <https://pypi.org/project/news-fetch/> to extract data.

## 2. Exploratory Data Analysis --

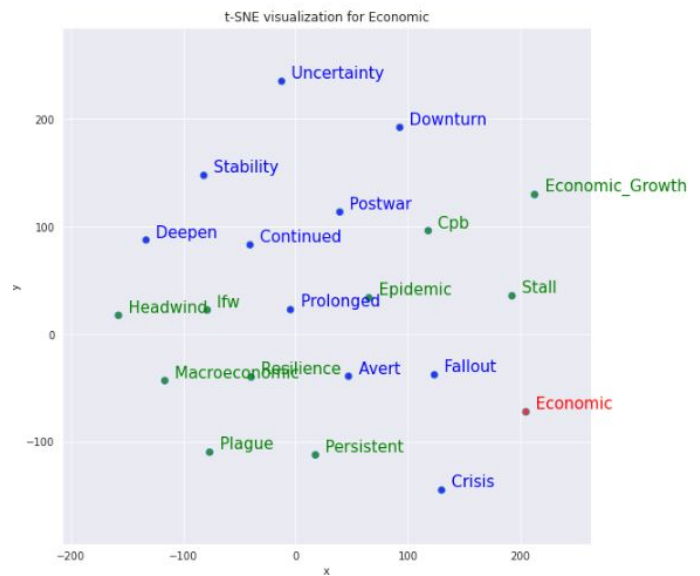
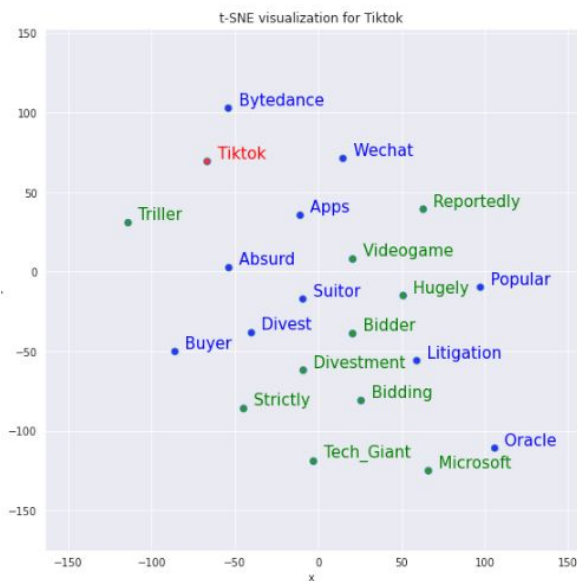
## What do you look for -- correlated data, similar phrases, and

**Conclusions drawn** -- rich corpus , meaningful, bigrams, trigrams, similar tags .

To make the visualizations more relevant, we will look at the relationships between a query word (in **red**), its most similar words in the model (in **blue**), and other words from the vocabulary (in **green**).



These are 2 examples to show how rich the corpus is that if we type just one word like “India” or “Covid” it is able to show similar/ correlated data. This is required as these words will be used as keywords to define a document and later used for whichever recommender system architecture we choose for the current problem statement.



The above data shows that the dataset has also covered variety of news articles.

