Taylor Herb
INFSCI2480 Final Project
Spring 2021

# Content Based News Recommender System
https://github.com/t-4-h/infsci2480finalProject

## I. Introduction

News websites and blogs are a highly popular aspect of the internet we use today. With a decrease in printed newspapers, online news websites have become an increased source for people to access local, domestic and international news. Since these news sources are being used by people so frequently, they would likely benefit from including a recommender system to personalize the user's experience. For example, if a user reads an article about the stock market in Europe, they may want to read additional articles on the topic. This could mean they want to see more articles about the stock market in general, Europe, the stock market in China, or maybe they aren't sure where to look next. In this situation, a recommender system would be useful to show them other relevant articles on the site that they may be interested in reading.

In this project, I will begin the process of designing a news recommender system by first using content-based recommendation methods. Since there is no user data on this dataset yet, the recommender system must rely on the information about the articles themselves rather than the user. The goal of the current system is to allow the user to search for recommendations based on the article's text body or the article's tagged keywords.

## II. Dataset and Pre-processing

The original dataset was pre-processed in order to remove unnecessary columns, missing data and duplicate entries. The final dataset included the five columns to be used in the recommender system (id, title, text, keyword, link) and 1720 rows that did not include any missing data or duplicates.

Fig 1: Original data (2190 rows, 9 columns)



Fig 2: Final dataset after initial pre-processing (1720 rows, 5 columns)

Additionally, text cleaning methods were used to clean the article text and keywords data so that the text could be implemented efficiently within the recommender system. The two text columns were cleaned by converting the text to lowercase, removing special characters, and removing stopwords. In combination with cleaning, the text from the keywords column were stored in a string array.

III. Technologies

For this project I have primarily used Python and Jupyter lab. Since the database is small and static, the .csv database is stored within the Jupyter notebook. The interactive portion works using the

ipywidgets interact library. A pseudo front-end is created using the Voila server extension, which creates a simple web application that allows the user to enter an article title and receive the top 20 recommendations based on either the article text or article keywords. The user can view the recommended article titles and then click links to visit the article online.

IV. Recommender Design

This recommender will use content based filtering methods to give recommendations based on the articles text or the articles keywords. The reason for doing both is to understand which system performs more accurately in general or which system may perform best under certain conditions. Logically, it would seem that keywords may be helpful when a user is interested in a broad topic while text may be more helpful for more specific topics (including multiple terms). Therefore, I decided it would be useful to build and test both in order to develop a more well-rounded system.

To build both recommender types (text-based and keyword-based), Term Frequency-Inverse Document Frequency (TF-IDF) was used on the text to evaluate term importance to the document within the corpus. To accomplish this, I used the TfidfVectorizer method from the sklearn.feature library. For the keyword-based recommender an ngram range of 1 to 2 was used (because keywords typically exist of a maximum of 2 terms) with a min_df of 0 (to include all keywords regardless of document frequency). For the text-based recommender, an ngram range of 1 to 3 was used to allow for longer combinations of terms within the article. Although stopwords were removed, a max_df of 0.8 was also selected to ignore terms that occur too often in the corpus because they are unlikely useful for recommendations (i.e. words similar to stop words that do not add meaning). Additionally, multiple min_df values were tested (0.0, 0.01, 0.02, 0.1) which yielded different matrix dimensions based on the amount of terms used [(1720, 1023226), (1720, 4705), (1720, 2546), (1720, 321)]. The results from min_df 0.1 were not tested due to such a small number of terms. However, the other 3 matrices all yielded seemingly accurate and relevant results in the recommendation system. It was impossible to know which performed best without user

feedback. So, min_df = 0 was ultimately chosen as it contains the most terms but user feedback should be included to understand which parameters perform best.

```
[4]:  tf = TfidfVectorizer(analyzer='word',ngram_range=(1, 3),max_df=0.8, min_df=0, stop_words='english')
```

```
[5]:  matrix = tf.fit_transform(df['cleaned'])
      ##show matrix dimensions
      matrix.shape
```

```
[5]:  (1720, 1023226)
```

min document freq = 0.01

```
[6]:  tf2 = TfidfVectorizer(analyzer='word',ngram_range=(1, 3),max_df=0.8, min_df=0.01, stop_words='english')
```

```
[7]:  matrix2 = tf2.fit_transform(df['cleaned'])
      ##show matrix dimensions
      matrix2.shape
```

```
[7]:  (1720, 4705)
```

min document freq = 0.02

```
[8]:  tf3 = TfidfVectorizer(analyzer='word',ngram_range=(1, 3),max_df=0.8, min_df=0.02, stop_words='english')
```

```
[9]:  matrix3 = tf3.fit_transform(df['cleaned'])
      ##show matrix dimensions
      matrix3.shape
```

```
[9]:  (1720, 2546)
```

min document freq = 0.1

```
[10]: tf4 = TfidfVectorizer(analyzer='word',ngram_range=(1, 3),max_df=0.8, min_df=0.1, stop_words='english')
```

```
[11]: matrix4 = tf4.fit_transform(df['cleaned'])
      ##show matrix dimensions
      matrix4.shape
```

```
[11]: (1720, 321)
```

Fig. 3: Matrix dimensions with different min_df values for text-based recommender

Once the article text or keywords were encoded into tf-idf matrices, I used cosine similarity (sklearn cosine_similarity method) to obtain a similarity measure based on proximity between the two matrices. Cosine similarity is useful here because it is a measure of distance and orientation rather than magnitude (i.e. Euclidean distance).

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|\|\mathbf{B}\|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2}\sqrt{\sum_{i=1}^{n} B_i^2}},$$

Fig. 4: Cosine similarity equation showing similarity ranging from -1 (opposite) to 1 (equal) where 0 indicates no relationship.

The recommender function was built to create indices for titles then find recommendations based on the cosine similarity (cosSim). This function returns the top 20 most relevant items and shows the user the title of the article as well as the link to the article.

```python
def recommender(title):
    ##create indices for titles/links
    titles = df['title']
    link = df['link']
    indices = pd.Series(df.index, index=df['title'])

    ##find recommendations using cosine similarity
    idx = indices[title]
    similarity = list(enumerate(cosSim[idx]))
    similarity = sorted(similarity, key=lambda x: x[1], reverse=True)
    similarity = similarity[1:21]
    news_indices = [i[0] for i in similarity]
    newsTitle = titles.iloc[news_indices]
    newsLink = link.iloc[news_indices]
    result = newsTitle, newsLink
    return result
```

Fig. 5: Recommender function based on cosine similarity

```
[1620                             Bank of England policymakers warn UK economy facing bigger risks
 1619                              UK economy might take years to recover from COVID hit-BoE's Vlieghe
 1717                                             Bank of England gears up for next stimulus push
 1526                         Britons a bit more upbeat on finances but worried about economy, GfK says
 1654                                 UK economy extends recovery from COVID crash, growth seen fading
 305               UK's Sunak considers sweeping tax hikes to plug COVID-19 hole, newspapers say
 1509                                      UK economy rebounding for now, as public borrowing mounts
 1671             Exclusive: BOJ to offer brighter view on economy as COVID crisis eases: sources
 1426                             Britain plans hiring spree to harness big data in pandemic recovery
 1713                             British business calls for green recovery, policies to meet net zero
 1701                           BOJ holds fire, offers brighter view of economy as pandemic impact eases
 1649                                      UK economy extends recovery from COVID crash
 1658                                      ECB must keep up support for the economy - Villeroy
 1634                          Japan's second-quarter capex falls most in decade on pandemic blow
 1688     German economic recovery to continue in second-half, third-quarter to show strong growth: ministry
 1691             German economic recovery to continue in H2, Q3 to show strong growth: ministry
 1552               French central banker says any 2020 GDP forecast revision would be up
 1613             India's recovery to take time after economy shrinks 24% in June quarter
 1202                                  UK plans to drop 'Facebook tax', Mail on Sunday says
 255               UK says always reviewing quarantine data, no comment on Portugal shift
Name: title, dtype: object,
 1620                             https://www.reuters.com/article/britain-boe-idUSL9N2ER032
 1619                         https://in.reuters.com/article/us-britain-boe-vlieghe-idINKBN25T261
 1717                             https://www.reuters.com/article/us-britain-boe-idUSKBN2673NS
 1526                 https://www.reuters.com/article/britain-economy-consumersentiment-idUSL8N2FM2YN
 1654               https://www.reuters.com/article/health-coronavirus-britain-economy-idUSKBN2620MS
 305                 https://in.reuters.com/article/health-coronavirus-britain-tax-idINL8N2FW03T
 1509               https://www.reuters.com/article/healthcoronavirus-britain-economy-idUSL8N2FN2UX
 1671                             https://www.reuters.com/article/japan-economy-boj-idUSL4N2G714U
 1426                         https://uk.reuters.com/article/uk-britain-politics-data-idUKKBN25Z3CD
 1713                         https://www.reuters.com/article/us-climate-change-britain-idUSKBN26410M
 1701                     https://www.reuters.com/article/japan-economy-boj-decision-idUSKBN268090
 1649                                 https://www.reuters.com/video/watch/idOVCV929AJ
 1658                         https://uk.reuters.com/article/uk-ecb-policy-villeroy-idUKKBN2612WC
 1634                         https://in.reuters.com/article/us-japan-economy-capex-idINKBN25S35N
 1688                         https://www.reuters.com/article/us-germany-economy-idUKKBN2650ZK
 1691                         https://www.reuters.com/article/us-germany-economy-idUSKBN2650ZK
 1552               https://www.reuters.com/article/us-france-economy-villeroy-idUSKBN25N2RE
 1613                     https://www.reuters.com/article/us-india-economy-gdp-idUSKBN25R1MT
 1202                         https://www.reuters.com/article/us-britain-usa-tax-idUSKBN25J09E
 255         https://in.reuters.com/article/health-coronavirus-portugal-britain-quar-idINL8N2FY3VF
```

Fig 6: Results for text-based recommender (title "Bank of England policymakers warn of bigger risks for UK economy"

```
(1620                    Bank of England policymakers warn UK economy facing bigger risks
 1619                  UK economy might take years to recover from COVID hit-BoEs Vlieghe
 1717                            Bank of England gears up for next stimulus push
 1196                          JPMorgan to launch UK digital lender in early 2021: Sky News
 1300                    Experts warn: High-tech tools to fight COVID-19 pose their own risks
 1658                             ECB must keep up support for the economy - Villeroy
 1654                     UK economy extends recovery from COVID crash, growth seen fading
 178                    COVID generation risks child marriage, forced labour, ex-leaders warn
 1552                  French central banker says any 2020 GDP forecast revision would be up
 1604           Australias central bank has limited options as economy sinks into steep slump
 1671            Exclusive: BOJ to offer brighter view on economy as COVID crisis eases: sources
 1560    Indian government consumption key to growth in economy amid pandemic, central bank says
 1629                   Feds Mester says central bank wont let inflation run rampant
 1618              Brazil economy back to 2009 size after record 9.7% slump in second quarter
 1652                            ECB sees "strong rebound" signs, monitoring FX
 1097                              Britains Tesco to trial drone deliveries
 1642                Climate change may wreck economy unless we act soon, federal report warns
 142                               Column: Megacities after coronavirus
 769                  The GoPro Hero9 Is a Little Bigger and a Lot Better in Every Possible Way
 1660     Global economy seeing sharper V recovery, raising case for inflation - Morgan Stanley
Name: title, dtype: object,
 1620                                    https://www.reuters.com/article/britain-boe-idUSL9N2ER032
 1619                                  https://in.reuters.com/article/us-britain-boe-vlieghe-idINKBN25T261
 1717                                    https://www.reuters.com/article/us-britain-boe-idUSKBN2673NS
 1196                                  https://www.reuters.com/article/jp-morgan-digital-banking-idUSKBN25H1H0
 1300                            https://www.reuters.com/article/us-health-coronavirus-technology-trfn-idUSKBN25S3X1
 1658                                    https://uk.reuters.com/article/uk-ecb-policy-villeroy-idUKKBN2612WC
 1654                              https://www.reuters.com/article/health-coronavirus-britain-economy-idUSKBN2620MS
 178                            https://in.reuters.com/article/health-coronavirus-education-childlabour-idINL8N2FJ49L
 1552                                 https://www.reuters.com/article/us-france-economy-villeroy-idUSKBN25N2RE
 1604                                  https://www.reuters.com/article/australia-economy-rba-idUSL4N2FY29G
 1671                                    https://www.reuters.com/article/japan-economy-boj-idUSL4N2G714U
 1560                                  https://www.reuters.com/article/india-cenbank-report-idUSL4N2FR1QW
 1629                            https://www.reuters.com/article/us-usa-fed-mester-inflation-idUSKBN25T2TW
 1618                                  https://www.reuters.com/article/us-brazil-economy-gdp-idUSKBN25S507
 1652                                https://www.reuters.com/article/us-europe-ecb-instantview-idUSKBN261288
 1097                                    https://in.reuters.com/article/tesco-drone-delivery-idINKBN2602QB
 1642    https://arstechnica.com/tech-policy/2020/09/climate-change-may-wreck-economy-unless-we-act-soon-federal-report-warns/
 142                                    https://in.reuters.com/article/global-cities-kemp-column-idINKBN25L1ZL
 769                           https://gizmodo.com/the-gopro-hero9-is-a-little-bigger-and-a-lot-better-in-1845058869
 1660                                https://uk.reuters.com/article/us-markets-growth-morgan-stanley-idUKKBN25Z0SD
```

Fig 7: Results for keyword-based recommender (title "Bank of England policymakers warn of bigger risks for UK economy"

V. Front-End

This recommender system was developed across three Jupyter notebooks. To test and further develop the system, the recommender-textBased.ipynb and recommender-keywordBased.ipynb notebooks (i.e. development notebooks) can be edited, tested and extended.  These notebooks allow for testing of new parameters such as the min_df or editing the results to include a different number of recommendations. When using Jupyter lab, these notebooks can be quickly tested using the interact widget at the bottom of the screen.

The third notebook (i.e. production notebook) contains the code for both recommender systems but is formatted to be presented via Voila.  This notebook does not contain any output or markdown, but is useful for presenting the recommender systems.

Fig. 8: Screenshot of part of the keyword-based development notebook (recommender-keywordBased.ipynb)
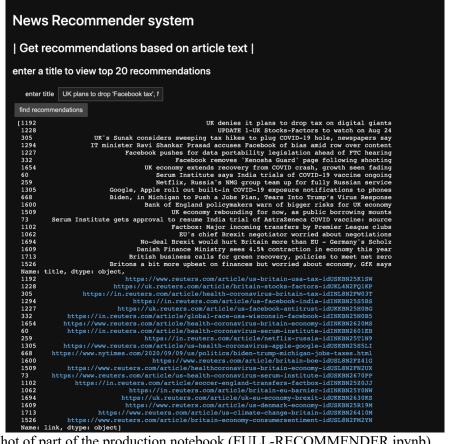


Fig. 9: Screenshot of part of the production notebook (FULL-RECOMMENDER.ipynb)

VI. Conclusions and Future Work

In this project, I have developed a news recommender system based on article text and article keywords.  Since there is not yet any user data, content based methods were used to provide the recommendations about the articles.  TF-IDF feature extraction and cosine similarity methods were used to create similarity measures for the recommendations.

Currently, it seems that both recommendation systems perform quite well despite the absence of any additional user data. At first glance it seems as though the text based recommender is more consistently accurate, however, it is impossible to guess which model performs better in the absence of user feedback. As the project develops, collaborative filtering methods can be employed based on user behavior data that is gathered while they interact with the system. Ideally, future users will be able to give feedback about their recommendations in order to fine tune the algorithms and personalize their experience. Additional features such as creating profiles, "favoriting" articles and searching by other terms (i.e. search by keyword or author) will be implemented during continued front-end development to help make the system more adaptive and customizable.