

Fighting Misinformation with Natural Language Processing and Machine Learning

Sean Griffin, Honey Rattanakasem, Bryan Anderson, Harsh Sangani

Department of Information Science, University of North Texas, Denton

INFO 5810: Data Analysis and Knowledge Discovery

Dr. Haihua Chen

11/1/22

Fighting Misinformation with Natural Language Processing and Machine Learning

This writing describes attempts to process large volumes of articles to locate misinformation, disinformation, and fake news. Numerous techniques are described detailing preprocessing, identification and classification of target data through labeling and tagging, filtering to include or exclude data, data reduction techniques, and final processing to produce final aggregated results for further review. During our research, several techniques were repeatedly used by other researchers and yielded clues to best and easiest to implement techniques. Several techniques such as Naive Bayes, Support Vector Machine, and logistic regression, were highly regarded by researchers and yielded highly accurate results. Within these techniques TF-IDF, count vectorization, stemming, stop words. We will utilize the aforementioned techniques against our dataset in order to identify misinformation.

Introduction

We live in a society where it's difficult to decipher what is factual and what is not. Our viewpoint or perspective on life determines what we believe and acknowledge as the truth. Plain and simple, we are biased (as much as many of us believe we are not) in what we want to believe is the truth. We can't handle the truth, or can we? Many online social media platforms have been the epicenter of information transforming into misinformation, and attempting to control the narrative has always been a difficult task to handle. For starters, how can news be identified as fake and who makes that determination? The idea of manipulating individuals into believing in something that is not true can be powerful and dangerous. However, there have been studies suggesting that natural language and machine learning techniques can be used as a method to identify fake news.

Fake news refers to misinformation and disinformation which gained attention significantly during the digital age. Although fake news exists as long as humans live, it has spread tremendously through social media and the internet. The term "information pollution" and "information disorder" also emerged with the aim of finding a comprehensive term for fake news (Baptista & Gradim, 2022). Lazer et al., 2018 defined fake news as "fabricated information that mimics news media content in form but not in organizational process or intent. Fake-news outlets, in turn, lack the news media's editorial norms and processes for ensuring the accuracy and credibility of information." (p.1).

Literature Review

Although information is accessible anywhere at the fingertips via the development of the internet, it is still very challenging to verify the source of credibility, primarily through social media. News from social media travels at an optimized speed, making it extremely dangerous. Ozbay & Alatas, 2020 introduced two-step methods to fight fake news by pre-processing to convert unstructured data into structured data sets using document term matrix and implemented twenty-three artificial intelligence algorithms to transform the data with text mining methods. During the data pre-processing step,

tokenization, stop-words removal, and stemming were performed to remove unrelated and redundant data to improve model accuracy. In this research, the authors selected twenty-three AI algorithms such as BayesNet, JRip, OneR, Decision Stump, ZeroR, Stochastic gradient descent, CV Parameter Selection, Randomizable filtered classifier, Logistic model tree, and locally weighted learning. The research result shows that the Decision Tree algorithm performed best mean values in terms of accuracy, precision, and F-measure. In contrast, ZeroR, CVPS, and WIHW algorithms are the best algorithms in terms of recall metric. Ozbay & Alatas, 2020 suggested that future work to explore new algorithms and combine testing methods will improve the performance of the models.

While the existing studies primarily focus on utilizing information extracted from the news content, Kaliyar et al., 2021 suggested detecting user-based engagements and the context-related group of people (echo-chamber) sharing the same opinions for fake news detection. The research is done by validating against a real-world fake news dataset from BuzzFeed and PolitiFact by infusing the tensor-factorization with news content from a social network with the combination of deep neural network by employing optimal hyper-parameters to demonstrate the EchoFakeD model. The method achieved a 92.30% accuracy rate, whereas traditional machine learning only achieved 81.30%.

Fake news is often used in different terms, such as misinformation or disinformation (Al-Rawi, 2019). This article focuses on comparing the coverage of fake news in mainstream news outlets (MSN) and the discourse on fake news on social media (SNS). Gatekeeping is an information filtering process deciding whether or not messages should be edited, filtered, blocked, or disseminated. Since there is no filter or centralized gatekeeper in social media, fake news is widespread in just a few seconds by the users sharing the information. This research analyzes a unique data set collected from over 8 million tweets and about 1,350 news stories from different media outlets that reference fake news. This study investigated the discourse on fake news on Twitter and the role of influential or popular online users rather than identifying all the original creators of fake news stories by using topic modeling with a computational method that has limitations in decontextualizing the data set.

Usai, et al. (2018) The author implemented a manual labeling approach in this study, combining each document with a label or tag acquired from the term extraction method. Hence, the author used two critical titles, which are, in taxonomy terms, interpreted as text mining and knowledge discovery.

In this study, the authors reexamined 85 various research studies published between 1998-2017. Next, they attempted to classify the results of the text-reviewing procedure, which discloses that the papers be in two main clusters with distinct and common attributes. Later, it processed and distributed the articles per year, implementing a literature analysis that sets and forecasts the opportunities and challenges for the organization in taking part in social & consumer-based data analysis.

The research result shows that from 1998-2009, the research on discovering knowledge by implementing text mining was addressed in a more technical way. For example, more studies were in the biomedical and engineering fields, and later the spread became more focused on reducing bias in patients' diseases. From 2010-2017, new fields in the area were involved. For instance, business scholars were conducting research to understand the behavior of customers more effectively to increase their

profits by reducing their costs and boosting their sales. Based on the results of the text-mining analysis, the author classifies the paper reviews into three main branches based on frequent and distinct characteristics that are technical algorithms, framework analysis, and performance platforms.

Ahmed, et al. (2022) discuss how detecting fake news is new compared to other research. Many researchers are interested in these fields due to concerns raised by people. Machine Learning is found to help simplify and resolve a complex set of problems. The author used machine learning algorithms to discover fake news and implemented three classifiers, Passive Aggressive, Naïve Bayes, and Support Vector Machine, on two open-source datasets collected from Kaggle that were issued during the election in 2016. The dataset comprised 18,000 news articles set apart through binary labels 0 and 1. The database was already categorized qualitatively into fake, non-fake, and not clear labels.

During preprocessing, the objective was to decrease the size of the database by excluding unrelated information, which is not required for classification. Also, for processing, the author reclassifies the data by including the first half of the set with fake data and the other half of the dataset with non-fake data. The words were changed into their base form for more clarity. After that, the author implemented a stemming algorithm to reduce the number of words based on word type and class.

Different models for classification were introduced, which can be implemented for this study. Still, to select one, the author ran various experiments and the model, which shows satisfactory results in related classification tasks. Support Vector Machines, Naïve Bayes, and Passive Aggressive showed promising results, so they continued exploring the methods.

The results highly encouraged the approach since the implementation helps to classify fake news and determine important features which can be used to detect misleading information. The system developed in this study shows an accuracy of up to 93%.

Chen, et al. (2022) analyzed the issue of how to choose a machine-learning model for the classification of legal text. To solve this problem, the author conducted a detailed study of the type of legal text with the help of two approaches: classification based on domain concept using random forests and work embeddings using deep neural networks. Domain concept using random forest implements the ATCT framework to discover helpful terminologies. The classifier was trained with the help of random forests algorithm. This procedure has also been compared with other deep-learning models formed upon various pre-trained models such as BERT, GloVe, and Word2vec. The experiment's result on the sub-dataset of SigmaLaw shows that the domain concept-based random forests classifier improves the accuracy, recall, precision, and F1 score by 35.26%, 21.17%, 26.54%, and 24.27%. The author also analyzes how the domain concept's size influences classification performance; the top 5% of the domain concepts can bring out the most robust and effective random forest classifier.

Granik & Mesyura (2017) utilize a Naïve Bayes classifier and describe a set of criteria for classification to identify fake news articles. Their criteria are based upon recognition of grammatical mistakes, emotionally colored language, attempt to manipulatively influence opinions, blatantly untrue information, and similar word sets. Posts were initially subject to Human verification to determine if they are “mostly true”, “mostly false”, “mixture of true and false”, or “no factual content.” If reviewers were uncertain, they could mark a post for review by another reviewer. If the two reviewers disagreed on the post, a third person became a tie breaker. As a final preparatory step, all false posts were verified. Posts

with a mixture or no factual content were discarded and the data was split into three datasets to accommodate training, validation, and test scenarios.

Abdullah-All-Tanvir, et al. (2019) make use of Naive Bayes, logistic regression, and support vector machine (SVM) in their Twitter tweet analysis to discover falsities in their dataset. They also employ count vectors and TF-IDF to train their classifiers. Count vector considers a row of text and the number of words within it. Each word is considered a column and the frequency of any given word in a row is captured. For TF-IDF, the TF is the number of times a term appears in a document divided by the total terms in the document. The IDF is a logarithmic function that considers the total number of documents divided by the number of documents with the term in it. These researchers consider Naive Bayes and SVM to have more accurate results.

Ahmed, et al. (2021) detail their approach to inclusion and exclusion preprocessing as follows: Include only English articles where the full text is accessible and the paper contains the terms machine learning, fake, or false news. This is valuable advice for our article as a basis to filter or exclude articles we cannot read or access, and that drill down into our specific search terms. Further, like other authors, mention removing stop words and the application of stemming to generate better results.

To believe or not to believe, what is the factual information? That is the question I am sure most of the population is asking themselves frequently these days. With the overwhelming amount of information circling all over the internet one must ask what is factual and what is not? It seems like everyone claims to be a “trusted source.” Technology, as always, is being created and upgraded to help humans solve a problem. In this case, machine learning and text mining technologies are helping combat the misinformation (or fake news) problem that exists in our world.

For machine learning and text mining, the role consists of a combined effort that mixes computational science, machine learning, and computational linguistics to fight the battle against infodemic. Machine learning and natural language processing are not new concepts in the fight against misinformation and fake news. In fact, prior to the overabundance of information from the COVID-19 pandemic, machine learning and language processing played a significant role in trying to reduce the spread of fake news as much as possible. Although different technologies are making a valiant effort to restrict the flow of misinformation, there is still a long way to go and the invitation for innovative ideas is always welcome.

Machine learning and natural language processing approaches have been used by many different researchers who have used several methods in a variety of ways to look for positive outcomes. In this article, few of the researcher’s algorithms were trained by bigram features that are a sequence of two words that occur in the text (Abdeen et al., 2021). The outcome was to achieve maximum performance through the training of the specific algorithms chosen.

Even though machine learning and text mining are considered the foundation for text classification and anomaly detection, being trained from static sets like bigrams leaves more questions than answers as no further work has been developed as a result. While this is true, in this model, the features can connect the bigrams in a natural way to the text. As a result, the entire training is not

necessary as the model is made extensible due to the new addition of the datasets. By using inherent centrality measures, the reduction of noise is allowable through the use of pruning. If necessary, the network training model used in this exercise gives the option of using multi-label classification, which involves the prediction of one or more class labels by utilizing network clustering techniques (Abdeen et al., 2021).

The total collaboration for this experiment can be summarized due to a network model that can be trained with new datasets without keeping any of the old datasets. Also, binary classification can be enabled that classifies objects into two groups.

One of the machine learning algorithms used was NeoNet, which was designed specifically for COVID-19 news classification. For this research, the part of the workflow describes how the methods and the various processes, beginning with a news dataset, transitions to being examined for bigrams using TF-IDF, making a network model and then training the algorithm with a model to predict the outcome from the tests (Abdeen et al., 2021).

NeoNet algorithm was measured from a series of classification experiments using several datasets and configurations. The reason for doing this was to find a threshold that could generate the best outcome available. Bigrams were the factor resulting in either a high or low threshold. The higher the number of bigrams, the better the chances of classifying an article. If it is a low threshold, there is a great chance of a slight overlap making the algorithms less accurate.

The overabundance of misinformation and disinformation pertaining to the COVID-19 pandemic crowding the online outlets, such as social media platforms, has become a serious threat for the public. In this article, the authors provided evidence of disinformation taken from a publication, which was confirmed to be an unreliable source. The reason it was deemed an untrustworthy source was because within the article, health issues were linked to vitamin D, and the authors themselves were not associated with any expert health group that could substantiate these claims. Furthermore, these false claims were highlighted by DailyMail on Facebook and Twitter platforms, which can be seen as a global threat to the world's public health.

There were some recommendations that were implemented due to suggestions from the science community. For instance, information monitoring and knowledge refinement solutions were used to address the problem directly from the source. Also, the research documented specifically what tools and research methods needed to be used.

Based on the results, NeoNet was implemented to highlight textual contents as safe or contested for COVID-19 for the public to read (Abdeen et al., 2021). NeoNet has the potential to be readily available in the battle against COVID-19 infodemic when compared to top machine learning algorithms such as support vector machines and decision trees. This research confirmed the need for allocating new defense mechanisms and practices against misinformation and other debated content online.

All of humanity relies on information to make decisions and guide their life choices. Unfortunately, it is becoming increasingly difficult to determine whether the information we are receiving is factual or not. One's perception is their reality, our biases and the way we perceive things will determine what information we will determine as truth and what we will discard as a falsehood. Many platforms – social media, news outlets, hearsay, you name it – provide information and it is up to us to determine if this information is dependable or if we are receiving fake news. While we may not be able to analyze some platforms for fake news, technology has made it easier to assess information in electronic formats.

This literature review will discuss some of the proposed text analysis methods used to detect false news including, count-vectorizer, N-gram model, TF-IDF Vectorizer, and machine learning algorithm. This study will further advance the issue with identifying false news and more importantly, how to combat it using data mining tools

One way to analyze documents is count-vectorizer, which is a method that takes a set of text documents and changes them into a token count matrix (Raja & Raj, 2022). In other words, this method basically breaks down a collection of texts into words, outlining the number of occurrences of each word or phrase to assist in recognizing patterns within the documents.

Additionally, N-gram models can be used to review text. There are several N-gram models to choose from when it comes to text order, but the most sought out models are word-based n-grams and character-based n-grams. The difference between the two is that a word-based n-gram is used to show the report's environment and to produce highlights for record grouping while character n-grams are used more for the identification of authors, language, speech analysis, text, and numerical classification (Raja & Raj, 2022). In general, n-gram models are employed to differentiate between genuine and fake news by finding many collections based on the data.

To measure the importance of a term in a dataset document, one would probably choose TF-IDF Vectorizer. The more times a term occurs in a document, the more substantial it is. TF-IDF is a more suitable model because it focuses mainly on the frequency of terms present in the corpus while offering the significance of the terms (Raja & Raj, 2022). This makes analysis of word importance less difficult due to removing words that are not as important for evaluation such as "an" and "the".

The machine learning algorithm divides datasets into two groups (train and test). Generally, the training set, which is usually larger in comparison to the testing set, is somewhere around the 70% threshold while the remaining 30% is used for testing the accuracy of data.

According to the test, a support vector machine with TF-IDF gives the best results with a 93% accuracy rating.

I think it is safe to say that regardless of who you are in this world, we have all been deceived by misinformation that we believed was factual. Beyond that most of us, if not all of us, can at least concede that the outcome of fake news can be detrimental. Once the news is out there, it's amazing how fast it

can reach an audience – there's almost no stopping the dissemination of this information. That can be troubling and disappointing because we all rely on information regardless of who we are. Nevertheless, there are ways to detect fake news to some degree. Using machine learning algorithms and various extraction methods such as TF-IDF and SVM classifiers can help achieve high accuracy percentages of detection.

The combination of the internet and social media platforms has allowed more people than ever before to receive news from a large variety of sources as opposed to the old format of news outlets. If you are an individual who spends a significant amount of time browsing through social media, you are highly likely to encounter a great deal of misinformation. Unfortunately, this has become all too common as social media has become a part of our daily lives. For this research, ensemble learning approaches have been planned to assist in recognizing gaps between the detection of fake news versus detecting deceptive reviews.

It's funny how humans believe they are right more often than not. We as humans think that we are so smart that we can't be tricked or fooled. Welcome to the 21st century. We all have been duped, conned, tricked, hoodwinked, and whatever other adjective you can think of. We believe what we want to believe and will discard anything that goes against our way of thinking because we believe it's not true. It's difficult to either accept information that goes against our perception of "popular belief," or refute information that supports our viewpoints. This is especially true since our thought processes have been ingrained in and influenced by social media platforms over and over. There is a saying, "that the more you tell a lie, the more you will start to believe it is true" and that happens a lot with fake news. Many people will say something is true but lack the factual evidence to support their statements. Combating fake news is a major concern because of the negative impact it can have.

There are methods that are being utilized to fight fake news such as machine learning tools. There are researchers who have used deep learning methods to recognize fake news and to this point, it has shown promising results. A lot of researchers have used different types of models that perform differently on data that is structured (Hansrajh et al., 2021). One thing to mention is that implementing a simple model such as logistic regression can deliver nice performance results, however it does not apply the same to a more sophisticated model turning in great results. Deep learning takes time, effort, and most importantly patience.

In this article, there were two datasets chosen to conduct this experiment. It included several categories that contained both fake and factual articles. The two datasets can be easily found and are readily available on the web. One of the interesting challenges is that when categorizing information as "fake news", it can be difficult and time-consuming because of its accessibility and issues related to copyright legal matters (Hansrajh et al., 2021). The datasets that were composed consisted of either "Liar" which contained statements from politifact.com or "ISOT" which contained statements both truthful and fake from a variety of domains.

To begin, the proposed model started out with the data being preprocess and then extracting the features. Once that was complete, blending ensemble techniques occurred which were used to predict and evaluate the performance results.

Blending ensemble model is a technique that is used to assist in improving the performance and increase accuracy. It is just like stacking and follows the same methods but when it comes to train set to make predictions, it only uses a validation. Ensemble methods mostly generate several accurate solutions as opposed to single models.

Logistic regression is a model that is used to measure the correlation between variables by calculating the probability scores.

A support vector machine is a machine learning algorithm that identifies a hyperplane that separates the data points from different classes. The key purpose of a support vector machine is to maximize the largest margin gaps that split the classes.

Linear discriminant analysis is an easy and effective method for classification. It estimates the probability that a new set of inputs are supposed to fit in every class.

The authors performed analysis from traditional machine learning models and blending ensemble for both the Liar and ISOT datasets (Hansrajh et al., 2021). There were six performance measurements used including ROC AUC, f1-score, AUC, precision, recall, and accuracy.

According to the results, logistic regression achieves the best scores for logistic regression with four out of six metrics. For the ISOT dataset, the linear support vector machine was labeled the best performing base model with a score of five out of six metrics. The overall consensus for both datasets was that the blending ensemble was the top performing model.

To conclude, there were six machine learning models featured with the goal to discover fake news. Based on the experience, the best performing model for both datasets is the blending ensemble. This was confirmed by utilizing several metrics. Although the strategy was to use other larger datasets and analyze trends on social media that are linked to fake news, the consensus is that often the data is not reliable, which can be prone to mistakes and/or anomalies for the model (Hansrajh et al., 2021).

Data Collection & Processing

A dataset of scholarly articles has been provided with the following 11 attributes: Title, Authors, Venue, Year, Citation Count, Fields of Study, Abstract, DOI, Query, Database, and Publication Type. Although the full article text is not included in the dataset, the entirety of the abstract is present, providing an opportunity for natural language processing and machine learning. Additionally, the DOI number may be used to access the full text of the original writing.

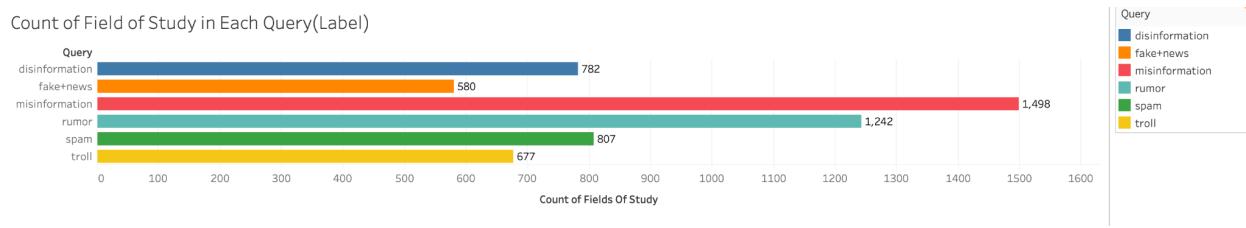
A significant portion of articles are foreign in origin and present challenges both in translation and understanding. For that reason, our research is limited to articles written in English, or those able to be translated to English. The venue attribute yields clues regarding the article's origin.

The Citation Count attribute is interesting in its potential to inform on the importance of a particular article as higher citations may indicate increased adoption of viewpoint, reach of the authors or organizations, level of commitment and involvement in the research community, or timing of the article.

Pulling the data file into Excel and applying filters allows us to sort and group by attribute. For example, sorting by Venue, it is easy to see six articles were released through the American Journal of Infection Control or that 91 articles in the dataset were released in 2011 and 1716 in 2021, a decade later. The articles from 2021 comprise just over 30% of the returns in the dataset.

There are total 5586 records in the dataset which were already labeled or identified in six categories: 782 Records classified as disinformation, 580 records which were labeled as fake+news , 1498 records identified as misinformation, 1242 records classified as rumor, 807 records classified as spam, and 677 records which were classified as troll.

Figure 1: Count of Field of Study in Each Query



Methodology

Our data processing begins with a rudimentary data analysis using four methods. First, we pulled the data into Excel in order to view the columns and number of rows. Additionally, we were able to apply filters, sort and group data by column to obtain a better view of the data and assist with processing ideation. Columns were inspected for missing data, special characters, and other anomalies. Second, the data was inspected using RapidMiner, again, verifying a complete data set, and providing the tools to inspect for duplications and missing data. The Replace operator may be useful to identify and replace special characters using a regular expression for pattern matching. Third, based upon the advice gleaned from other researchers we established the following criteria to include or exclude articles from our dataset: Include articles in English only, Include articles with full text only, established our search criteria to include to the keywords disinformation, fake+news, misinformation, rumor, spam, and troll. The fourth stage consists of removing stop words, and application of a stemming algorithm.

Next, we decided upon classifiers that determine if news is fake. One approach is a supervised training model. With this approach a researcher must manually supervise and train the system to recognize the classifiers. One well known example of this approach is CAPTCHA and reCAPTCHA. CAPTCHA is an acronym for “Completely Automated Public Turing test to tell Computers and Humans Apart.” (SOURCE!!!) CAPTCHA’s original intention was to automate machine text recognition. Therefore,

the system must be trained to recognize unknown variations of text which would be the classifiers. For this study our criteria to determine fake news or misinformation is(DECISIONS MUST BE MADE!!!)

After preprocessing we considered several processing and modeling strategies. For example, we wanted to discover the top 25 words (excluding stop words and including stemming) found in the data, and the top 10 occurring phrases.

NLP processing examples included translation software, chatbots, spam filters, and search engines, to grammar correction software, voice assistants, and social media monitoring tools.

Data analysis techniques include Logistic Regression, Naive Bayes (with bag of words), Machine learning algorithms like “CountVectorizer”, “TfidfVectorizer”, Support Vector Machine, and Passive Aggressive Classifiers for the identification of false news in public data sets.

Other possible techniques include Binary Classification & Regression, Data Mining Software - Rapid Miner, and other custom Python algorithms.

Table 1: Model Descriptions (Ahmed, S., Hinkelmann, K., & Corradini, F. (2022))

<u>Support Vector Machine</u> : It performs supervised learning on data for regression and classification. The SVM computes the data and converts into various categories. The advantages of Support Vector Machine are learning speed, accuracy, classification and tolerance to irrelevant features. Support Vector Machine is one of the most researched classifiers nowadays and it performs well in the fake news detection problem.
<u>Logistic Regression</u> : It is used to estimate the relationship between variables after using statistical methods. It performs well in binary classification problems because it deals with classes and requires a large sample size for initial classification.
<u>Passive Aggressive</u> : These algorithms are mainly used for classification. The idea is very simple and the performance has been proven with many other alternative methods like Online Perceptron and MIRA.

Experiment and Data Analysis Plan

Data Analysis Plan

- Human review data
- Preprocessing data
- Identify phrases or words for inclusion
- Analyze word count
- Look for trends (what type?)
- Year chart
- Citations chart
- Software: Rapidminer

Visualizations:

Figure 2: Citation Count per Year

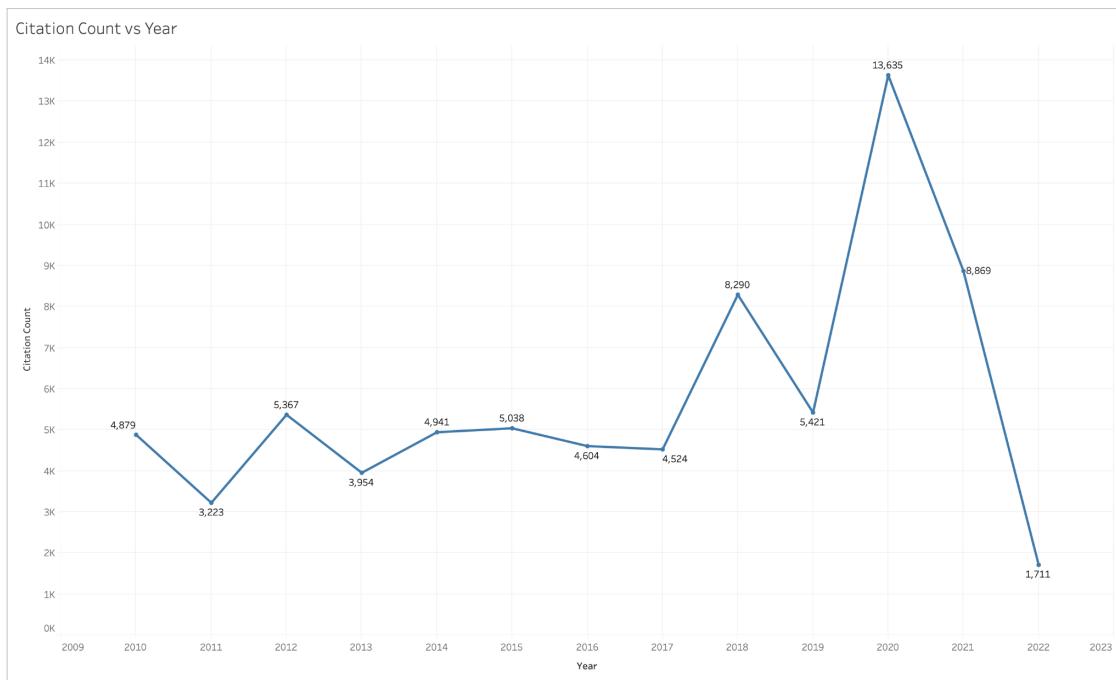


Figure 3: Count of Field of Study In Each Query Per Year

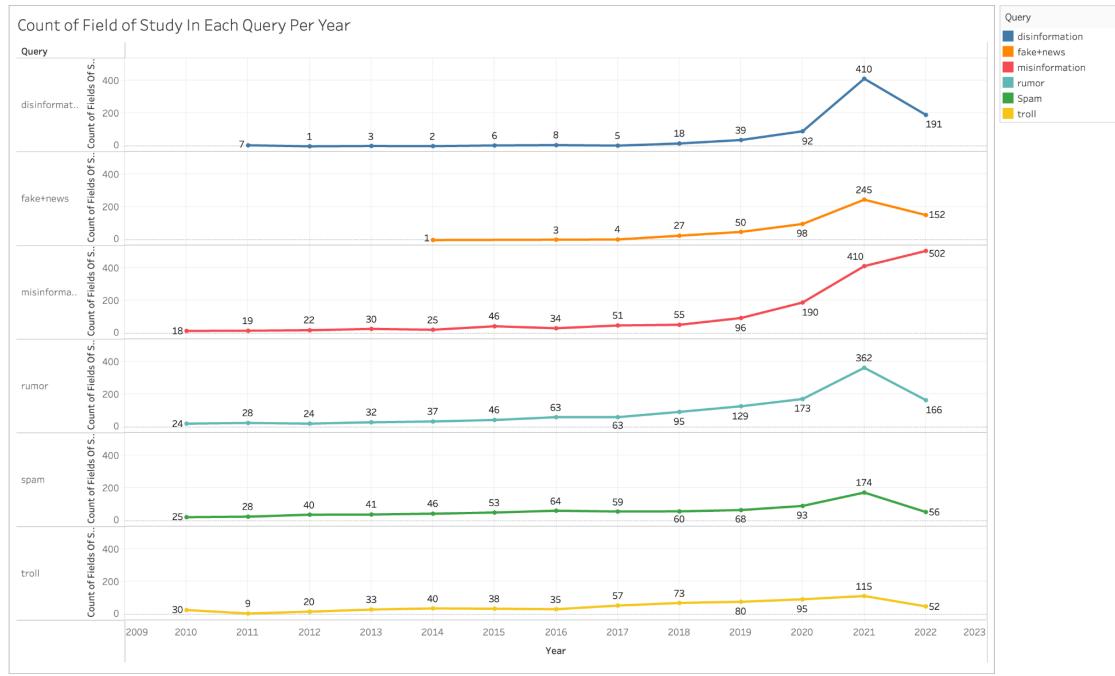


Figure 4: Distinct Count of Field of Study In Each Query Per Year

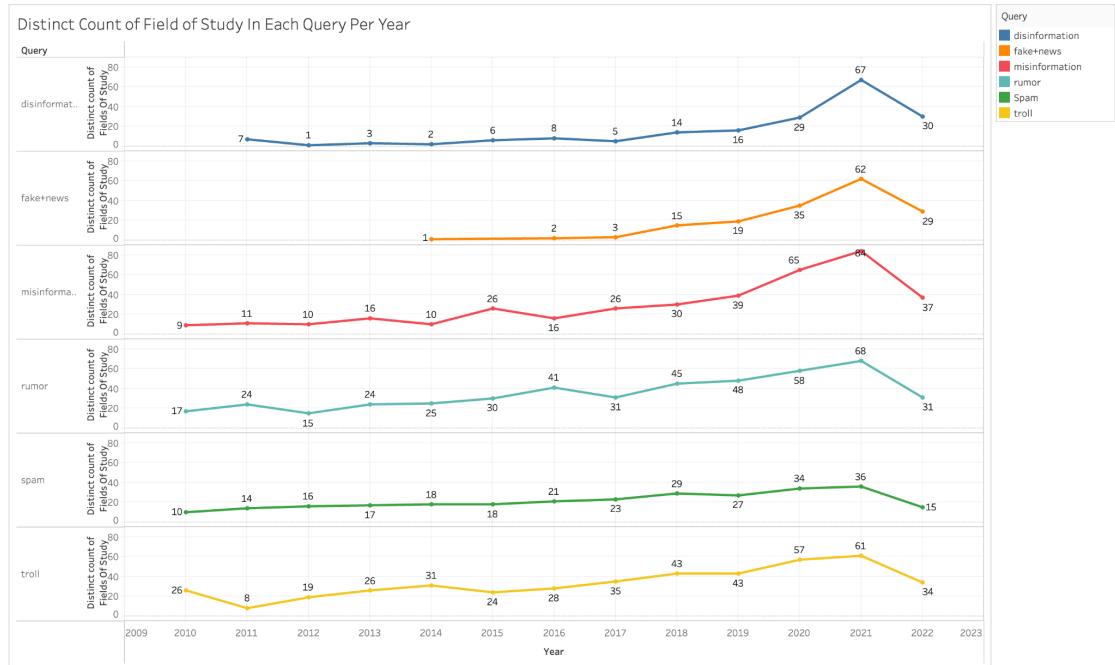


Figure 5: Count of Citations and Field of Study Per Year

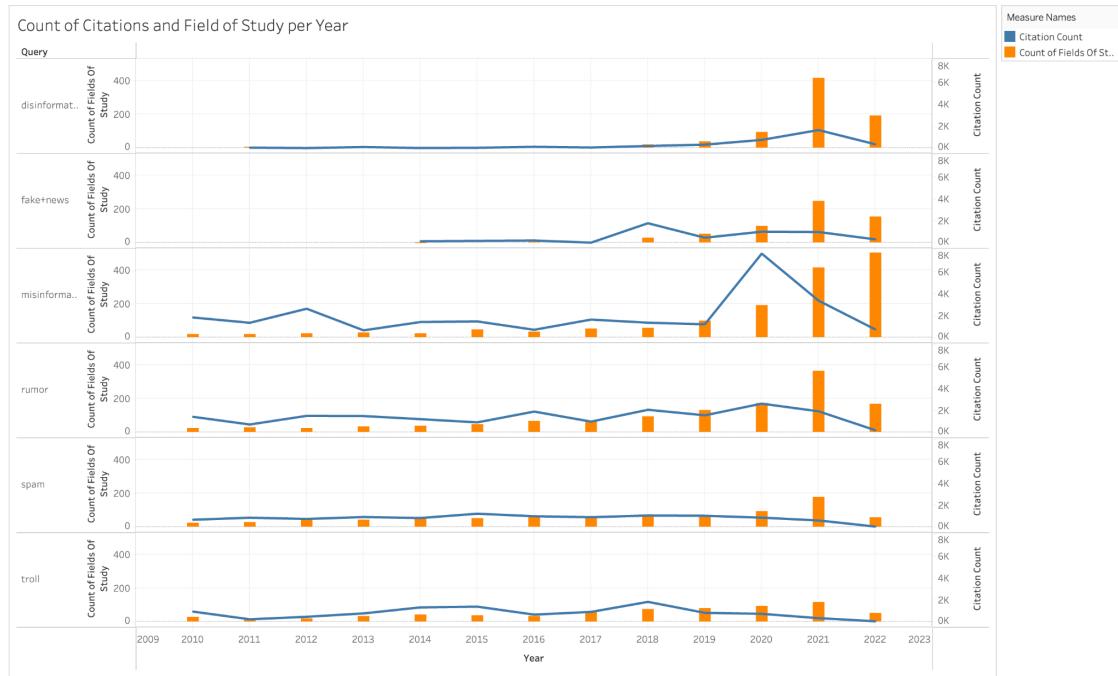


Figure 6: Total Citation Count In Each Field of Study

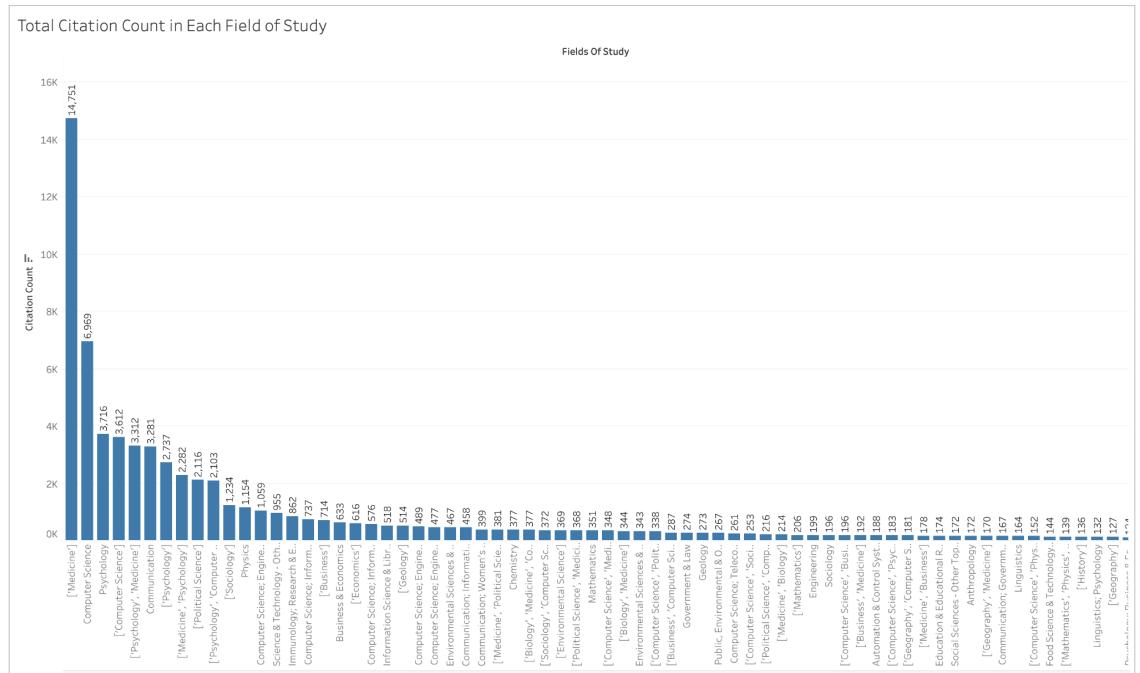


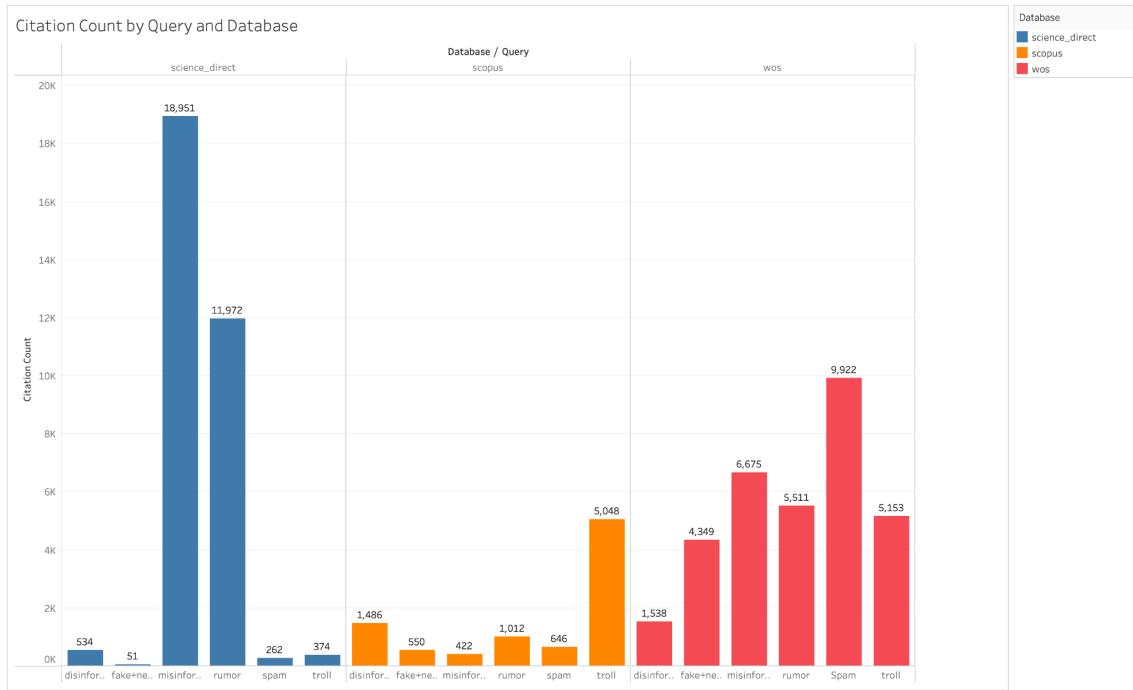
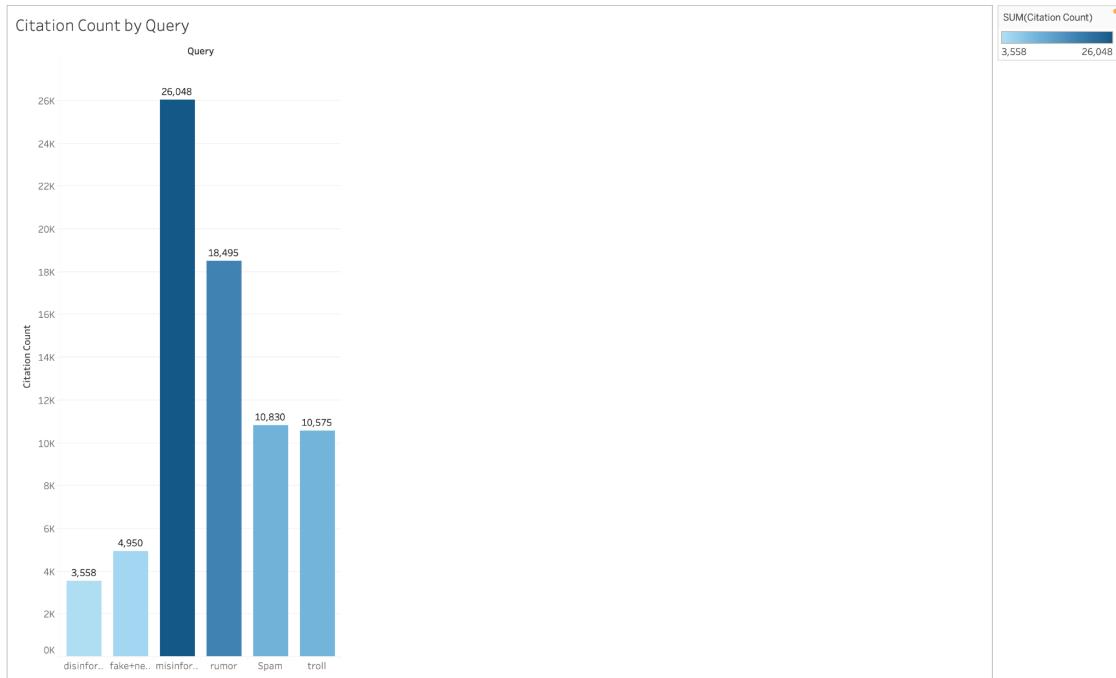
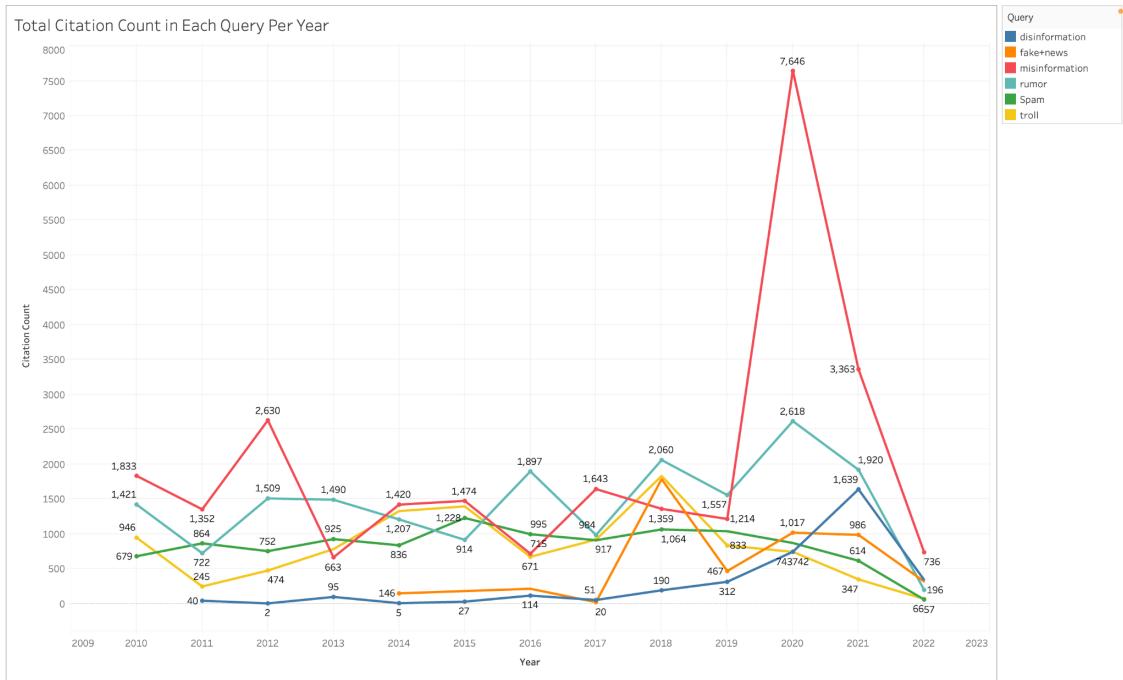
Figure 7: Citation Count by Query and Database**Figure 8: Citation Count by Query**

Figure 9: Total Citation Count in Each Query Per Year



Human Data Review

The trend of collecting and dissecting information on subjects and topics that concerns human interest, especially when it comes to informing and influencing human decisions and actions. This trend, as a result, can cause a variety of challenges for those who are tasked with trying to analyze certain data. With the lack of ability to understand and interpret large amounts of data, it can be difficult for researchers who are attempting to take the information and interpret it to the best of their ability. When looking at the data assigned to my team, we observed many inconsistencies that we thought may not be relevant or could contribute significantly to our end results. Some of those inconsistencies included columns we couldn't use, html tags, symbols, or special characters, and detecting values in ID columns that were inconsistent with the rest of the data.

Strategy

We have chosen to implement a two-pronged approach to preprocessing. First, we attempt to clean the data by removing duplicates, identifying, and replacing missing values. Second, we remove all non-text information in the Abstract field, remove all HTML markup, tokenize the data, normalize it by converting all text to lower case, then apply stemming and stop word algorithms to allow RapidMiner to create clusters later in our analysis.

For our next phase of analysis, we decided upon four initial process strategies. The first strategy involves clustering word associations and results in a visualization to aid analysis. The second strategy utilizes the FP-Growth algorithm in conjunction with association rules to generate frequently found token combinations across our Abstract field. The third strategy analyzes sentiment to rank each abstract as a net positive or net negative with regard to the intent of the communication. The third strategy employs cross validation using the K-NN. The goal amongst all strategies is to exploit data preprocessing operators to expose commonly occurring patterns with high support and confidence and provide a set of visuals to better interpret the data.

Data Preprocessing (Process Documents from Data):

For Data Preprocessing we built and reused the same subprocess for three of our larger processes. The the following operators are included:

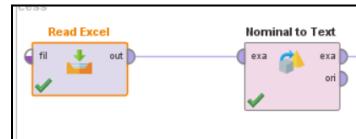
1. Remove ID Column: The ID column contains inconsistent data and needed to be removed prior to importing the data into RapidMiner.
2. Remove Duplicates: This operator removes duplicate examples from the dataset. The all_selected-data_automatic-analysis dataset that we will be using does not have any duplicate values.

Row No.	id	title	authors	venue	year	citationCo...	fieldsOfStu...	abstract	type	database	publication...
1	1	Identity, sta...	Freedman, J...	REVISTA ICO...	2021	1	Communicat...	The followin...	disinformation	wos	journal
2	2	The anatom...	Colom-Piella...	HISTORIA Y ...	2020	2	Film, Radio ...	Revealed wh...	disinformation	wos	journal
3	3	Meta-reflexi...	Golob, Teo...	COMUNICAR	2021	1	Communicat...	The rise of ...	disinformation	wos	journal
4	4	Creation, di...	Agarwal, Na...	ASLIB JOURN...	2021	4	Computer S...	Purpose Thi...	disinformation	wos	journal
5	5	Artificial Inte...	Kerysova, K...	SECURITY A...	2018	7	Government...	This article ...	disinformation	wos	journal
6	6	Disinformati...	Marshall, Jo...	COSMOPOLIS	2017	12	Sociology	This paper ...	disinformation	wos	journal
7	7	Disinformati...	Budaykova, ...	PSYCHOLOG...	2018	0	Psychology	The article ...	disinformation	wos	journal
8	8	Methodologic...	Puebla-Marti...	PUBLICATIONS	2021	0	Information ...	We live in a ...	disinformation	wos	journal
9	9	Artificial Inte...	Manfredi Sa...	REVISTA CID...	2020	10	International...	This paper i...	disinformation	wos	journal
10	10	Disinformati...	Valverde-Ber...	COMUNICAR	2022	1	Communicat...	Disinformati...	disinformation	wos	journal
11	11	Lacuna publ...	Kristina, Ari...	JOURNAL OF...	2021	0	Communicat...	The purpos...	disinformation	wos	journal
12	12	Gender disti...	Herrero-Díez...	REVISTA ICO...	2020	3	Communicat...	Information ...	disinformation	wos	journal
13	13	Novel Validit...	Jin, Michael...	CUREUS	2021	1	General & In...	In the past ...	disinformation	wos	journal
14	14	The handlin...	Hernandez, ...	ANALISIS-QL...	2021	0	Communicat...	Yes, Ministe...	disinformation	wos	journal
15	15	THE AGENDA...	Magallón-R...	MIGRACIONES	2021	0	Demography	The objectiv...	disinformation	wos	journal
16	16	Not falling ...	Gray, Kisho...	FEMINIST M...	2021	1	Communicat...	Black lives h...	disinformation	wos	journal
17	17	Infodemic in ...	López-Pujol...	REVISTA ESP...	2020	10	Information ...	This paper ...	disinformation	wos	journal
18	18	Safeguarding ...	Hanley, Mon...	REVISTA CID...	2020	0	International...	Those who c...	disinformation	wos	journal

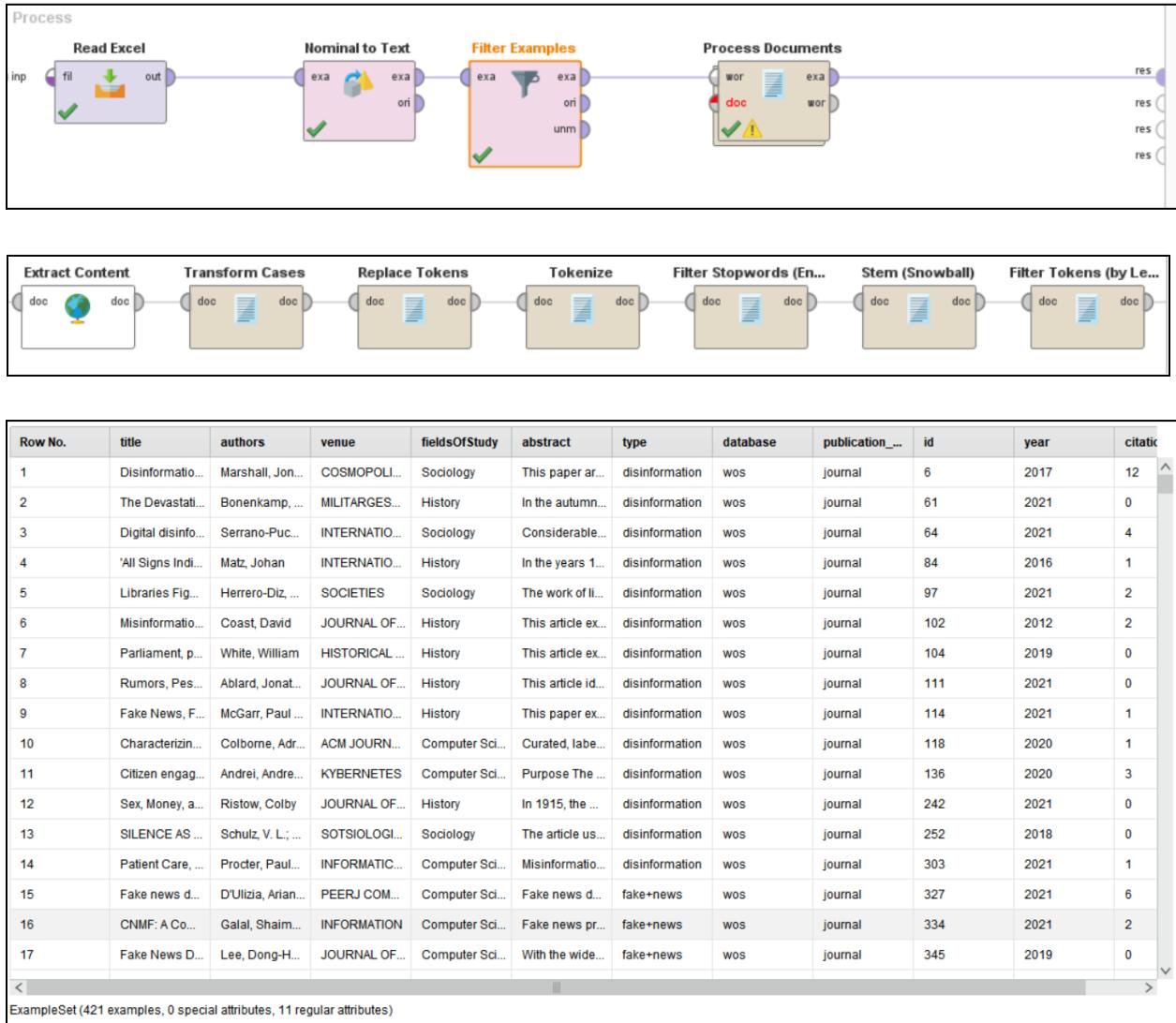
3. Replace Missing Values: This operator helps us by handling missing values. We can replace the missing values with average, none, min, max, zero, and with a particular value such as N/A. In our dataset, there are no missing values present.

Name	Type	Missing	Min	Max	Average
id	Integer	0	1	5586	2793.500
title	Polynomial	0	Least ÁúYou m [...] sions (1)	Most "It's li [...] toric (1)	Values "It's li [...] rhetoric (1), "Trollin [...] rtainment (1), ...[5584 more]
authors	Polynomial	0	Least =Überan [...] iu IA (1)	Most Szpitala [...] muald (5)	Values Szpitala [...], Romuald (5), Neo, Ric (4), ...[5494 more]
venue	Polynomial	0	Least Ááthique & Santv© (1)	Most Vaccine (62)	Values Vaccine (62), Internat [...] ic Health (51), ...[2682 more]
year	Integer	0	Min 2010	Max 2022	Average 2019.252
citationCount	Integer	0	Min 0	Max 3572	Average 13.329
fieldsOfStudy	Polynomial	0	Least [Social [...] ogy] (1)	Most ["Medicine"] (1256)	Values ["Medicine"] (1256), ["Computer Science"] (488), ...[454 more]
abstract	Polynomial	0	Least ÁúUnpre [...] ity. (1)	Most 2014. ht [...] vnpWY (2)	Values 2014. ht [...] fuy-vnpWY (2), A rumor [...] strategic. (2), ...[5547 mo
type	Polynomial	0	Least spam (209)	Most misinformation (1498)	Values misinformation (1498), rumor (1242), ...[5 more]
database	Polynomial	0	Least science_direct (1282)	Most wos (2715)	Values wos (2715), scopus (1589), ...[1 more]
publication_type	Polynomial	0	Least journal (5586)	Most journal (5586)	Values journal (5586)

4. Nominal to Text Operator: To convert the whole abstract column to text.

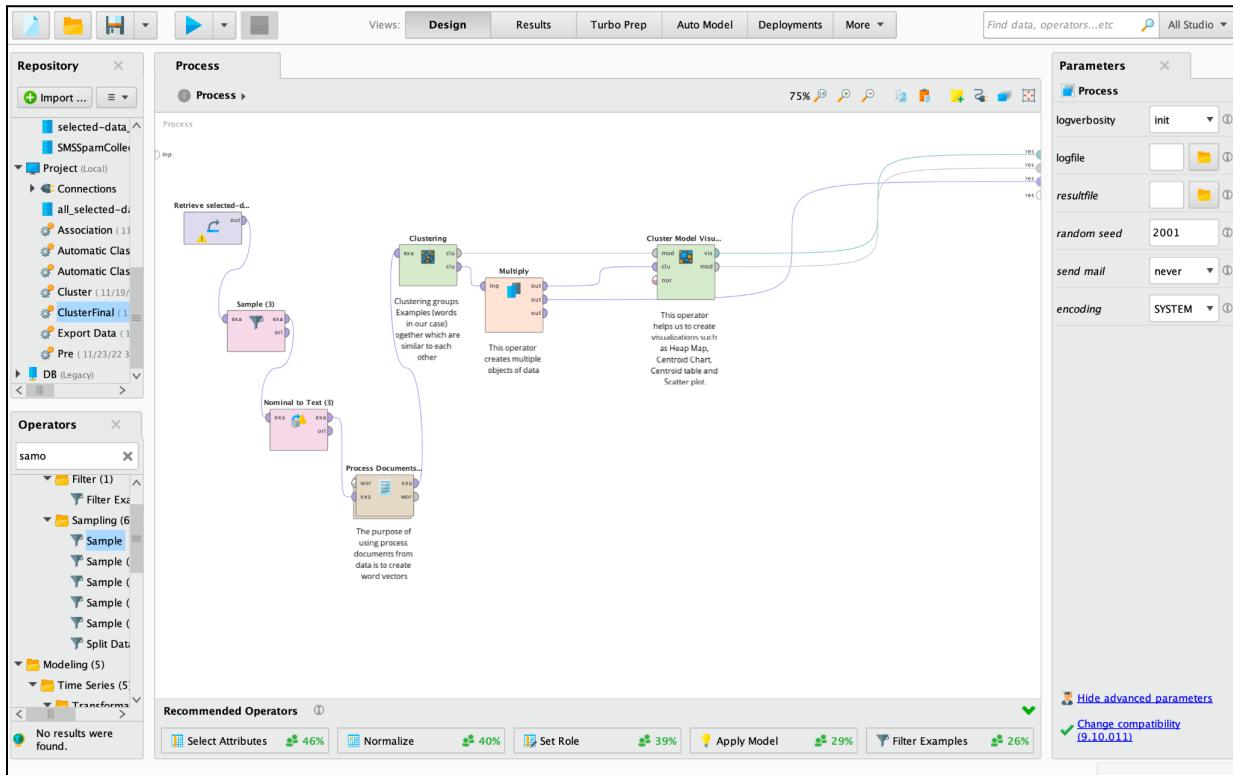


5. Extract Content: Filter out all HTML code and tags from the text.
6. Tokenize: To split the abstract into unique tokens (words).
7. Transform Case: To convert the text into lowercase
8. Filter Stop words: To remove simple English words such as, or, as, &, and.
9. Stemming: This operator stems English words using the stemming algorithm applying an iterative, rule-based replacement of word suffixes intending to reduce the length of the words until a minimum length is reached.
10. Filter Tokens (by Length): To filter tokens based on their length such as the number of the characters. In this case, we set up the min character to 4 and the max characters to 25.



Process 1: Document Similarity and Clustering Process

In this process we attempt to identify similarities in the Abstract field in our data. The similarities we are specifically targeting are term frequency, which allows us to cluster groups of abstracts when they have a high frequency of the same terminology.



Process Documents from Data Operator:

The purpose of using process documents from data is to create word vectors. Word Vectors are a series of new attributes containing numerical data for each token. The simplest type of word vector looks at all the words in our spreadsheet/document & creates a new attribute for each token. Then it looks in each abstract and puts 1 if that word exists in that paper and puts 0 if it does not exist. This type of word vector is called Binary Term Occurrences.

We have selected binary term occurrences for association analysis because we can only use 1 and 0 for association analysis.

We have selected the absolute method for pruning.

Absolute Pruning: To ignore words that already exist in one abstract.

For prune below absolute, we have selected value 2, so it will ignore words which is in only one abstract.

At this phase of the process, we are outputting a TF-IDF score which basically indicates the relative importance of a word in our current abstract field, compared to the relative importance of the word in all of our abstract fields. The higher the TDF score, the more important the word.

Clustering (K-Means)

At this phase of processing we have converted terms to numbers and ranked their relative weights. We now turn to the Clustering operator to create groups, also called clusters, that share some similarities. In our case, we are looking for abstracts that contain multiples of the same terminology. We have the option to specify the number of clusters, or groupings. For example, setting the tool to locate 40 clusters on a sample of 1000 records yields the following results. Each unique term is listed and a relative importance score is listed for term found more frequently.

Row No.	id	cluster ↑	text	abbreviated	abbreviations	abhandlung	aboriginal	abortifacient	aboveground	absenteeism	absorption	abstracting	academicians	accelerate
35	35	cluster_0	importance especially on mani...	0	0	0	0	0	0	0	0	0	0	0
51	51	cluster_0	exponential communication fa...	0	0	0	0	0	0	0	0	0	0	0
191	191	cluster_0	information information propa...	0	0	0	0	0	0	0	0	0	0	0
197	197	cluster_0	interaction individuals individu...	0	0	0	0	0	0	0	0	0	0	0.149
205	205	cluster_0	popularity environment wides...	0	0	0	0	0	0	0	0	0	0	0
210	210	cluster_0	preventing dissemination pro...	0	0	0	0	0	0	0	0	0	0	0
214	214	cluster_0	development information techn...	0	0	0	0	0	0	0	0	0	0	0
219	219	cluster_0	identifies propagation theoreti...	0	0	0	0	0	0	0	0	0	0	0
223	223	cluster_0	information propagation incre...	0	0	0	0	0	0	0	0	0	0	0
228	228	cluster_0	subjective aggressive disturb...	0	0	0	0	0	0	0	0	0	0	0
232	232	cluster_0	metropolises disturbing chara...	0	0	0	0	0	0	0	0	0	0	0
233	233	cluster_0	propagation multilingual envir...	0	0	0	0	0	0	0	0	0	0	0
270	270	cluster_0	techniques propagation propa...	0	0	0	0	0	0	0	0	0	0	0
271	271	cluster_0	propagation techniques evalu...	0	0	0	0	0	0	0	0	0	0	0
275	275	cluster_0	algorithms techniques algorit...	0	0	0	0	0	0	0	0	0.101	0	0
290	290	cluster_0	similarity propagation similarit...	0	0	0	0	0	0	0	0	0	0	0
321	321	cluster_0	increasing popularity reasona...	0	0	0	0	0	0	0	0	0	0	0
626	626	cluster_0	propagation interweaves infor...	0	0	0	0	0	0	0	0	0	0	0
663	663	cluster_0	traditional particular destructio...	0	0	0	0	0	0	0	0	0	0	0
831	831	cluster_0	propagation empirically inform...	0	0	0	0	0	0	0	0	0	0	0
852	852	cluster_0	propagation differential differ...	0	0	0	0	0	0	0	0	0	0	0
880	880	cluster_0	preventing dissemination pro...	0	0	0	0	0	0	0	0	0	0	0

Cluster Model

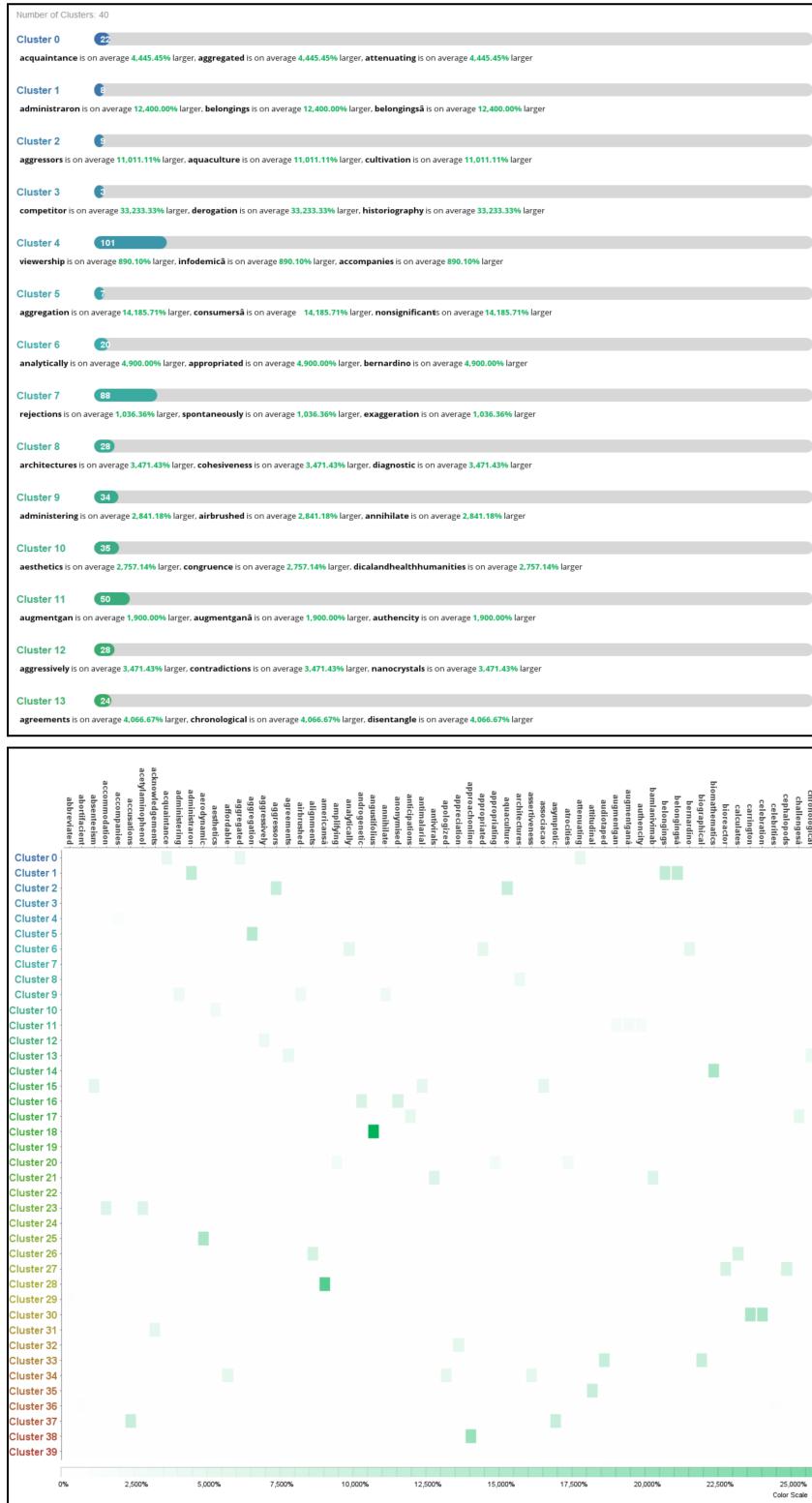
```

Cluster 0: 22 items
Cluster 1: 8 items
Cluster 2: 9 items
Cluster 3: 3 items
Cluster 4: 101 items
Cluster 5: 7 items
Cluster 6: 20 items
Cluster 7: 88 items
Cluster 8: 28 items
Cluster 9: 34 items
Cluster 10: 35 items
Cluster 11: 50 items
Cluster 12: 28 items
Cluster 13: 24 items
Cluster 14: 6 items
Cluster 15: 19 items
Cluster 16: 13 items
Cluster 17: 21 items
Cluster 18: 2 items
Cluster 19: 9 items
Cluster 20: 41 items
Cluster 21: 15 items
Cluster 22: 9 items
Cluster 23: 15 items
Cluster 24: 17 items
Cluster 25: 6 items
Cluster 26: 12 items
Cluster 27: 12 items
Cluster 28: 3 items
Cluster 29: 129 items
Cluster 30: 6 items
Cluster 31: 20 items
Cluster 32: 18 items
Cluster 33: 9 items
Cluster 34: 20 items
Cluster 35: 9 items
Cluster 36: 111 items
Cluster 37: 9 items
Cluster 38: 5 items
Cluster 39: 7 items
Total number of items: 1000

```

Cluster Model Visualizer

The Cluster Model Visualizer operator provides us with two useful pieces of information: The frequently found terms in the cluster and a heat map with frequently found terms.



Process 2: FP-Growth and Text Association Process

Process Documents from Data Operator:

The purpose of using process documents from data is to create word vectors, just as we performed in the previous process.

Row No.	text	aaai	abandon	abandoned	abandonm...	abatement	abbasi	abbreviated	abbreviation	abbreviatio...	abdominal	abi...
12	information disorder sa...	false	false	false	false	false	false	false	false	false	false	fals
13	decade disinformation ...	false	false	false	false	false	false	false	false	false	false	fals
14	minister series part coll...	false	false	false	false	false	false	false	false	false	false	fals
15	objective work analyze...	false	false	false	false	false	false	false	false	false	false	fals
16	black lives continually ...	false	false	false	false	false	false	false	false	false	false	fals
17	paper analyzes spread...	false	false	false	false	false	false	false	false	false	false	fals
18	control information spa...	false	false	false	false	false	false	false	false	false	false	fals
19	disinformation covid pa...	false	false	false	false	false	false	false	false	false	false	fals
20	article presents main f...	false	false	false	false	false	false	false	false	false	false	fals
21	previous research outli...	false	false	false	false	false	false	false	false	false	false	fals
22	recent scholarship indi...	false	false	false	false	false	false	false	false	false	false	fals
23	introduction scientific d...	false	false	false	false	false	false	false	false	false	false	fals
24	communication profess...	false	false	false	false	false	false	false	false	false	false	fals
25	analysis disinformation ...	false	false	false	false	false	false	false	false	false	false	fals
26	rise disinformation incr...	false	false	false	false	false	false	false	false	false	false	fals
27	research analyses imp...	false	false	false	false	false	false	false	false	false	false	fals
28	disinformation armed ...	false	false	false	false	false	false	false	false	false	false	fals
29	digital media environm...	false	false	false	false	false	false	false	false	false	false	fals

In our dataset, there are 5586 unique abstracts, and it contains 15,863 unique significant words. Process Documents creates new attributes for these words which can be seen in the table.

Row No.	text	aaai	abandon ↓	abandoned
2796	communicating health knowledge effectively community level essential shaping resilient urban g...	false	true	false
3175	literature attitudes wind power underpinned assumptions limit scope restrict findings assumpti...	false	true	false
3588	summary text addresses risks involving xenophobia chinese internet social networks covid pan...	false	true	false
4319	abstract habgood coote áústop talking fake news inquiry interdisciplinary journal philosophy ar...	false	true	false
4321	abstract response habgood coote áústop talking fake news inquiry interdisciplinary journal phil...	false	true	false

The abstract located at id 2796, 3175, 3588, 4319, and 4321 contain the word abandon so we can see true in the abandon column.

Numerical to Binomial Operator:

It will convert 0 and 1 to False and True. And the FP-Growth operator requires all variables in true and false only.

FP-Growth Operator:

To identify frequently occurring itemset. Here Support is the proportion of documents that include the itemset/word.

Frequent Itemset result:

No. of Sets: 754	Size	Support	Item 1	Item 2	Item 3	Item 4
Total Max. Size: 4	1	0.399	results			
Min. Size: 1	1	0.398	information			
Max. Size: 4	1	0.393	social			
Contains Item:	1	0.392	study			
	1	0.327	based			
	1	0.306	media			
	1	0.268	using			
	1	0.254	misinformation			
	1	0.254	research			
	1	0.248	analysis			
	1	0.248	data			
	1	0.237	paper			
	1	0.209	news			
	1	0.203	online			
	1	0.202	methods			
	1	0.198	health			
	1	0.197	model			
	1	0.192	covid			
	1	0.169	related			

The word results appear almost in 40% of all the abstracts.

No. of Sets: 754	Size	Support	Item 1	Item 2	Item 3	Item 4
Total Max. Size: 4	1	0.051	efforts			
Min. Size: 1	1	0.051	form			
Max. Size: 4	1	0.050	particularly			
Contains Item:	1	0.050	obtained			
	2	0.182	results	information		
	2	0.168	results	social		
	2	0.196	results	study		
	2	0.161	results	based		
	2	0.122	results	media		
	2	0.139	results	using		
	2	0.108	results	misinformation		
	2	0.105	results	research		
	2	0.112	results	analysis		
	2	0.120	results	data		
	2	0.099	results	paper		
	2	0.075	results	news		
	2	0.089	results	online		
	2	0.142	results	methods		
	2	0.101	results	health		

We got some itemset with the size of 2 terms. We can see that words results and information occur together in about 18% of the documents.

No. of Sets: 754
Total Max. Size: 4

	Size	Support	Item 1	Item 2	Item 3	Item 4
	2	0.053	elsevier	reserved		
	3	0.097	results	information	social	
	3	0.090	results	information	study	
	3	0.075	results	information	based	
	3	0.076	results	information	media	
	3	0.062	results	information	using	
	3	0.064	results	information	misinformation	
	3	0.053	results	information	research	
	3	0.054	results	information	analysis	
	3	0.056	results	information	data	
	3	0.067	results	information	methods	
	3	0.059	results	information	health	
	3	0.050	results	information	covid	
	3	0.085	results	social	study	
	3	0.069	results	social	based	
	3	0.100	results	social	media	
	3	0.058	results	social	using	
	3	0.054	results	social	misinformation	
	3	0.052	results	social	analysis	

We got some itemsets with the size of 3 terms. We can see words such as results, social, and media occur together in about 10% of the documents.

No. of Sets: 754
Total Max. Size: 4

	Size	Support	Item 1	Item 2	Item 3	Item 4
	3	0.057	study	news	fake	
	3	0.057	study	health	covid	
	3	0.066	study	covid	pandemic	
	3	0.051	media	misinformation	covid	
	3	0.081	media	news	fake	
	3	0.052	media	health	covid	
	3	0.062	media	covid	pandemic	
	3	0.050	misinformation	health	covid	
	3	0.059	misinformation	covid	pandemic	
	3	0.072	health	covid	pandemic	
	3	0.053	rights	elsevier	reserved	
	4	0.063	results	information	social	media
	4	0.055	results	social	study	media
	4	0.062	information	social	study	media
	4	0.055	information	social	media	misinformation
	4	0.056	information	social	media	news
	4	0.053	information	social	media	covid
	4	0.061	social	media	news	fake
	4	0.053	social	media	covid	pandemic

There are some itemsets with the size of 4 words/terms which appear together. We can see words such as social, media, news, and fake occur together in about 6% of the documents.

Create Association Rules Operator

Association rules help us to determine which words appear together more frequently.
I have selected 80% for confidence.

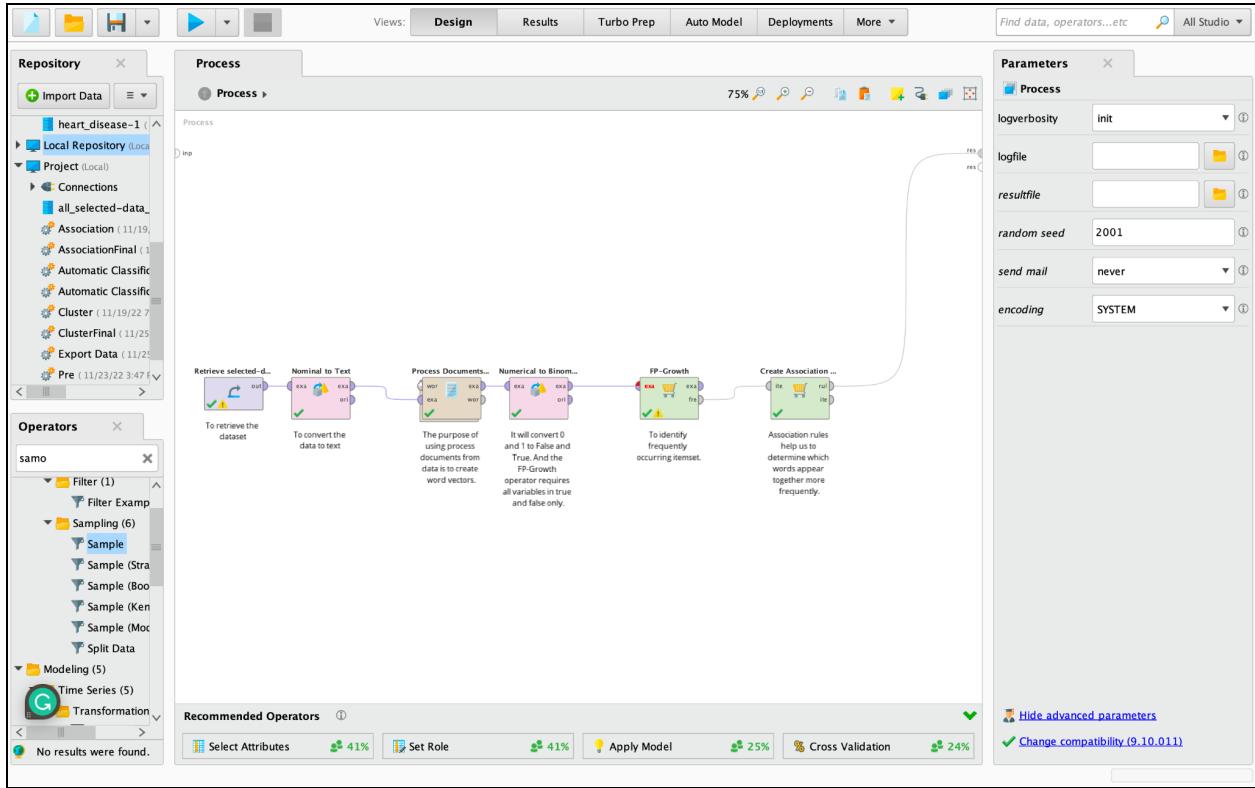
Results of Create Association Rule Operator:

No.	Premises	Conclusion	Support	Confidence	LaPlace	Gain	p-s
63	reserved	rights	0.059	0.991	0.999	-0.060	0.055
62	study, pandemic	covid	0.066	0.948	0.997	-0.073	0.053
61	health, pandemic	covid	0.072	0.937	0.996	-0.082	0.057
60	results, pandemic	covid	0.054	0.929	0.996	-0.062	0.043
59	media, fake	news	0.081	0.926	0.994	-0.094	0.062
58	social, media, fake	news	0.061	0.924	0.995	-0.071	0.047
57	information, pandemic	covid	0.070	0.922	0.995	-0.081	0.055
56	media, platforms	social	0.056	0.921	0.995	-0.066	0.032
55	social, pandemic	covid	0.072	0.920	0.994	-0.084	0.057
54	misinformation, pandemic	covid	0.059	0.919	0.995	-0.069	0.046
53	methods, background	results	0.052	0.917	0.996	-0.061	0.029
52	study, fake	news	0.057	0.917	0.995	-0.068	0.044
51	social, media, pandemic	covid	0.053	0.914	0.995	-0.063	0.042
50	media, pandemic	covid	0.062	0.911	0.994	-0.074	0.049
49	pandemic	covid	0.125	0.910	0.989	-0.150	0.099
48	rights, reserved	elsevier	0.053	0.909	0.995	-0.064	0.050
46	reserved	elsevier	0.053	0.900	0.994	-0.065	0.050
47	reserved	rights, elsevier	0.053	0.900	0.994	-0.065	0.050
45	social, fake	news	0.080	0.900	0.992	-0.098	0.062
44	results, background	methods	0.052	0.895	0.994	-0.064	0.040
43	results, fake	news	0.052	0.893	0.994	-0.065	0.040
42	information, fake	news	0.074	0.892	0.992	-0.092	0.057
41	media, users	social	0.060	0.892	0.993	-0.075	0.034

It created about 65 association rules.

Here as we can see, set of words/premises such as study, and pandemic often occur together.

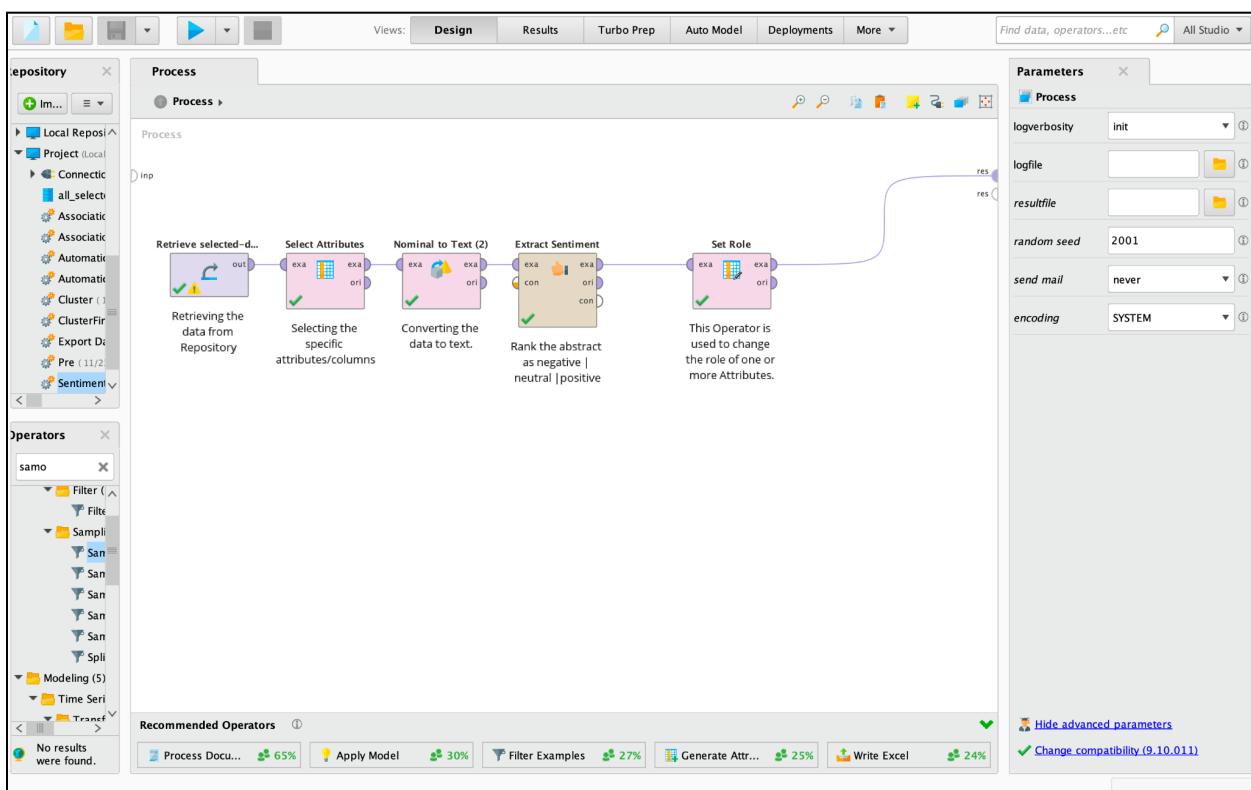
No.	Premises	Conclusion	Support	Confidence	LaPlace	Gain	p-s
49	pandemic	covid	0.125	0.910	0.989	-0.150	0.099
50	media, pandemic	covid	0.062	0.911	0.994	-0.074	0.049
51	social, media, pandemic	covid	0.053	0.914	0.995	-0.063	0.042
54	misinformation, pandemic	covid	0.059	0.919	0.995	-0.069	0.046
55	social, pandemic	covid	0.072	0.920	0.994	-0.084	0.057
57	information, pandemic	covid	0.070	0.922	0.995	-0.081	0.055
60	results, pandemic	covid	0.054	0.929	0.996	-0.062	0.043
61	health, pandemic	covid	0.072	0.937	0.996	-0.082	0.057
62	study, pandemic	covid	0.066	0.948	0.997	-0.073	0.053



Process 3: Sentiment Analysis Process

Sentiment analysis is the process of looking into blocks of text, inspecting each term in the block, and assigning a weight to that word. Sentiment analysis is commonly used with social media platforms to determine sentiment, or the intent of a user post. For example, the terms bad, evil, hate, and similar terms could be described as negative and assigned a negative value. Likewise, terms like good, wholesome, awesome, helpful, and terrific are generally considered positive and could be assigned a positive value. So, if a term like “great” is ranked at 1.5 and “horrific” at -2.0 are added then the average is calculated, the value is -0.25. By ranking all of the terms in this manner and calculating a net result, text may be inferred to be positive, negative, or even neutral in nature. This is one version of sentiment analysis.

Our basic sentiment analysis process can be seen below, along with its output.

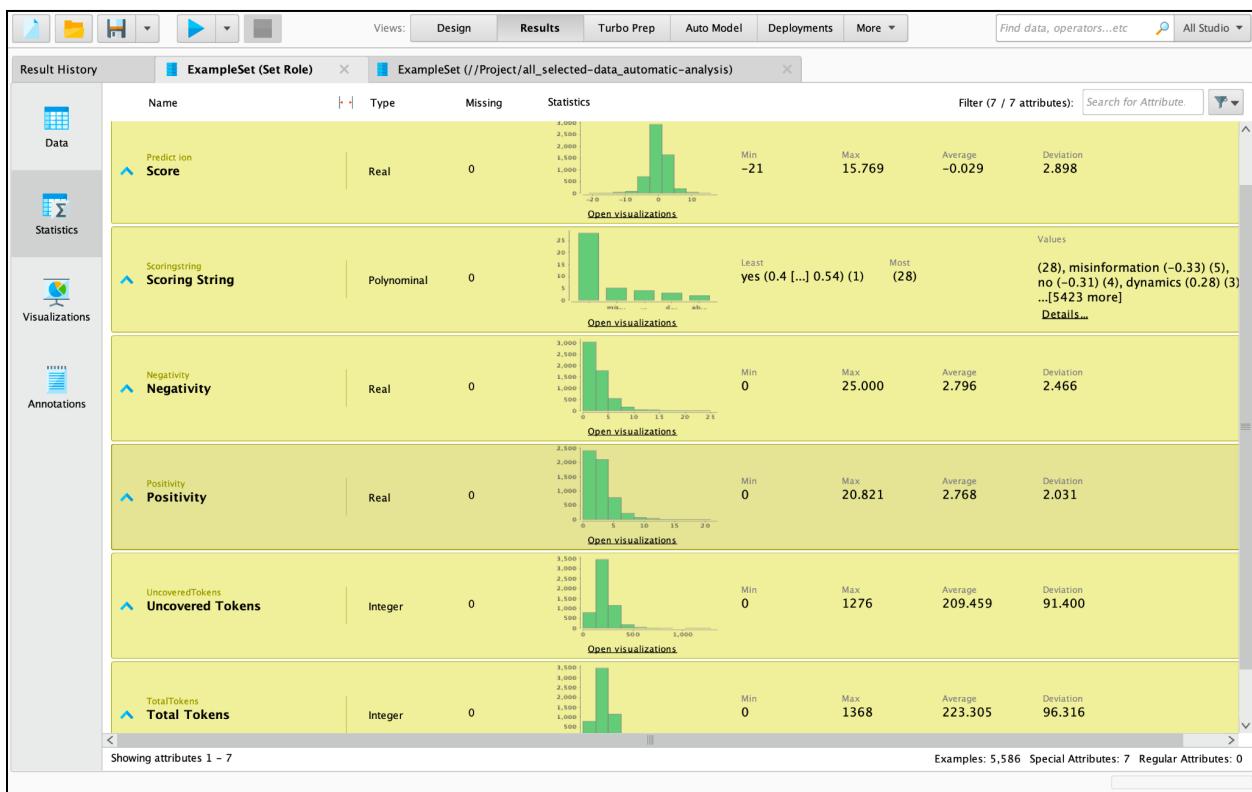


Sentiment Analysis Process

The Sentiment Analysis operator has several models to choose from, among which are Aylien, MeaningCloud, and VADER. Aylien and MeaningCloud require users to sign up and provide a connection to their services. However, VADER is the default and can run as a stand alone process. VADER is an acronym for Valence Aware Dictionary and sEntiment Reasoner, according to the RapidMiner help file. In the results below, the first abstract is scored at -2.359 versus line 5 which scored at 0.744. The Scoring String column details individual terms and their weights. Once all terms are ranked and aggregated, the Negativity and Positivity columns are revealed for each abstract. The net result is the ability to weigh unstructured text and determine negative or positive sentiment. This may ultimately help discover false

intentions and false information as users attempt to fabricate information and influence through the use of negative terminology.

ExampleSet (Set Role) ExampleSet (//Project/all_selected-data_automatic-analysis)							
	abstract	Score	Scoring String	Negativity	Positivity	Uncovered Tokens	Total Tokens
1	The following article examines the relevance of gender and inte...	-2.359	better (0.49) threats (-0.46) threats (-0.46) warfare (-...	2.846	0.487	147	155
2	Revealed with the invasion of Crimea and popularized after the ...	-0.256	popularized (0.49) propaganda (-0.26) manipulation (...	1.308	1.051	88	93
3	The rise of digital media contributes to fake news and disinfo...	-0.667	fake (-0.54) actively (0.33) fake (-0.54) argued (-0.38)...	1.872	1.205	218	226
4	Purpose This paper seeks to disambiguate the phenomenon by...	-3.692	importance (0.38) critical (-0.33) fight (-0.41) fake (-0...	6.256	2.564	225	249
5	This article explores the challenges and opportunities presente...	0.744	challenges (0.08) opportunities (0.41) intelligence (0.54)...	1.641	2.385	110	122
6	This paper argues that 'fake news' is endemic to 'information s...	-0.718	argues (-0.41) fake (-0.54) disruptions (-0.36) dynami...	2.359	1.641	104	113
7	The article presents the results of an experimental study of psy...	-1.949	errors (-0.36) misinformation (-0.33) focused (0.41) mi...	2.359	0.410	175	182
8	We live in a hyper-informed society that is constantly being fed ...	1.333	vulnerable (-0.23) weaknesses (-0.38) creation (0.28) ...	1.282	2.615	198	211
9	This paper investigates the impact of artificial intelligence (AI) o...	3.641	intelligence (0.54) ethical (0.59) opportunity (0.46) impr...	0	3.641	115	122
10	Disinformation is a serious problem for democratic systems in o...	0	serious (-0.08) problem (-0.44) vision (0.26) problem (...	1.282	1.282	221	229
11	The purpose of this study was twofold. First, this study sought t...	0.949	validate (0.38) anti (-0.33) controversial (-0.21) advanc...	0.769	1.718	141	149
12	Information disorder (satire or parody, false connections, misle...	-5.974	disorder (-0.44) misleading (-0.44) manipulated (-0.41) ...	7.051	1.077	278	296
13	In the past decade, disinformation has become an increasingly ...	-0.513	dangerous (-0.54) enemy (-0.04) risk (-0.28) validated...	1.692	1.179	90	97
14	Yes, Minister is a series that has been part of the collective ima...	1.897	yes (0.44) gained (0.41) special (0.44) yes (0.44) num...	0.949	2.846	322	332
15	The objective of this work is to analyze how rumors and hoaxes...	-0.487	problem (-0.44) reach (0.03) problem (-0.44) value (0...	0.872	0.385	141	145
16	Black lives have continually been subject to historical and conte...	-3.103	harassment (-0.64) engagement (0.51) anxieties (-0.15) ...	4.513	1.410	177	189
17	This paper analyzes the spreaded disinformation about the cor...	0.436	well (0.28) misleading (-0.44) clearly (0.44) critical (-0...	0.769	1.205	141	147
18	Those who control the information space control society's ability...	-1.538	ability (0.33) riots (-0.59) warning (-0.36) challenge (0...	2.744	1.205	120	130
19	Disinformation about the COVID-19 pandemic has reached suc...	-1.744	reached (0.10) emotional (0.15) negative (-0.69) emotio...	3.179	1.436	247	257
20	This article presents the main features of the phenomenon of di...	-0.949	false (-0.54) well (0.28) misinformation (-0.33) combat ...	1.231	0.282	95	99
21	In previous research, I outlined some of the new challenges that...	-1.308	challenges (0.08) increase (0.33) alarmed (-0.36) fake ...	2.667	1.359	184	196
22	Recent scholarship indicates that populism rhetoric can profound...	-4.359	enerov (0.28) recession (-0.46) losses (-0.44) lobby (0...	5.923	1.564	206	222



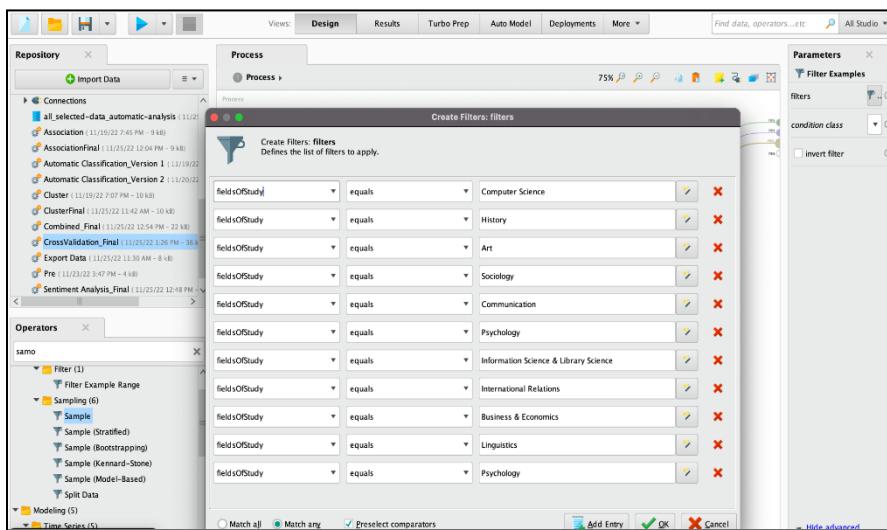
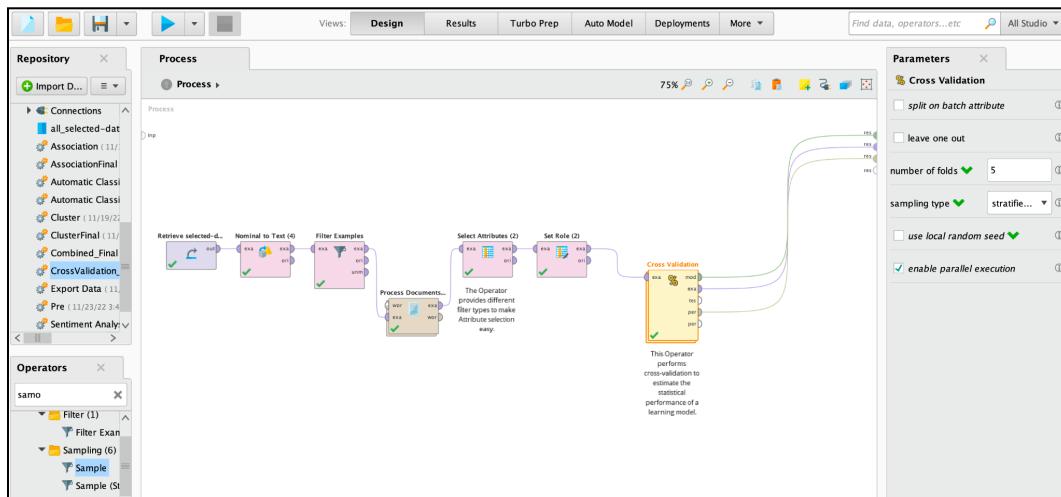
Process 4: Cross Validation and Text Classification Process

Cross Validation

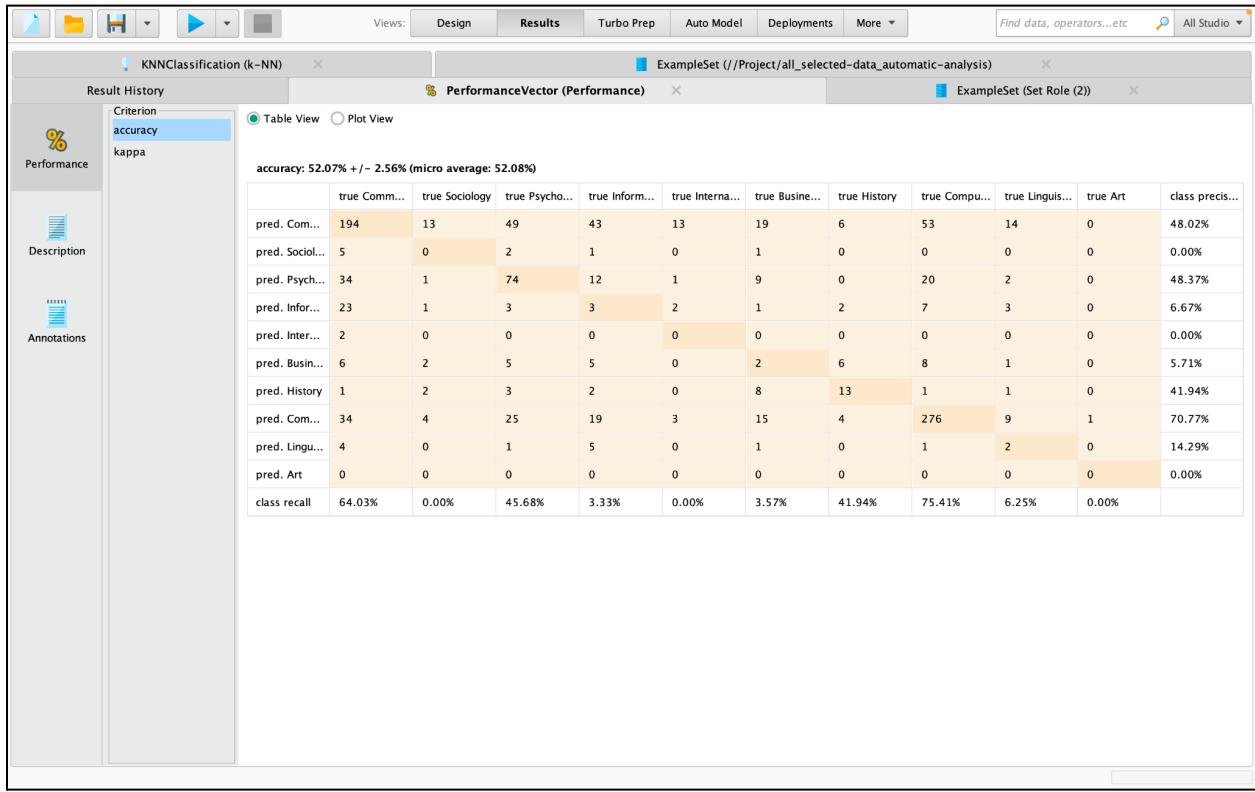
We used the cross-validation operator to check the performance of the models for training and testing. Cross-validation is usually the preferred method because it gives the model the opportunity to train on multiple train-test splits and gives the better indication of how well the model will perform on the unseen data.

From the dataset we will use attributes such as field of study and abstracts. From the abstracts we will calculate the word frequency using an algorithm called k nearest neighbour to automatically categorize the abstract into the categories. So we tried to train the algorithm to categorize the abstract from previous examples to be able to automatically categorize other ones into correct categories.

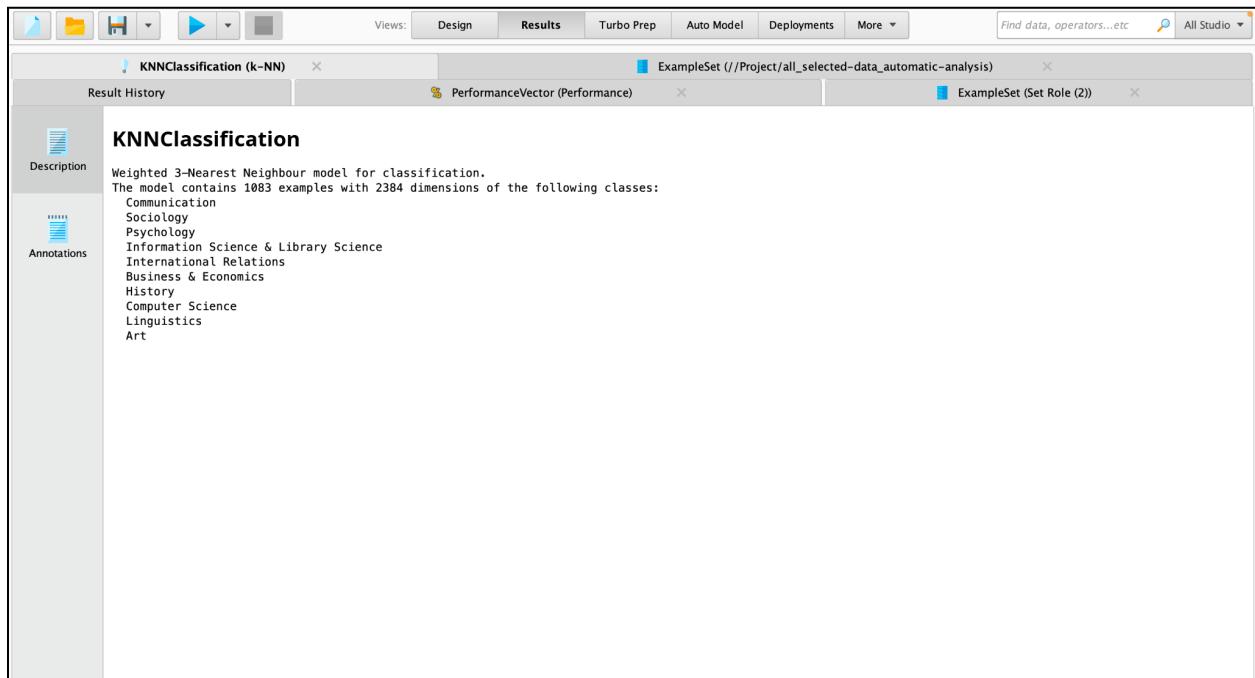
Overall, the model was able to categorize the abstracts with accuracy of 52.07%



We have used Filter Examples Operator to select only certain fields of study from the dataset. Matrix Result From Cross Validation:

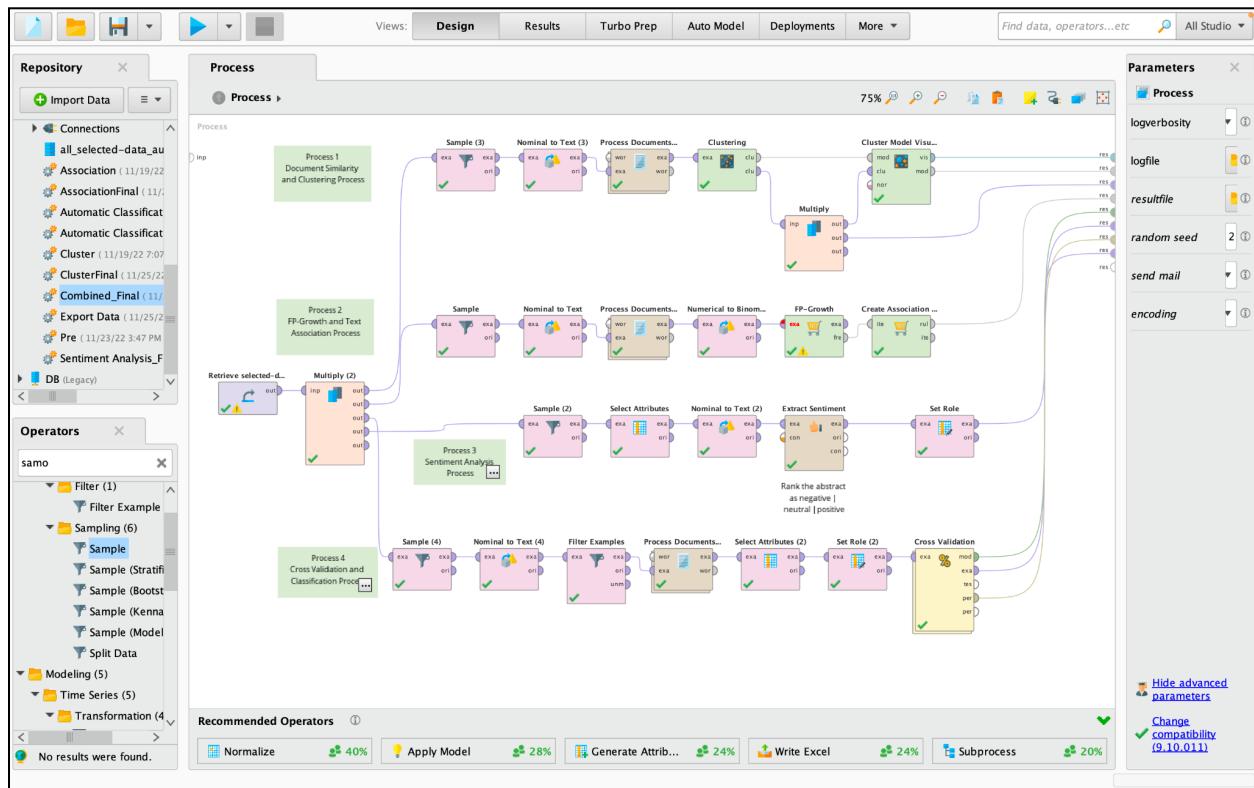


The model performed best for the field of study: Computer Science.



Final Staging

Although RapidMiner can run individual processes, it can also be used to run several subprocesses in parallel within the same process. After developing and testing three independent processes, we decided to combine them into a single process with all three processes running consecutively. In the figure above, we retrieve the data once, then use the Multiply operator to fork the data out to three subprocesses. Although the Sample operator is used three times, it could be included once and moved before the Multiply operator. However, we thought it prudent to maintain the ability to adjust the sample size for each of the subprocesses due to the resource intensive nature of these data analysis processes, especially when combined into a single process. Likewise, any process duplicated in multiple subprocesses could be consolidated into a single process for use across multiple parallel processes. For example, one candidate is the Process Document Data subprocess which is found three times. Further, complex subprocesses may be saved and grouped into Building Blocks for reuse. In this way, we are able to create complex, reusable processes, saved to our repository, are available to drag and drop, and introduce an efficient operating strategy.



Team Member Contributions

Harsh Sangani: Discussion and meeting / Literature Review / Experiment and Data Analysis Plan / Sample of Visualizations / Methodology / Data Pre-processing / Data Analysis

Sean Griffin: Discussion and meeting / Literature Review / Experiment and Data Analysis Plan / Sample of Visualizations / Methodology / Data Pre-processing / Data Analysis

Honey Rattanakasem: Discussion and meeting / Literature Review / Experiment and Data Analysis Plan / Sample of Visualizations / Methodology / Data Pre-processing / Data Analysis

Bryan Anderson: Discussion and meeting / Literature Review / Experiment and Data Analysis Plan / Sample of Visualizations / Methodology / Data Pre-processing / Data Analysis

References

- Abdeen, M. A., Hamed, A. A., Wu, X. (2021). Fighting the COVID-19 infodemic in news articles and false publications: The neonet text classifier, a supervised machine learning algorithm. *Applied Sciences* 2021, 11(16), 7265. <https://doi.org/10.3390/app11167265>
- Abdullah-All-Tanvir, Mahir, E. M., Akhter, S., & Huq, M. R. (2019). Detecting Fake News using Machine Learning and Deep Learning Algorithms. Paper presented at the 1-5. 10.1109/ICSCC.2019.8843612
- Ahmed, A. A. A., Aljabouh, A., Donepudi, P. K., & Choi, M. S. (2021). Detecting Fake News Using Machine Learning : A Systematic Literature Review. *CoRR, abs/2102.04458* <https://arxiv.org/abs/2102.04458>
- Ahmed, S., Hinkelmann, K., & Corradini, F. (2022). Development of fake news model using machine learning through natural language processing Retrieved from https://ui.adsabs.harvard.edu/abs/2022arXiv220107489A*
- Al-Rawi, A. (2019). Gatekeeping Fake News Discourses on Mainstream Media Versus Social Media. *Social Science Computer Review*, 37(6), 687-704. <https://10.1177/0894439318795849>
- Baptista, J. P., & Gradim, A. (2022). A Working Definition of Fake News. *Encyclopedia*, 2(1), 645. <https://10.3390/encyclopedia2010043>
- Chen, H., Wu, L., Chen, J., Lu, W., & Ding, J. (2022). A comparative study of automated legal text classification using random forests and deep learning. *Inf.Process.Manage.*, 59(2) doi:10.1016/j.ipm.2021.102798
- Granik, M., & Mesyura, V. (2017). Fake news detection using naive Bayes classifier. Paper presented at the 900-903. 10.1109/UKRCON.2017.8100379
- Hangloo, S., & Arora, B. (2021). Fake News Detection Tools and Methods--A Review. arXiv Preprint arXiv:2112.11185,
- Hansrajh, A., Adeliyi, T. T., Wing, J. (2021). Detection of online fake news using blending ensemble learning. *Scientific Programming*, 2021(Article ID 3434458). <https://doi.org/10.1155/2021/3434458>
- Kaliyar, R. K., Anurag, G., & Pratik, N. (2021). EchoFakeD: improving fake news detection in social media with an efficient deep neural network. *Neural Computing & Applications*, 33(14), 8597-8613. <https://doi.org/10.1007/s00521-020-05611-1>
- Lazer, D. M. J., Baum, M. A., Benkler, Y., Berinsky, A. J., Greenhill, K. M., Menczer, F., Metzger, M. J., Nyhan, B., Pennycook, G., Rothschild, D., Schudson, M., Sloman, S. A., Sunstein, C. R., Thorson, E. A., Watts, D. J., & Zittrain, J. L. (2018). The science of fake news. *Science*, 359(6380), 1094-1096. <https://10.1126/science.aaq2998>

- Liu, Y., Yu, K., Wu, X., Qing, L., & Peng, Y. (2019). Analysis and detection of health-related misinformation on Chinese social media. *IEEE Access*, 7, 154480-154489.
- Oshikawa, R., Qian, J., & Wang, W. Y. (2018). A survey on natural language processing for fake news detection. *arXiv Preprint arXiv:1811.00770*
- Ozbay, F. A., & Alatas, B. (2020). Fake news detection within online social media using supervised artificial intelligence algorithms. *Physica A: Statistical Mechanics and its Applications*, 540, 123174. <https://doi.org/10.1016/j.physa.2019.123174>
- Raja, M. S., Raj, L. A. (2022). Fake news detection on social networks using machine learning techniques. *MaterialsToday: Proceedings*, 62(7), 4821-4827. <https://doi.org/10.1016/j.matpr.2022.03.351>
- Reis, J. C. S., Correia, A., Murai, F., Veloso, A., & Benevenuto, F. (2019). Supervised Learning for Fake News Detection. *IEEE Intelligent Systems*, 34(2), 76-81. 10.1109/MIS.2019.2899143
- Usai, A., Pironti, M., Mital, M., & Aouina Mejri, C. (2018). Knowledge discovery out of text data: A systematic review via text mining. *Journal of Knowledge Management*, 22(7), 1471-1488. doi:10.1108/JKM-11-2017-0517