
Predicting the Impact of Research Papers in the Covid-19 Open Research Dataset

Tyler Ashoff

Department of Interdisciplinary Data Science
Duke University
Durham, NC 27708
`tyler.ashoff@duke.edu`

Jennifer Wilson

Department of Statistical Science
Duke University
Durham, NC 27708
`jennifer.wilson994@duke.edu`

Abstract

The volume of research papers and articles relating to Coronaviruses combined with the rapid increase at which new information is being published makes it nearly impossible for those relying on this research to keep up [1]. We propose the use of citation count prediction [2] to identify impactful articles which may have been overlooked.

1 Introduction

As the world community became engulfed in the Covid-19 pandemic in early 2020, the scientific community sprang into action. The rate at which new results and research are still being published at an alarming rate. In an effort to help identify potentially influential research, we turn to citation count prediction.

Citation count prediction is one of a number of methods for predicting the importance and popularity of a paper [2]. While numerous modeling methods have already been produced to predict this, we look at citation count prediction specifically as it relates to the Covid-19 Open Research Dataset (COORD) [1]. This dataset compiled by the Whitehouse and leading researchers provides information on over 57,000 articles, including the full text. This dataset provides some unique challenges related to the steep rise in articles recently published resulting in a dataset that is highly skewed. Additionally, the does not include information heavily relied on for citation count prediction such as the h-index [3]. Therefore we focus on extracting data from the documents and abstracts themselves in an attempt to predict citation count. Numerous different modeling techniques were attempted, including neural nets, zero-inflated poisson regression, hierarchical modeling and gradient boosting. Many of these models had little to no significance.

2 Data Wrangling

Our main data source is the COVID-19 Open Research Dataset (CORD-19) assembled by the White House and leading research groups. It is a collection of scientific articles about COVID-19, SARS-CoV-2, and related coronaviruses. Our 12372 sample subset of this dataset is chosen by first selecting observations available on the NCBI website and second ensuring that the full body of the text is available in the dataset. Unfortunately citation data was not supplied in the CORD-19, to gather

these values we scraped available citation data from the NCBI website. The information available in our compiled dataset included, journal of publication, full article text, authors, number of articles cited in the text, and number of citations received by the article. Since we are interested in understanding paper importance in terms of contribution to the field, these citations are key to our inquiry.

To transform the articles' text into usable features, each abstract and text body are run separately through gensim's python implementation of Doc2Vec. This resulted in a 100 dimensional vectorization of each text. We chose to analyze the abstracts independently because they are designed as a succinct summarization of the paper's content and could help focus the vectorization. Each of these vectorizations were then embedded into 5 dimensional space using a variational autoencoder, Isomap, and Spectral embedding. Using these embedded documents, we create three new features for the abstracts and the full text: mean Euclidean distance of 5 nearest neighbors, mean Euclidean distance of 5 nearest neighbors from the preceding six months, and a count of texts within an open ball of radius $0.1 * \text{standard deviation of all distances}$. These were designed to measure where the work sits within the topic's literature, where it sits in relation to prior work, and density of work in the specific topic respectively.

The final features used in the model are as shown below:

Features
Journal name
Publication year
Average abstract distance
Six month average abstract distance
Count of close abstracts
Average body distance
Six month average body distance
Count of close bodies

Due to heavy right tails and non-linearities, some of the variables are transformed. A log transform is applied to the average abstract distance, six month average abstract distance, average body distance, and six month average body distance.

Because of their total citation numbers. Our predicted ranking of these papers is reproduced in the appendices. Our interest is in predicting paper influence for recent papers, we hold out all papers published in 2020. In addition to being the papers of interest, we felt the recent nature of their publication may have an influence.

3 Methods

3.1 Transformations

Transformations of the response and covariates proved to have significant impact on model performance. Most gains in improving results came from transformations as opposed to tuning parameters.

3.1.1 Response

A significant number of articles in this dataset have no citations. While some of this is simply due to the fact that they have not been published long enough for others to have cited them yet, this phenomenon is still present when restricting to articles published in previous years. Additionally, there are some articles that receive a significantly large number of citations. We opt to use a $\log(x + 1)$ transformation on the response which allows us to normalize the non zero counts and further separate articles with citations from those without. While not among the final models we present, we did attempt a zero inflated poisson model without the transformation to try and deal with the high mass at zero. This model had poor results.

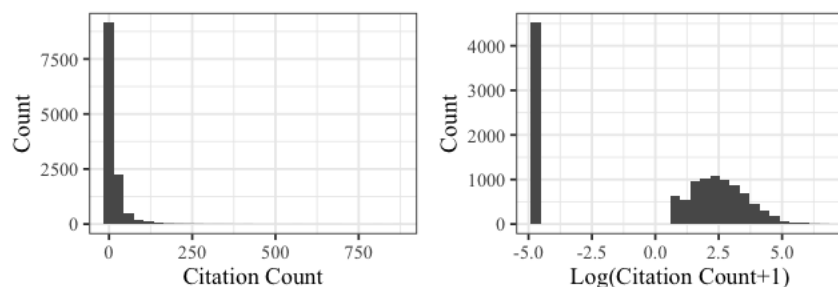


Figure 1: Transformation of Citation Count

3.1.2 Covariates

Many of the covariates from the embedding task were not linearly related to the citation counts. Most often this appeared as a high citation count and lots of density near a value of 0 with low citation counts in the positive and negative limits with low density. We use a $\log(x^2)$ transformation of these covariates. While somewhat convoluted and uninterpretable, this transformation does a good job at spreading out the distribution of the data and linearizing the relationship. Additionally, any count related covariates were given a log transformation.

4 Results

The most successful model is an xgboost model with max depth of 5, learning rate of 0.5, and objective functions squared error and hinge loss for regression and binary classification respectively. Also shown below are the results from a random forests model. In addition to these models we tested linear models, a neural network, and Gaussian process regression and classification. However, the tree based methods performed best on the training data.

The results show that while the choice of embedding had an effect on accuracy, with the variational autoencoder embedding performing best, these improvements were only marginal. Similarly, the effect of the text based features, while helpful, were less important than date of publication and journal of publication. Inspection of feature importance from the xgboost model showed that far and away the most predictive feature for citations is the journal in which the work is published.

—	AUC	Accuracy	F1
Spectral Embedding			
XGBoost	0.71	0.71	0.77
Random.Forests	0.70	0.70	0.72
Isomap Embedding			
XGBoost	0.71	0.72	0.78
Random.Forests	0.69	0.70	0.71
Variational Autoencoder			
XGBoost	0.72	0.73	0.77
Random.Forests	0.70	0.70	0.72
No Embedding			
XGBoost	0.72	0.73	0.78
Random.Forests	0.71	0.71	0.73

References

References follow the acknowledgments. Use unnumbered

[1] Allen Institute For AI. “COVID-19 Open Research Dataset Challenge (CORD-19).” Kaggle, CDC, 25 Apr. 2020, www.kaggle.com/allen-institute-for-ai/CORD-19-research-challenge.

—	Variance Explained	MSE	R2
Spectral Embedding			
XGBoost	0.40	18.36	0.63
Random.Forests	0.35	19.94	0.60
Isomap Embedding			
XGBoost	0.39	18.73	0.62
Random.Forests	0.33	20.46	0.59
Variational Autoencoder			
XGBoost	0.38	19.10	0.61
Random.Forests	0.36	19.70	0.60
No Embedding			
XGBoost	0.40	18.52	0.63
Random.Forests	0.35	19.98	0.60

[2] Yan, Rui & Tang, Jie & Liu, Xiaobing & Shan, Dongdong & Li, Xiaoming. (2011). Citation count prediction: Learning to estimate future citations for literature. Conference on Information and Knowledge Management, Proceedings 1247-1252. 10.1145/2063576.2063757.

[3] J. Hirsch. An index to quantify an individual's scientific research output. Proceedings of the National Academy of Sciences of the United States of America, 102(46):16569, 2005.