

# “How Prolific Are You?”: A Scientometric Analysis of Researchers

Rahul Yedida<sup>1</sup> and Snehanshu Saha<sup>1</sup>

<sup>1</sup>Department of Computer Science & Engineering, PESIT South Campus

## Abstract

Scholars in academia regularly write papers showing the results of their work. To gauge how relevant their work is, several measures are used, including citation count and the h-index. This paper describes the methods used to identify “prolific” authors, whose metrics are significantly higher than their peers, and provides a general discussion of the relationships between different metrics. We also discuss how this analysis is performed even with the high dimensionality of the dataset.

## 1 Introduction

Give an introduction here.

The rest of the paper is structured as follows. Section 2 describes the data and shows correlations between the attributes. Section 3 discusses the methods we used to analyse this data.

## 2 Data

The data contains 618 records of scholars with 25 attributes. There are no missing values, and eight of the attributes are integers. For ease of analysis, we have labeled the attributes as *v1* through *v25*. [maybe show what these are in a table].

### 2.1 Preliminary Analysis

For a preliminary analysis of the relationships between the attributes, we sought to discover monotonic relationships. This was done by computing pairwise Spearman’s rank correlation coefficient. This revealed that there were only positive monotonic relationships. 42.08% of these coefficients were above 0.8, and 52% were above 0.7.

A logical next step, then, was to investigate the percentage of linear relationships among these. This was done by computing pairwise Pearson correlation coefficients. To visually understand pairs that had high linear correlations, we displayed this information in a matrix, maintaining only the upper half, and avoiding diagonal entries. This avoids counting self-correlations and

Multiple regression analysis was performed as follows. For each attribute that the variable being analysed has a high correlation with, we ran an Ordinary Least Squares (OLS) regression model in two cases—with a constant term and without (that is, in a forward step-wise fashion).

In each case, for each variable, we check the p-value to ensure that all variables are statistically significant. Finally, we choose the model with the highest adjusted  $R^2$  value, which avoids the pitfalls of “kitchen sink regression”. We also make sure that there are no significant autocorrelations, using the Durbin-Watson test. [1] suggests that a conservative range for the acceptable values for this test is between 1 and 3. In some cases, we chose to pick a model that has a slightly lesser adjusted  $R^2$  value in favor of using lesser variables. This results in a model that still has a high predictive power, yet is simpler. In addition, we also check the residuals vs. fits plots and check that they are not very correlated.

The full regression results can be found in the Supplementary Materials section at the end of this paper. Here, we present a summary of the results of the regression analysis that was performed.

**Figure 1:** Matrix showing pairs with Pearson R values 0.8 or higher

These correlations provide a starting point to find relationships between the attributes.

## 3 Methods

This section discusses in detail the methods used to analyse the data.

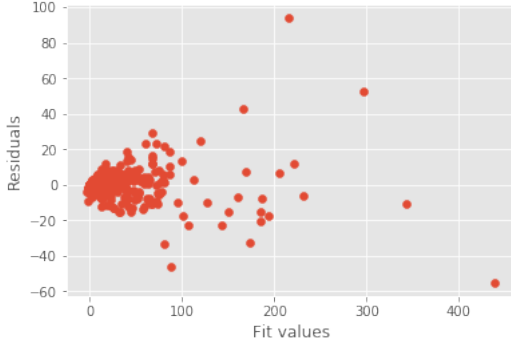
### 3.1 Regression Analysis

Our first analysis investigated the precise relationships between the attributes. We looked at linear relationships. Because v4 had the most number of correlations, we started with it.

### 3.1.1 Analysis on the full data

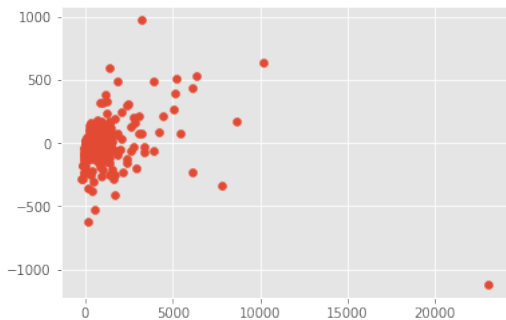
We first looked at v4, and did a step-wise regression analysis.  $v6$ ,  $v9$ ,  $v11$ ,  $v13$ ,  $v14$ ,  $v17$ ,  $v19$ ,  $v24$ , and  $v25$  were statistically significant, and the adjusted  $R^2$  was 0.972. For this model, the Durbin-Watson statistic was 2.06, which suggests very minor negative autocorrelations. The correlations between the residuals and the fitted values was -0.0047. Figure 2 shows the residuals

vs. fits plot for  $v_4$ .



**Figure 2:** Residuals vs. fits plot for  $v_4$

Our next variable of interest was  $v_2$ , having high Pearson R correlations with nine other variables. Only  $v_9$  was not statistically significant. However, we also chose to discard  $v_{24}$  because it added no predictive power to the model, and the adjusted  $R^2$  for this model was 0.993, and the Durbin-Watson statistic was 1.853, suggesting only minor positive autocorrelations. All the p-values were less than  $10^{-3}$ . The correlation between the residuals and the fits was 0.023. Figure 3 shows the residuals vs. fits plot for  $v_2$ .



**Figure 3:** Residuals vs. fits plot for  $v_2$

We then looked at  $v_6$ , also having high

correlations with nine other variables.  $v_9$  and  $v_{14}$  were not statistically significant. We also discarded  $v_{25}$  to obtain a simpler model with an adjusted  $R^2$  value 0.001 less. This model had an adjusted  $R^2$  of 0.913, and a Durbin-Watson statistic of 1.64. All p-values were below  $10^{-3}$ .

Next, we looked at  $v_9$ . Only  $v_{13}$  was not statistically significant, but we also discarded  $v_{10}$  and  $v_{11}$ , which resulted in a model whose adjusted  $R^2$  was 0.986 (a reduction of 0.002 from when the two variables were included). This model had a Durbin-Watson statistic of 1.99.

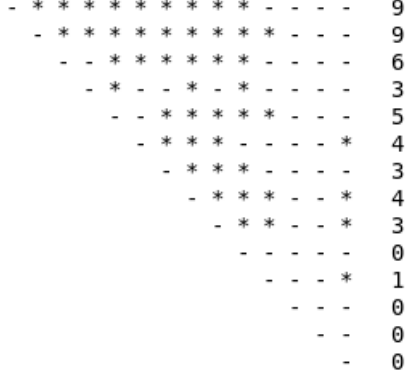
Finally, for  $v_{10}$ , all the predictor variables were statistically significant, and we discarded  $v_{19}$ , which resulted in a model with the same adjusted  $R^2$ . All the p-values were below  $10^{-3}$ , except  $v_{13}$ , whose p-value was 0.043 (95% C.I [0.07, 0.425]). The Durbin-Watson statistic was 1.861.

It should be noted that in all of the regression results of the above variables, the constant was not statistically significant, and the results reported are those for the models that did not include the constant.

### 3.1.2 Analysis on subset of the data

Our next step of analysis was to pick a subset of the variables, and redo the same analysis discussed above on this subset. Once again, it turns out that for all the regression equations obtained, the constant was not statistically significant.

The variables we chose were  $v_9$  through  $v_{21}$  and  $v_{24}$ . maybe explain why these. Again, we created a matrix of Pearson R correlations, and picked out those that were higher than 0.8. Figure 4 shows this result.



**Figure 4:** Pearson R correlation matrix for subset of the data

Using this, we redid the same analysis as in the previous subsection. Again, we only discuss the results here, but the full list of experiments is in the Supplementary Materials section.

We first looked at  $v9$ .  $v13$  was not significant when all the variables it was correlated to were considered. However, when  $v10$  was not considered, it was statistically significant. We chose to sacrifice 0.2% adjusted  $R^2$  in favor of a model that had two fewer predictor variables. The Durbin-Watson value was 1.99.

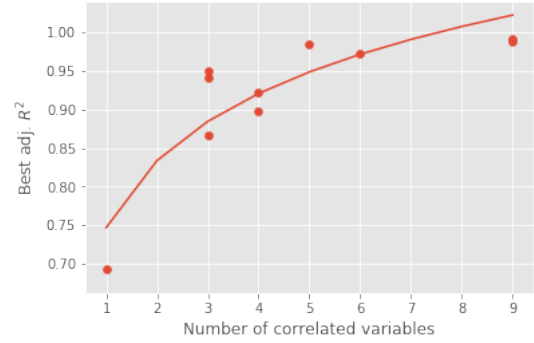
While predicting  $v10$ , dropping  $v19$  resulted in no loss of adjusted  $R^2$ . We further dropped  $v18$ , and dropped  $v12$  because it was not significant. While dropping  $v18$  caused a 0.3% drop in adjusted  $R^2$ , we were able to use two less variables. The Durbin-Watson statistic was 1.991.

We only looked at two other variables, because the adjusted  $R^2$  for the others was below 0.95. For  $v11$ , dropping  $v17$  and  $v18$  resulted in a 0.2% drop in adjusted  $R^2$ , and we were able to use only four variables to

get an adjusted  $R^2$  of 0.971 and a Durbin-Watson of 2.081. Finally, for  $v13$ , all variables were statistically significant, and we discarded  $v19$  with no loss of adjusted  $R^2$ .

Perhaps more interesting than the above results was when we plotted the maximum adjusted  $R^2$  (without discarding any variables) against the number of predictor variables used (on this subset of the data). Figure 5 shows this result, where a logarithm function fits quite well to the data. The equation obtained was

$$\text{Adj. } R^2 = 0.125465 \log n + 0.746654 \quad (1)$$



**Figure 5:** Adjusted  $R^2$  vs. number of predictor variables

The sum of the squared residuals for this fit was 0.01456.

## References

- [1] A. Field, *Discovering statistics using SPSS*. Sage publications, 2009.