
Bibliometrics: Predicting Publication Success

Team 4:
Tamanna Baig and Jon Oakley

Prepared for: Dr. Wang
CPSC 8650
April 15th, 2019

Abstract—Determining the impact of researchers has always fascinated scholars. Citation count is one of the most primitive bibliometrics that has been used to assess the success of a scientific publication. We presented an approach to predicting the success of an academic publication based on the features described in Table I. We reduced our initial corpus of 22,588 documents to 4,914 documents by removing any record that did not contain the required features. These documents were clustered into 27 clusters, and one cluster was omitted due to its size. Five of these clusters were chosen for analysis. The papers were marked as successful and unsuccessful based on whether or not they had more (successful) or less (unsuccessful) citations than the median number of citations for that cluster. The numeric features (Table I) were classified using a Random Forest classifier, and the text-based features (title, abstract, and keywords) were classified using three independent Multinomial Naive Bayes classifier. The predictions from these four classifiers were combined using an ensemble Random Forest classifier. Through this methodology, we were able to achieve an overall classification accuracy of 66.5% for the optimal clusters, and an overall classification accuracy of 57.7% for all of the clusters.

Index Terms—Bibliometrics, Prediction, Citation Count, Success

I. INTRODUCTION

Determining the impact of researchers has always fascinated scholars. Citation count is one of the most primitive bibliometrics that has been used to assess the success of a scientific publication. The h -index incorporates this idea and scores researchers based on the number of highly cited papers they have [1]. In [2] the h -index of a researcher was predicted using different attributes from their curriculum vitae (CV). Other predictors have used Twitter to predict the h -index change in authors with an 80% success rate [3]. However, it has been noted that there are many biases in the h -index, such as the age of the researcher [4].

An overview of recent bibliometric analyses is presented in [5]. A model for the long-term impact of a publication is presented by [6]. A study of accumulated advantage in academic endeavors was presented in [7]. By analyzing relationships of authors using a co-authorship network, [8] were able to predict future success of authors. Collaboration relationships were also found to be a predictor of success in [9]. Scholar trajectory was mapped in [10] in order to predict researchers migrating to different institutions. In [11] author characteristics were used to predict the success of academic publications.

Other studies used paper characteristics to predict their success. In [12], it was found that shorter titles were a positive predictor of success. Other research found that declarative titles didn't bias readers [13]. Author count and abstract complexity was also identified as a positive predictor [14]. Abstract complexity as a predictor was corroborated by [15]. Keywords were also identified as a predictor in [16]. Alternatively, full-text features were also useful in predicting the success of publications [17].

Figure 1 shows how the following sections fit together. Section II details the steps taken to clean and cluster the dataset. Section III discusses the visualizations used to inform our prediction methodology. Section IV discusses the

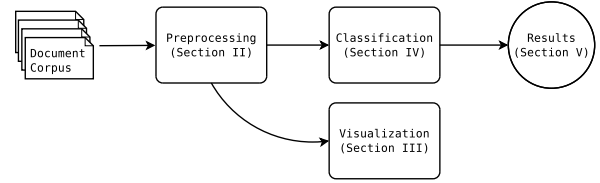


Fig. 1. Graphical abstract illustrating the tasks performed.

TABLE I
FEATURES USED FOR PREDICTING THE SUCCESS OF A PUBLICATION.

Feature	Type	Meta?
title	str	✗
abstract	str	✗
keywords	str	✗
title_length	int	✓
abstract_length	int	✓
keyword_length	int	✓
year	int	✗
author_count	int	✓
page_count	int	✓
n_citation	int	✗

various classification approaches used to predict the success of an academic publication, and Section V shows the results achieved with those classifier. Section VI provides closing thoughts and discusses future work. Code and supplemental results are provided in the Appendix.

II. PREPROCESSING

The preprocessing stage involved cleaning the data and making it usable for our analyses. The original data set was sparse—many papers had missing fields, and there was no guarantee that a paper contained every field in the complete set of fields, shown in Table VIII in the Appendix. We selected the fields (and meta-fields) shown in Table I based on availability and the existing literature discussed in Section I.

Processing the original dataset was broken down into two sub-tasks: cleaning and clustering the dataset. The process used to clean and reduce the original dataset is discussed in Section II-A, and the process used to cluster relevant papers is described in Section II-B.

A. Cleaning

The cleaning processes is depicted in Figure 2. All existing features were extracted from the document corpus. The original document corpus was a Python dictionary, so Python was used to extract and process the papers. These features were analyzed for the relevance to success using the existing literature. The features shown in Table I were chosen, and all papers that did not contain those attributes were dropped. The initial corpus contained 22,588 papers, but only 4,914 papers contained all of the relevant features. Once the dataset contained only papers that had all the relevant

features, the meta-features were extracted. We define meta-features as features that depend on one or more features. For instance, `page_count` is considered a meta-feature since it depends on `page_start` and `page_end`. The inclusion of some meta-features reduced the number of documents in the final corpus more than regular features since they required several features to be present (otherwise the paper would be removed from the corpus). Other meta-features only depended on features already required (such as `title_length`). The final corpus was again stored in a Python dictionary for clustering. The script to perform this task is shown in Listing 1.

B. Clustering

In order to ensure our results were meaningful, we used k-means clustering to group similar papers for classification. There were two primary challenges: 1) determining the appropriate field to cluster based on, and 2) determining the appropriate value of K . Empirically, we found the `abstract` was the most effective field to use for clustering. In order to chose an appropriate value of K , we iterated over a number of possible values and chose the most appropriate based on the *inertia* (the sum-squared error in each of the centroids). The process for clustering the papers is shown in Figure 3. Python’s `sklearn` machine learning library was used for all data processing in this work. The `TfidfVectorizer` was used to compute the TF-IDF value for each abstract, and the `KMeans` clustering function was used to find the clusters given the TF-IDF values.

After iterating over a number of values for K , we plotted the values, as shown in Figure 4. The optimal K was chosen because it occurred before the jitter occurred in the inertia values. Additionally, only one cluster needed to be excluded due to size. Cluster 0 contained only 15 papers and was therefore excluded.

With an appropriate K chosen, the papers were clustered according to the process depicted in Figure 3. The complete list of clusters is shown in the Appendix in Table IX. In this work, we consider the five clusters in Table II. Based on the keywords, we can see that these clusters are relatively compact. We also present the overall accuracy across all classes for comparison, but note that Cluster 0 was omitted due to insufficient sample size. The script to perform this task is shown in Listing 2.

III. VISUALIZATION

Visualization is the graphical representation of information and data. By using visual elements like graphs and charts, data visualization tools provide an accessible way to see and understand any trends, outliers or patterns present in the dataset.

For this project, various features like total number of the publications each year, authors with most publications, the papers to get the most citations and the institutions with the most affiliations are used to show any trends and patterns in the form of bar charts.

TABLE II
CLUSTERS AND A SAMPLE OF THEIR RESPECTIVE KEYWORDS.

Cluster	Count	Keywords
4	106	video, videos, 3D, quality, streams, users
15	170	query, XML, search, data, databases
17	869	design, people, user, information, research
20	93	internet, TCP, network, protocol, congestion
22	60	privacy, private, data, information, awareness

Figure 5 shows the total number of papers published each year. As depicted by the bars, it appears to be normally distributed, with 2010 seeing the highest publications. The number of publications by an author also play an important role in the success of an academic publication. Keeping this in mind the plot in Figure 6 depicts the Top 10 authors from the dataset with the highest number of publications. Another direct factor to predict the success is the citation count, and Figure 7 shows the highest number of citations received by a paper each year. Although the assumption would be that earlier the papers more the number of citations, but the plot tends to shed light that the year is not supposed to be the only factor in predicting so. Lastly, Figure 8 represents the top ten institutions/organizations in the world with highest number of affiliations, Carnegie Mellon University leading the line with around 874 affiliations.

IV. CLASSIFICATION

We break each cluster into two classes: unsuccessful and successful. We define a successful publication any publication that receives more than the median number of citations for that cluster. There are two important points to this: 1) the use of median as the average allows us to ensure the classes are balanced, and 2) we use the median relative to the cluster to control for some clusters having inherently more citations (due to popularity). The classification workflow is shown in Figure 9.

From Table I, we see that there are three string-based fields, and the rest are numeric. We use a Random Forest Classifier for the numeric fields, and we use a Multinomial Naive Bayes classifier for the text-based fields (title, abstract, and keywords). We use an independent classifier for each text-based field, so there is a title classifier, abstract classifier, and keyword classifier. In order to reach a consensus among the classifiers, we use another Random Forest classifier as an ensemble classifier. We give this classifier the prediction from all the other classifiers, and let it predict the resulting class. We train and optimize the classifiers for each cluster, so the parameters are not generalizable and are specific only to the designated cluster.

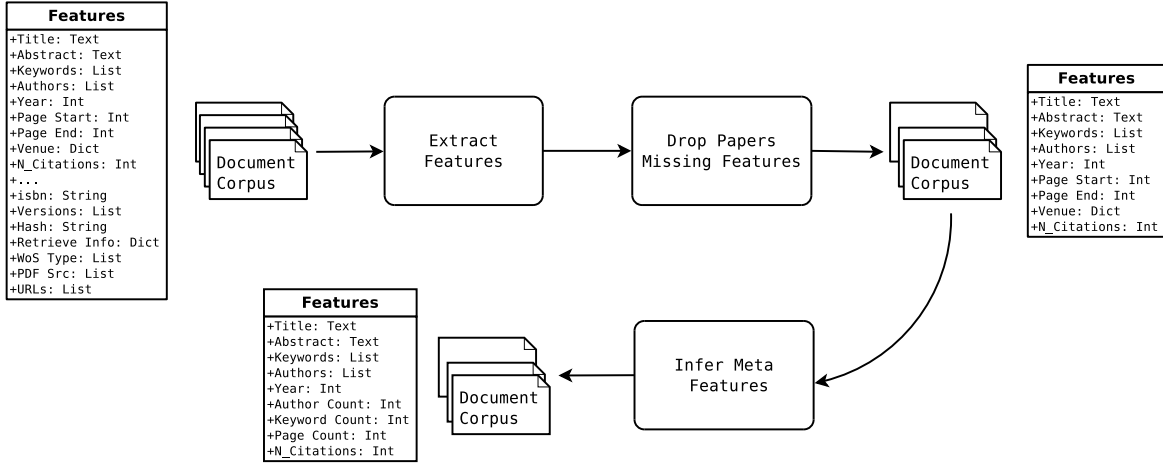


Fig. 2. Process used to clean the dataset.

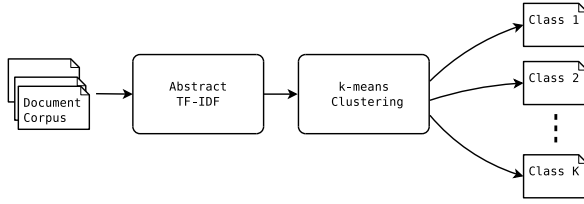


Fig. 3. Process used to cluster the papers.

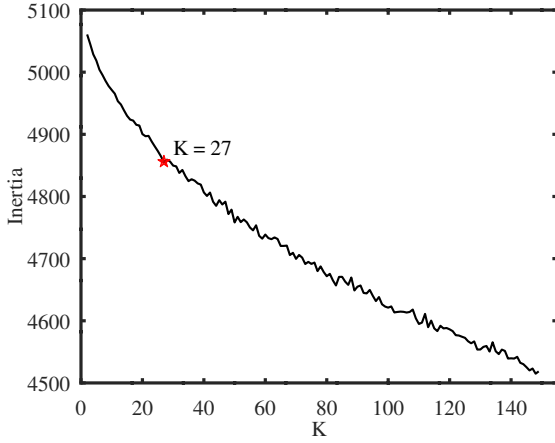


Fig. 4. Centroid inertia values for different K .

For evaluation purposes, we consider accuracy as the primary metric. We define accuracy as shown in Equation (1). This effectively measures the percent of correctly predicted values (true positive plus true negative divided by the sum of all occurrences).

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (1)$$

A. Training Method

Figure 10 shows our training methodology. For the text-based data, we consider the term frequency inverse document

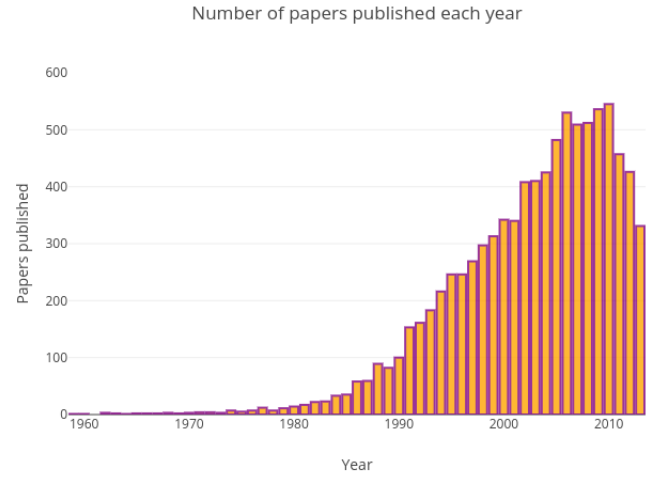


Fig. 5. Number of papers published each year.

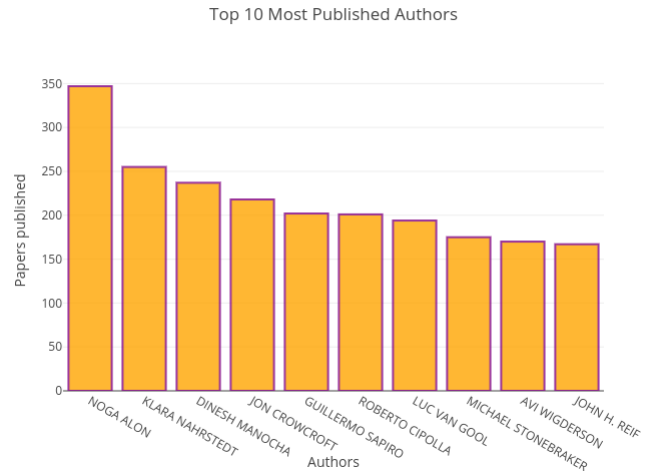


Fig. 6. Most prolific authors from the dataset.

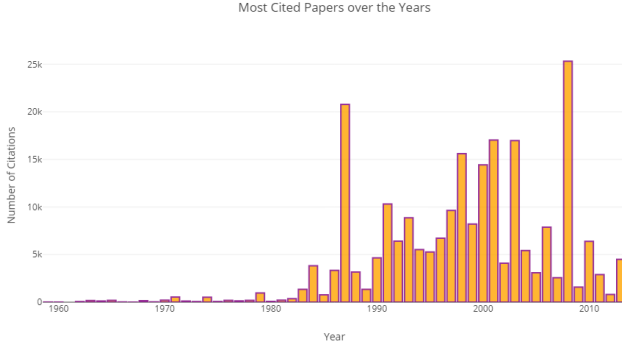


Fig. 7. Most citations of a paper published in a given year.

TABLE III
PARAMETERS FOR THE NUMERIC RANDOM FOREST CLASSIFIER.

Cluster	Numeric Classifier Parameters	
	n_estimators	max_features
4	10	2
15	10	6
17	8	2
20	12	2
22	6	4

frequency (TF-IDF) with stop-word removed. This is obtained using the `TfidfVectorizer` from `sklearn`. The numeric attributes are scaled using the `StandardScaler` from the `sklearn` library.

Before training the classifiers, the data from the clusters was formatted into a form that could be used by the classifier. The script to perform this task is shown in Listing 3. The data are split using the `test_train_split` function from `sklearn`, and the data are randomly split in half (50%). As shown in Figure 10, the training half of the data are used to train and optimize each classifier. The optimal parameters for the Random Forest classifier are found using the `GridSearchCV` function from `sklearn` and iterating over the following parameters: `n_estimators` and `max_features`. A five-fold cross validation is used with the grid search to ensure the parameters are in fact optimal. Once the classifiers are trained, the predicted values for the training data are used to train the ensemble classifier. As previously mentioned, this is also a Random Forest classifier, so the same `GridSearchCV` function is used to find the optimal parameters. The script to perform this task is shown in Listing 4.

V. RESULTS

The optimal parameters for the numeric Random Forest classifier are shown in Table III, and the optimal parameters for the ensemble Random Forest classifier are shown in Table IV.

With both classifiers, we can see there is some variation in the optimal parameters between clusters, emphasizing the importance of training the classifiers for each cluster independently. After training the Random Forest classifiers, we can see

TABLE IV
PARAMETERS FOR THE ENSEMBLE RANDOM FOREST CLASSIFIER.

Cluster	Ensemble Classifier Parameters	
	n_estimators	max_features
4	1	1
15	4	1
17	1	1
20	1	1
22	1	1

TABLE V
FEATURE WEIGHTS FOR THE NUMERIC RANDOM FOREST CLASSIFIER.

Cluster	Numeric Feature Weight					
	Title (length)	Abstract (length)	Keyword (length)	Year	Author Count	Page Count
4	0.089	0.117	0.143	0.240	0.188	0.223
15	0.160	0.157	0.121	0.321	0.065	0.175
17	0.136	0.253	0.174	0.183	0.089	0.165
20	0.087	0.214	0.121	0.242	0.109	0.228
22	0.140	0.038	0.066	0.556	0.155	0.045

how the feature importance weights the classification. Table V shows the weights for the numeric Random Forest classifier, and Table VI shows the weights for the ensemble Random Forest classifier.

Again, we can see the feature weights change between each cluster for both classifiers, emphasizing the importance of individual training. It is worth noting that year is the most important feature in Clusters 4, 15, 20, and 22, while abstract length is the most important feature in Cluster 17. It is also worth noting that in most cases, the numeric classifier is a sufficient predictor, but in some cases, there is noticeable improvement by including the other classifiers. The classification accuracy for each cluster are shown in Table VII.

The most important column in Table VII is the ensemble accuracy, as that is the effective accuracy for the cluster. As we can see, each of these clusters had an overall accuracy greater than 60%, but in all cases, there is room for improvement. As a final metric, we also consider the overall accuracy. We calculate this accuracy by summing the true positives, true negative, false positive, and false negatives from the ensemble classifiers for each cluster. This gives us an effective accuracy

TABLE VI
FEATURE WEIGHTS FOR THE ENSEMBLE RANDOM FOREST CLASSIFIER.

Cluster	Ensemble Feature Weight			
	Numeric Classifier	Title Classifier	Abstract Classifier	Keyword Classifier
4	1	0	0	0
15	0.420	0.271	0.059	0.250
17	0.917	0	0.083	0
20	1	0	0	0
22	1	0	0	0

A bar chart titled 'Number of affiliations' on the y-axis and 'Organization/Institutions' on the x-axis. The y-axis ranges from 0 to 900 in increments of 100. The x-axis lists ten organizations/institutions. The bars are blue with black outlines. The data is as follows:

Organization/Institutions	Number of affiliations
CMU	880
Stanford	490
UW, Seattle	440
Microsoft	440
UIUC	375
UNC, Chapel Hill	365
UC, UK	300
MIT	280
Cornell	275
UC, Berkeley	265

```
graph LR; A[Document Corpus] --> B[Preprocessing]; B --> C[Classification]; C --> D((Results)); C --> E[Numeric Classification (Random Forest)]; C --> F[Title Classification (Multinomial NB)]; C --> G[Abstract Classification (Multinomial NB)]; C --> H[Keyword Classification (Multinomial NB)]; E --> I[Ensemble (Random Forest)]; F --> I; G --> I; H --> I; I --> D;
```

The flowchart illustrates the ensemble classification process. It begins with a 'Document Corpus' (represented by a stack of papers) which undergoes 'Preprocessing'. The output of preprocessing is then fed into a 'Classification' step. From 'Classification', the data is distributed to four parallel classifiers: 'Numeric Classification (Random Forest)', 'Title Classification (Multinomial NB)', 'Abstract Classification (Multinomial NB)', and 'Keyword Classification (Multinomial NB)'. The outputs of these four classifiers are then combined in an 'Ensemble (Random Forest)' step. Finally, the ensemble's output leads to the 'Results' (represented by a circle).

```

graph LR
    ND[Numeric Data] -- Train --> NDC[Numeric Classification<br/>(Random Forest)]
    ND -- Test --> T1(( ))
    T1 --> T2((+))
    ND -- Test --> T3((+))
    T2 --> E[Ensemble<br/>(Random Forest)]
    T3 --> E
    
    T1 --> T4((+))
    T4 --> E
    
    T1 --> T5((+))
    T5 --> E
    
    T1 --> T6((+))
    T6 --> E
    
    T1 --> T7((+))
    T7 --> E
    
    T1 --> T8((+))
    T8 --> E
    
    T1 --> T9((+))
    T9 --> E
    
    T1 --> T10((+))
    T10 --> E
    
    T1 --> T11((+))
    T11 --> E
    
    T1 --> T12((+))
    T12 --> E
    
    T1 --> T13((+))
    T13 --> E
    
    T1 --> T14((+))
    T14 --> E
    
    T1 --> T15((+))
    T15 --> E
    
    T1 --> T16((+))
    T16 --> E
    
    T1 --> T17((+))
    T17 --> E
    
    T1 --> T18((+))
    T18 --> E
    
    T1 --> T19((+))
    T19 --> E
    
    T1 --> T20((+))
    T20 --> E
    
    T1 --> T21((+))
    T21 --> E
    
    T1 --> T22((+))
    T22 --> E
    
    T1 --> T23((+))
    T23 --> E
    
    T1 --> T24((+))
    T24 --> E
    
    T1 --> T25((+))
    T25 --> E
    
    T1 --> T26((+))
    T26 --> E
    
    T1 --> T27((+))
    T27 --> E
    
    T1 --> T28((+))
    T28 --> E
    
    T1 --> T29((+))
    T29 --> E
    
    T1 --> T30((+))
    T30 --> E
    
    T1 --> T31((+))
    T31 --> E
    
    T1 --> T32((+))
    T32 --> E
    
    T1 --> T33((+))
    T33 --> E
    
    T1 --> T34((+))
    T34 --> E
    
    T1 --> T35((+))
    T35 --> E
    
    T1 --> T36((+))
    T36 --> E
    
    T1 --> T37((+))
    T37 --> E
    
    T1 --> T38((+))
    T38 --> E
    
    T1 --> T39((+))
    T39 --> E
    
    T1 --> T40((+))
    T40 --> E
    
    T1 --> T41((+))
    T41 --> E
    
    T1 --> T42((+))
    T42 --> E
    
    T1 --> T43((+))
    T43 --> E
    
    T1 --> T44((+))
    T44 --> E
    
    T1 --> T45((+))
    T45 --> E
    
    T1 --> T46((+))
    T46 --> E
    
    T1 --> T47((+))
    T47 --> E
    
    T1 --> T48((+))
    T48 --> E
    
    T1 --> T49((+))
    T49 --> E
    
    T1 --> T50((+))
    T50 --> E
    
    T1 --> T51((+))
    T51 --> E
    
    T1 --> T52((+))
    T52 --> E
    
    T1 --> T53((+))
    T53 --> E
    
    T1 --> T54((+))
    T54 --> E
    
    T1 --> T55((+))
    T55 --> E
    
    T1 --> T56((+))
    T56 --> E
    
    T1 --> T57((+))
    T57 --> E
    
    T1 --> T58((+))
    T58 --> E
    
    T1 --> T59((+))
    T59 --> E
    
    T1 --> T60((+))
    T60 --> E
    
    T1 --> T61((+))
    T61 --> E
    
    T1 --> T62((+))
    T62 --> E
    
    T1 --> T63((+))
    T63 --> E
    
    T1 --> T64((+))
    T64 --> E
    
    T1 --> T65((+))
    T65 --> E
    
    T1 --> T66((+))
    T66 --> E
    
    T1 --> T67((+))
    T67 --> E
    
    T1 --> T68((+))
    T68 --> E
    
    T1 --> T69((+))
    T69 --> E
    
    T1 --> T70((+))
    T70 --> E
    
    T1 --> T71((+))
    T71 --> E
    
    T1 --> T72((+))
    T72 --> E
    
    T1 --> T73((+))
    T73 --> E
    
    T1 --> T74((+))
    T74 --> E
    
    T1 --> T75((+))
    T75 --> E
    
    T1 --> T76((+))
    T76 --> E
    
    T1 --> T77((+))
    T77 --> E
    
    T1 --> T78((+))
    T78 --> E
    
    T1 --> T79((+))
    T79 --> E
    
    T1 --> T80((+))
    T80 --> E
    
    T1 --> T81((+))
    T81 --> E
    
    T1 --> T82((+))
    T82 --> E
    
    T1 --> T83((+))
    T83 --> E
    
    T1 --> T84((+))
    T84 --> E
    
    T1 --> T85((+))
    T85 --> E
    
    T1 --> T86((+))
    T86 --> E
    
    T1 --> T87((+))
    T87 --> E
    
    T1 --> T88((+))
    T88 --> E
    
    T1 --> T89((+))
    T89 --> E
    
    T1 --> T90((+))
    T90 --> E
    
    T1 --> T91((+))
    T91 --> E
    
    T1 --> T92((+))
    T92 --> E
    
    T1 --> T93((+))
    T93 --> E
    
    T1 --> T94((+))
    T94 --> E
    
    T1 --> T95((+))
    T95 --> E
    
    T1 --> T96((+))
    T96 --> E
    
    T1 --> T97((+))
    T97 --> E
    
    T1 --> T98((+))
    T98 --> E
    
    T1 --> T99((+))
    T99 --> E
    
    T1 --> T100((+))
    T100 --> E
    
    T1 --> T101((+))
    T101 --> E
    
    T1 --> T102((+))
    T102 --> E
    
    T1 --> T103((+))
    T103 --> E
    
    T1 --> T104((+))
    T104 --> E
    
    T1 --> T105((+))
    T105 --> E
    
    T1 --> T106((+))
    T106 --> E
    
    T1 --> T107((+))
    T107 --> E
    
    T1 --> T108((+))
    T108 --> E
    
    T1 --> T109((+))
    T109 --> E
    
    T1 --> T110((+))
    T110 --> E
    
    T1 --> T111((+))
    T111 --> E
    
    T1 --> T112((+))
    T112 --> E
    
    T1 --> T113((+))
    T113 --> E
    
    T1 --> T114((+))
    T114 --> E
    
    T1 --> T115((+))
    T115 --> E
    
    T1 --> T116((+))
    T116 --> E
    
    T1 --> T117((+))
    T117 --> E
    
    T1 --> T118((+))
    T118 --> E
    
    T1 --> T119((+))
    T119 --> E
    
    T1 --> T120((+))
    T120 --> E
    
    T1 --> T121((+))
    T121 --> E
    
    T1 --> T122((+))
    T122 --> E
    
    T1 --> T123((+))
    T123 --> E
    
    T1 --> T124((+))
    T124 --> E
    
    T1 --> T125((+))
    T125 --> E
    
    T1 --> T126((+))
    T126 --> E
    
    T1 --> T127((+))
    T127 --> E
    
    T1 --> T128((+))
    T128 --> E
    
    T1 --> T129((+))
    T129 --> E
    
    T1 --> T130((+))
    T130 --> E
    
    T1 --> T131((+))
    T131 --> E
    
    T1 --> T132((+))
    T132 --> E
    
    T1 --> T133((+))
    T133 --> E
    
    T1 --> T134((+))
    T134 --> E
    
    T1 --> T135((+))
    T135 --> E
    
    T1 --> T136((+))
    T136 --> E
    
    T1 --> T137((+))
    T137 --> E
    
    T1 --> T138((+))
    T138 --> E
    
    T1 --> T139((+))
    T139 --> E
    
    T1 --> T140((+))
    T140 --> E
    
    T1 --> T141((+))
    T141 --> E
    
    T1 --> T142((+))
    T142 --> E
    
    T1 --> T143((+))
    T143 --> E
    
    T1 --> T144((+))
    T144 --> E
    
    T1 --> T145((+))
    T145 --> E
    
    T1 --> T146((+))
    T146 --> E
    
    T1 --> T147((+))
    T147 --> E
    
    T1 --> T148((+))
    T148 --> E
    
    T1 --> T149((+))
    T149 --> E
    
    T1 --> T150((+))
    T150 --> E
    
    T1 --> T151((+))
    T151 --> E
    
    T1 --> T152((+))
    T152 --> E
    
    T1 --> T153((+))
    T153 --> E
    
    T1 --> T154((+))
    T154 --> E
    
    T1 --> T155((+))
    T155 --> E
    
    T1 --> T156((+))
    T156 --> E
    
    T1 --> T157((+))
    T157 --> E
    
    T1 --> T158((+))
    T158 --> E
    
    T1 --
```

TABLE VII
CLASSIFICATION ACCURACY FOR EACH CLUSTER.

	Accuracy (%)				
Cluster	Numeric	Title	Abstract	Keywords	Ensemble
4	60.4	45.8	56.3	58.3	60.4
15	63.0	64.2	65.4	55.6	61.7
17	64.2	59.1	59.6	62.5	68.5
20	64.3	42.9	35.7	38.1	64.3
22	64.3	0.5	0.5	0.5	64.3

VI. CONCLUSIONS

We presented an approach to predicting the success of an academic publication based on the features described in Table I. We reduced our initial corpus of 22,588 documents to 4,914 documents by removing any record that did not contain the required features. These documents were clustered into 27 clusters, and one cluster was omitted due to its size. Five of these clusters were chosen for analysis. The papers were marked as successful and unsuccessful based on whether or not they had more (successful) or less (unsuccessful) citations than the median number of citations for that cluster. The numeric features (Table I) were classified using a Random Forest classifier, and the text-based features (title, abstract, and keywords) were classified using three independent Multinomial Naive Bayes classifier. The predictions from these four classifiers were combined using an ensemble Random Forest classifier. Through this methodology, we were able to achieve an overall classification accuracy of 66.5% for the optimal clusters, and an overall classification accuracy of 57.7% for all of the clusters.

We expect this accuracy to drastically improve with a better dataset. We were limited to the available data set, which contained a wide array of papers with no particular focus. However, it is likely these results would be generalizable to a particular journal, given the paper fell within the same cluster. Future work will investigate alternative clustering algorithms and ways to incorporate title and keywords into the clustering. Future work will also consider additional numeric metrics, authors, and institutions.

Our accuracy for the optimal clusters is promising. While there is still room for improvement, we show our classification algorithm can accurately predict every two out of three successful publications, which offers realistic insight into the future success of an academic publication.

REFERENCES

- [1] J. E. Hirsch, "An index to quantify an individual's scientific research output," *Proceedings of the National academy of Sciences*, vol. 102, no. 46, pp. 16 569–16 572, 2005.
- [2] D. E. Acuna, S. Allesina, and K. P. Kording, "Future impact: Predicting scientific success," *Nature*, vol. 489, no. 7415, p. 201, 2012.
- [3] Z. Luo, J. Chen, and X. Liu, "Real-time scientific impact prediction in twitter," in *CCF Conference on Big Data*. Springer, 2018, pp. 108–123.
- [4] O. Penner, R. K. Pan, A. M. Petersen, K. Kaski, and S. Fortunato, "On the predictability of future impact in science," *Scientific reports*, vol. 3, p. 3052, 2013.
- [5] F. Radicchi and C. Castellano, "Understanding the scientific enterprise: citation analysis, data and modeling," in *Social Phenomena*. Springer, 2015, pp. 135–151.
- [6] D. Wang, C. Song, and A.-L. Barabási, "Quantifying long-term scientific impact," *Science*, vol. 342, no. 6154, pp. 127–132, 2013.
- [7] M. Perc, "The matthew effect in empirical data," *Journal of The Royal Society Interface*, vol. 11, no. 98, p. 20140378, 2014.
- [8] E. Sarigöl, R. Pfitzner, I. Scholtes, A. Garas, and F. Schweitzer, "Predicting scientific success based on coauthorship networks," *EPJ Data Science*, vol. 3, no. 1, p. 9, 2014.
- [9] T. Amjad, Y. Ding, J. Xu, C. Zhang, A. Daud, J. Tang, and M. Song, "Standing on the shoulders of giants," *Journal of Informetrics*, vol. 11, no. 1, pp. 307–323, 2017.
- [10] Z. Shao, J. Tang, Y. Zhang, B. Gao, and Y. Wang, "Scholar trajectory: Visualizing the migration of talents," 2018.

- [11] W. F. Laurance, D. C. Useche, S. G. Laurance, and C. J. Bradshaw, "Predicting publication success for biologists," *BioScience*, vol. 63, no. 10, pp. 817–823, 2013.
- [12] A. Letchford, H. S. Moat, and T. Preis, "The advantage of short paper titles," *Royal Society open science*, vol. 2, no. 8, p. 150266, 2015.
- [13] E. Wager, D. G. Altman, I. Simera, and T. P. Toma, "Do declarative titles affect readers' perceptions of research findings? a randomized trial," *Research integrity and peer review*, vol. 1, no. 1, p. 11, 2016.
- [14] J. Sienkiewicz and E. G. Altmann, "Impact of lexical and sentiment factors on the popularity of scientific papers," *Royal Society open science*, vol. 3, no. 6, p. 160140, 2016.
- [15] A. Letchford, T. Preis, and H. S. Moat, "The advantage of simple paper abstracts," *Journal of Informetrics*, vol. 10, no. 1, pp. 1–8, 2016.
- [16] S. Uddin and A. Khan, "The impact of author-selected keywords on citation counts," *Journal of Informetrics*, vol. 10, no. 4, pp. 1166–1177, 2016.
- [17] K. McKeown, H. Daume III, S. Chaturvedi, J. Paparrizos, K. Thadani, P. Barrio, O. Biran, S. Bothe, M. Collins, K. R. Fleischmann *et al.*, "Predicting the impact of scientific concepts using full-text features," *Journal of the Association for Information Science and Technology*, vol. 67, no. 11, pp. 2684–2696, 2016.

APPENDIX

TABLE VIII
ALL FIELDS PRESENT IN THE ORIGINAL DATASET.

Feature	Type
title	str
authors	list
n_citation	int
year	int
id	str
r	int
labels	list
volume	str
keywords	list
issue	str
abstract	str
versions	list
hash	str
lang	str
doi	str
venue	dict
page_start	str
isbn	str
pdf	str
issn	str
page_str	str
title_zh	str
pdf_src	list
hash_zh	str
page_end	str
keywords_zh	list
urls	list
abstract_zh	str
v_versions	list
retrieve_info	dict
wos_type	list
classes	list
n_wos_citation	int
i	str

TABLE IX
CLUSTERS AND A SELECTION OF THEIR MOST PREVALENT FEATURES.

Cluster	Features
0	alice, bob, girls, programming
1	storage, memory, data, server
2	translation, word, chinese, english
3	peer, p2p, peers, streaming
4	video, videos, 3d, quality
5	social, users, content, user
6	web, search, pages, page
7	data, algorithms, large, paper
8	learning, training, text, classification
9	algorithm, collision, present, compute
10	qos, service, multimedia, applications
11	traffic, network, flows, packet
12	motion, surface, 3d, algorithm
13	image, object, images, recognition
14	detection, malware, attacks, attack
15	query, queries, xml, search
16	graph, graphs, vertices, edges
17	design, people, user, information
18	retrieval, recommender, systems, information
19	log, time, algorithm, size
20	internet, tcp, network, protocol
21	problem, algorithm, algorithms, time
22	privacy, private, data, information
23	security, secure, protocols, protocol
24	devices, mobile, device, users
25	routing, network, networks, nodes
26	ontology, owl, semantic, logic


```

1#!/usr/bin/env python3
2import pickle
3
4ATTR_LIST = ['title', 'abstract', 'keywords', 'year', 'authors', 'page_start', 'page_end', '
    venue', 'n_citation']
5NAME_ATTR_LIST = ['name', 'aff']
6VENUE_ATTR_LIST = ['name']
7
8def process_paper(p):
9    # Returns None if the paper is missing ANY of the attributes
10    if not all(a in p.keys() for a in ATTR_LIST):
11        return None
12    if not all(p[k] for k in ATTR_LIST):
13        return None
14    # Affiliation attribute present and not empty
15    if not all(k in author.keys() for author in p['authors'] for k in NAME_ATTR_LIST):
16        return None
17    if not all(author[k] for author in p['authors'] for k in NAME_ATTR_LIST):
18        return None
19    # Venue Attribute present and not empty
20    if not all(a in p['venue'].keys() for a in VENUE_ATTR_LIST):
21        return None
22    if not all(p['venue'][k] for k in VENUE_ATTR_LIST):
23        return None
24
25    # Copies only the attributes in the list from the old paper to the new paper
26    new_p = {k:v for (k,v) in p.items() if k in ATTR_LIST and v}
27    if not all(a in new_p.keys() for a in ATTR_LIST):
28        return None
29    else:
30        return new_p
31
32def transform(p):
33    # Authors
34    authors = p['authors']
35    author_list = [a['name'].upper() for a in authors]
36    p['authors'] = author_list
37    # Affiliations
38    if not any(not 'aff' in a.keys() for a in authors):
39        aff_list = [a['aff'].upper() for a in authors]
40        p['aff'] = aff_list
41    # Page count
42    if p['page_start'].isnumeric() and p['page_end'].isnumeric():
43        page_cnt = int(p['page_end']) - int(p['page_start'])
44        del p['page_start']
45        del p['page_end']
46        p['page_count'] = page_cnt + 1
47    else:
48        p['page_count'] = None
49    p['venue'] = p['venue']['name'].upper()
50
51def hist_attr(papers, attr):
52    vals = {}
53
54    for p in papers.values():
55        if not p[attr] in vals.keys():
56            vals[p[attr]] = 1
57        else:
58            vals[p[attr]] += 1
59
60    for k,v in sorted(vals.items(), key=lambda kv: kv[1], reverse=True):
61        print("%s: %d" % (k,v))
62
63if __name__ == "__main__":
64    filename = 'data/trajectory_data.list'
65    with open(filename, 'rb') as data_source:
66        data = pickle.load(data_source)
67

```

```

68 papers = {}
69 # List of publications that meet the criteria
70 for author in data:
71     for p in author['pubs']:
72         pID = p['id'].upper()
73         # Create a deep copy with the specified attributes
74         new_p = process_paper(p)
75         if new_p and not pID in papers.keys():
76             transform(new_p)
77             papers[pID] = new_p
78
79 print("Processed: %d papers" % len(papers))
80 with open('data/database.pickle', 'wb') as f:
81     pickle.dump(papers, f)

```

Listing 1. Python script to extract the papers from the original dataset and store papers containing all the required features in a new database.

```

1  #!/usr/bin/env python3
2  import pickle
3  from sklearn.feature_extraction.text import TfidfVectorizer
4  from sklearn.cluster import KMeans
5  import numpy as np
6
7  if __name__ == "__main__":
8      filename = 'data/database.pickle'
9      with open(filename, 'rb') as data_source:
10         papers = pickle.load(data_source)
11
12     dataset = [p['abstract'] for p in papers.values()]
13     titles = [p['title'] for p in papers.values()]
14
15     # Clustering!
16     vectorizer = TfidfVectorizer(max_df=0.5, max_features=10000,
17                                min_df=2, stop_words='english',
18                                use_idf=True)
19     X = vectorizer.fit_transform(dataset)
20
21     print("n_samples: %d, n_features: %d" % X.shape)
22     print()
23
24     # Determined experimentally
25     exp_k = 27
26
27     km = KMeans(n_clusters=exp_k, init='k-means++', max_iter=100, n_init=1,
28                verbose=True, random_state=5)
29
30     print("Clustering sparse data with %s" % km)
31
32     km.fit(X)
33     print(km.random_state)
34     groups = km.predict(X)
35
36     order_centroids = km.cluster_centers_.argsort()[:, ::-1]
37     terms = vectorizer.get_feature_names()
38
39     with open('data/pp-fit.pickle', 'wb') as f:
40         pickle.dump(groups, f)
41
42     for g in range(min(groups), max(groups)+1):
43         cluster = {}
44         citations = []
45         indices = np.argwhere(groups==g)
46
47         # Print the relevant keywords and titles in the cluster
48         print("Class: %d | Len: %d" % (g, indices.shape[0]))
49         print("Terms:", end='')
50         for t_idx in order_centroids[g, :10]:
51             print(' %s' % terms[t_idx], end='')
52         print()

```

```

53     for idx in np.nditer(indices):
54         print('\t' + titles[idx])
55
56     # Add all the relevant papers to the cluster
57     for idx in np.nditer(indices):
58         k = list(papers.keys())[idx]
59         if papers[k]['n_citation'] < 1000:
60             cluster[k] = papers[k]
61             citations.append(papers[k]['n_citation'])
62
63     print("Median: " + str(np.median(citations)))
64
65     with open('data/cluster-%d.pickle'%g, 'wb') as f:
66         pickle.dump(cluster, f)

```

Listing 2. Python script to cluster the papers.

```

1  #!/usr/bin/env python3
2  import pickle
3  from sklearn.cluster import KMeans
4  import numpy as np
5  import sys
6
7  if __name__ == "__main__":
8      if not len(sys.argv) == 2:
9          print("Usage: pp-transform <cluster-max>")
10         sys.exit()
11
12     c_max = int(sys.argv[1])
13
14     # Iterate over all clusters to process
15     for c in range(0, c_max+1):
16         print("Cluster: %d" % c)
17         # Load the papers in that cluster
18         filename = 'data/cluster-%d.pickle' % c
19         with open(filename, 'rb') as f:
20             papers = pickle.load(f)
21
22         # Calculate the median citations in each cluster
23         cite_list = [p['n_citation'] for p in papers.values()]
24         median = np.median(cite_list)
25
26         # Features: title_count, abstract_count, keyword_count, year, author_count,
27         # page_count
28         X = np.empty([0, 6], dtype=int)
29         # Features: Title, Abstract, Keywords
30         X_text = []
31         # Target
32         Y = np.empty(0, dtype=int)
33         for idx, p in enumerate(papers.values()):
34             if not p['page_count'] == None:
35                 # Extracting the numerica data
36                 data = []
37                 data.append(len(p['title'].split(' ')))
38                 data.append(len(p['abstract'].split(' ')))
39                 data.append(len(p['keywords']))
40                 data.append(p['year'])
41                 data.append(len(p['authors']))
42                 data.append(p['page_count'])
43                 X = np.append(X, [data], axis=0)
44                 # Extracting the text data
45                 X_text.append([p['title'], p['abstract'], ' '.join(p['keywords'])])
46                 # Generating the ground truth
47                 if p['n_citation'] < median:
48                     Y = np.append(Y, 0)
49                 else:
50                     Y = np.append(Y, 1)
51
52         with open('data/X-%d.pickle'%c, 'wb') as f:

```

```

52         pickle.dump(X, f)
53     with open('data/X-text-%d.pickle'%c, 'wb') as f:
54         pickle.dump(X_text, f)
55     with open('data/Y-%d.pickle'%c, 'wb') as f:
56         pickle.dump(Y, f)

```

Listing 3. Python script to create a training set from the clusters.

```

1  #!/usr/bin/env python3
2  import pickle
3  import sys
4  from sklearn.model_selection import train_test_split
5  from sklearn.model_selection import GridSearchCV
6  from sklearn.preprocessing import StandardScaler
7  from sklearn.feature_extraction.text import CountVectorizer
8  from sklearn.feature_extraction.text import TfidfTransformer
9  # Import Classifiers
10 from sklearn.ensemble import RandomForestClassifier
11 from sklearn.naive_bayes import MultinomialNB
12 from sklearn.svm import SVC
13
14 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
15 import numpy as np
16
17 RAND_STATE = 5
18
19 if __name__ == "__main__":
20     # Input clusters to process
21     if not len(sys.argv) == 2:
22         print("Usage: pp-transform <cluster-file>")
23         sys.exit()
24
25     class_txt = ''
26     with open(sys.argv[1], 'r') as f:
27         class_txt = f.read()
28
29     total_tn = 0
30     total_fp = 0
31     total_fn = 0
32     total_tp = 0
33     # Iterate over clusters
34     for c_str in class_txt.strip().split('\n'):
35         # Load the cluster data
36         c = int(c_str)
37         print("Class %d" % c)
38         with open('data/X-%d.pickle'%c, 'rb') as f:
39             X = pickle.load(f)
40         with open('data/X-text-%d.pickle'%c, 'rb') as f:
41             X_text = pickle.load(f)
42         with open('data/Y-%d.pickle'%c, 'rb') as f:
43             Y = pickle.load(f)
44
45         # Split numeric data
46         X0_train, X0_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.5, \
47             random_state=RAND_STATE)
48
49         # Standard scale
50         scaler = StandardScaler()
51         scaler.fit(X0_train) #only scale w/ respect to training data
52         X0_train = scaler.transform(X0_train)
53         X0_test = scaler.transform(X0_test)
54         X = scaler.transform(X) #for validation on whole group
55
56         # Title TF-IDF
57         count_vec = CountVectorizer()
58         X1_count = count_vec.fit_transform([x[0] for x in X_text])
59         tfidf_transformer = TfidfTransformer()
60         X1_tfidf = tfidf_transformer.fit_transform(X1_count)
61         X1_text_train, X1_text_test, Y_train, Y_test = train_test_split(X1_tfidf, \

```

```

62         Y, test_size=0.5, random_state=RAND_STATE)
63
64     # Abstract TF-IDF
65     count_vec = CountVectorizer()
66     X2_count = count_vec.fit_transform([x[1] for x in X_text])
67     tfidf_transformer = TfidfTransformer()
68     X2_tfidf = tfidf_transformer.fit_transform(X2_count)
69     X2_text_train, X2_text_test, Y_train, Y_test = train_test_split(X2_tfidf, \
70         Y, test_size=0.5, random_state=RAND_STATE)
71
72     # Keyword TF-IDF
73     count_vec = CountVectorizer()
74     X3_count = count_vec.fit_transform([x[2] for x in X_text])
75     tfidf_transformer = TfidfTransformer()
76     X3_tfidf = tfidf_transformer.fit_transform(X3_count)
77     X3_text_train, X3_text_test, Y_train, Y_test = train_test_split(X3_tfidf, \
78         Y, test_size=0.5, random_state=RAND_STATE)
79
80     # Random Forest Search Parameters
81     param_grid = [
82         {'n_estimators': [1,2,4,6,8,10,12], 'max_depth': [None], '
min_samples_split': [2], 'max_features': [2,4,6], 'random_state': [0]}
83     ]
84     param_grid1 = [
85         {'n_estimators': [1,2,4,6,8,10,12], 'max_depth': [None], '
min_samples_split': [2], 'max_features': [1,2,3,4], 'random_state': [0]}
86     ]
87
88     # Optimize numeric Random Forest classifier
89     clf0 = GridSearchCV(RandomForestClassifier(), param_grid, cv=5, scoring='
accuracy')
90     clf0.fit(X0_train, Y_train)
91
92     print("Detailed classification report:")
93     print(clf0.best_params_)
94     print()
95     print("The model is trained on the full development set.")
96     print("The scores are computed on the full evaluation set.")
97     print()
98
99     Y0_train_pred = clf0.predict(X0_train)
100    Y_true, Y0_pred = Y_test, clf0.predict(X0_test)
101
102    print(classification_report(Y_true, Y0_pred))
103    print("ACCURACY: ", end='')
104    print(accuracy_score(Y_true, Y0_pred))
105
106    importances = clf0.best_estimator_.feature_importances_
107    indices = np.argsort(importances)[::-1]
108    for f in range(X.shape[1]):
109        print("%d. feature %d (%f)" % (f + 1, indices[f], importances[indices[f]]))
110
111    # Title MNB Classifier
112    clf1 = MultinomialNB().fit(X1_text_train, Y_train)
113    Y1_train_pred = clf1.predict(X1_text_train)
114    Y_true, Y1_pred = Y_test, clf1.predict(X1_text_test)
115    print(classification_report(Y_true, Y1_pred))
116    print("ACCURACY: ", end='')
117    print(accuracy_score(Y_true, Y1_pred))
118
119    # Abstract MNB Classifier
120    clf2 = MultinomialNB().fit(X2_text_train, Y_train)
121    Y2_train_pred = clf2.predict(X2_text_train)
122    Y_true, Y2_pred = Y_test, clf2.predict(X2_text_test)
123    print(classification_report(Y_true, Y2_pred))
124    print("ACCURACY: ", end='')
125    print(accuracy_score(Y_true, Y2_pred))
126

```

```

127     # Keywords MNB Classifier
128     clf3 = MultinomialNB().fit(X3_text_train, Y_train)
129     Y3_train_pred = clf3.predict(X3_text_train)
130     Y_true, Y3_pred = Y_test, clf3.predict(X3_text_test)
131     print(classification_report(Y_true, Y3_pred))
132     print("ACCURACY: ", end='')
133     print(accuracy_score(Y_true, Y3_pred))
134
135     # Ensemble Classifier
136     clf4 = GridSearchCV(RandomForestClassifier(), param_grid1, cv=5, scoring='
accuracy')
137     Y4_train = np.vstack((Y0_train_pred, Y1_train_pred, Y2_train_pred, Y3_train_pred)
).T
138     Y4_pred = np.vstack((Y0_pred, Y1_pred, Y2_pred, Y3_pred)).T
139     clf4.fit(Y4_train, Y_train)
140
141     print("Detailed classification report:")
142     print(clf4.best_params_)
143     print()
144     print("The model is trained on the full development set.")
145     print("The scores are computed on the full evaluation set.")
146     print()
147
148     Y_true, Y_pred = Y_test, clf4.predict(Y4_pred)
149
150     print(classification_report(Y_true, Y_pred))
151     print("ACCURACY: ", end='')
152     print(accuracy_score(Y_true, Y_pred))
153
154     importances = clf4.best_estimator_.feature_importances_
155     indices = np.argsort(importances)[::-1]
156     for f in range(Y4_pred.shape[1]):
157         print("%d. feature %d (%f)" % (f + 1, indices[f], importances[indices[f]]))
158
159     # Accuracy calculation
160     tn, fp, fn, tp = confusion_matrix(Y_true, Y_pred).ravel()
161     total_tn += tn
162     total_fp += fp
163     total_fn += fn
164     total_tp += tp
165
166     print("OVERALL ACCURACY: ", end='')
167     print((total_tp+total_tn)/(total_tp+total_fp+total_fn+total_tn))

```

Listing 4. Python script to classify the successful papers.