# Bibliometrics:
# Predicting Publication Success

Team 4: Tamanna Baig and Jon Oakley

April 18th 2019

Clemson University

# Introduction

# Motivation

" *Knowledge is power* "
– Sir Francis Bacon

## Motivation

– 2.5 million scientific papers published each year[1]

– Research trends

– Research *funding*

---

[1]The STM Report; Fourth Edition, March 2015

# Background: Bibliometrics

# Background: Bibliometrics

## Defining the Problem

- Given: a corpus of academic publications
- Goal: predict which papers will be "successful"
- Success Metric: citation count
- Success Definition: citation count > median citation count for that cluster
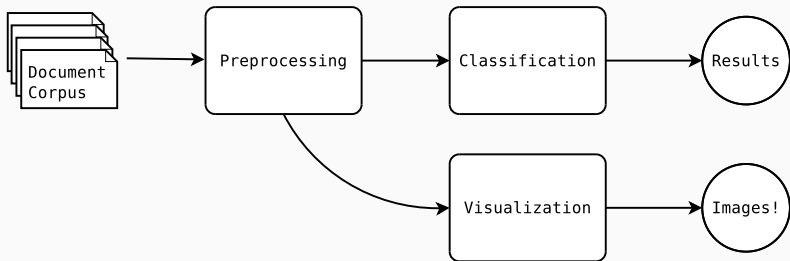- Evaluation Metric: accuracy

# Overview

# Preprocessing

# Workflow

# Extraction, Thinning, and Meta Features

Introduction

Preprocessing

Visualization

Classification

Results

Conclusion

– Initial Corpus: 22,588
– Removed:
  ☐ Papers after 2013
  ☐ Outliers
  ☐ Papers missing features
– Final Corpus: 4,914

# Clustering

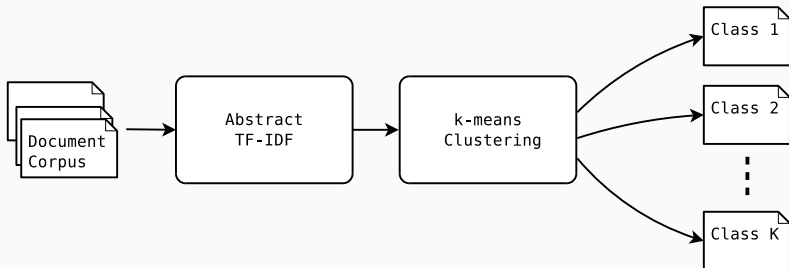## Clustering: Choosing K

# Clustering: Good Clusters

| Cluster | Count | Keywords |
|---------|-------|----------|
| ~~0~~ | ~~15~~ | ~~alice, bob, girls, programming, communication~~ |
| 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 |

# Visualization

# Yearly Publications

Number of papers published each year

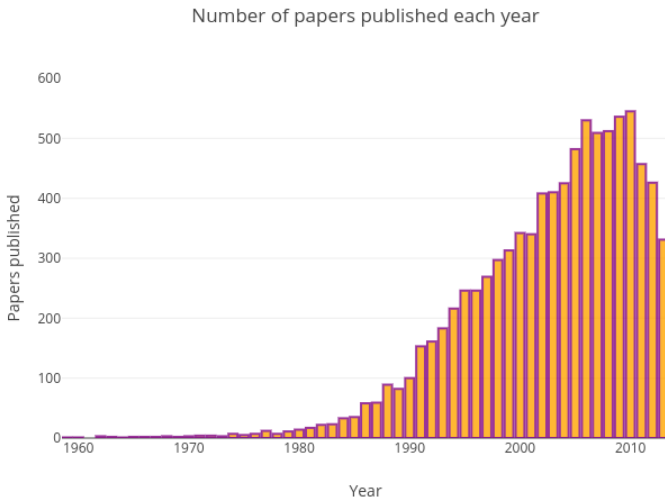# Author Publications

Top 10 Most Published Authors

# Popular Publications

Most Cited Papers over the Years

# Institution Publications

Most Affiliations to Institutions

# Classification

# Workflow

## Random Forest

# Multinomial Naive Bayes

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} \, P(d|c)P(c)$$

$$P(d|c) = P(f_1, f_2, ..., f_n|c) = P(f_1|c)P(f_2|c)...P(f_n|c)$$

$$P(c) = \frac{N_c}{N_d}$$

– $c$: Class

– $d$: Document

– $f_i$: Feature

– $N_c$: Number of words in class $c$

– $N_d$: Number of words in document $d$

# Classifier Accuracy

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

## Training Method

- test_train_split: 50%
- GridSearchCV
- CV = 5
- RF Params: n_estimators and max_features

20

# Results

# Workflow

# Classification Parameters

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

# Classification Parameters

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

## Feature Importance

| 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 |

## Feature Importance

| 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 |

## Classification Accuracy

| Cluster | Accuracy (%) | | | | |
|---|---|---|---|---|---|
| | 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 |

– Overall Accuracy (All Clusters): 57.7%

– Overall Accuracy (4,15,17,20,22): 66.5%

# Conclusion

# Conclusions

– Cleaning data is important

– Some clusters are better than others

– **66.5% prediction accuracy in optimal clusters**

– Future work: targeted dataset

– Controlling for year

Introduction

Preprocessing

Visualization

Classification

Results

Conclusion