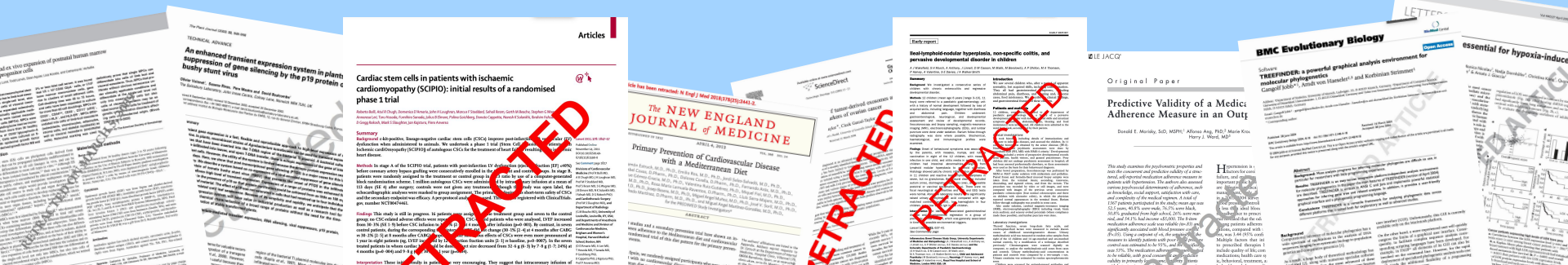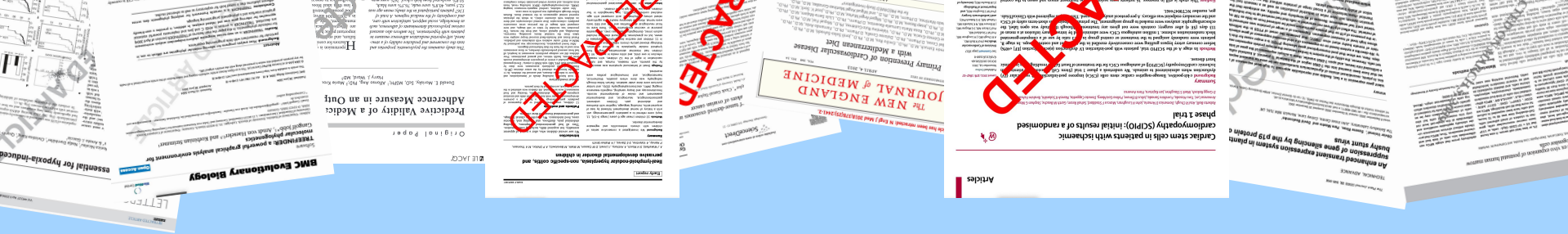# *Predicting Paper Retractions*

William Davis and Jack Kendrick

# *Motivation*

→ In theory, articles published in peer reviewed journals should hold up to scrutiny. When reading a paper, it is often taken for granted that the research is credible and trustworthy.

→ However, sometimes bad papers fall through the cracks and end up being retracted.

→ Retractions indicate seriously flawed and unreliable research, errors, fraud, ethical issues, or other serious concerns

# Can we identify papers at a high risk of retraction?

# *Starting Small: PLOS One*

→ We concentrated on the journal PLOS One
→ Relatively well known and respected, but has a high number of retractions

| Rank | Journal | Retractions* |
|:---:|:---:|:---:|
| 1 | 2011 International Conference on E-Business and E-Government | 1280 |
| 2 | 2011 5th International Conference on Bioinformatics and Biomedical Engineering | 1084 |
| 3 | PLoS One | 944 |
| 4 | Journal of Physics: Conference Series | 878 |

# The Dataset

➔ We collected data from the OpenAlex database using the PyAlex API

➔ **Huge** amounts of data available (including retractions)

➔ We used the raw data from OpenAlex to build features of interest

◆ Fraction of authors that have previous retractions, whether any authors come from institutions with many retractions, etc.

➔ Our dataset consists of papers published in 2010-2020: 424,223 papers, with 797 retractions

## Challenge: Retracted vs Non-Retracted Classes are massively unbalanced

# Our Approach

→ Baseline logistic regression model
→ Compare more baseline to nearest neighbour, random forest, and SVC classification methods.

## Key Performance Indicators: $F_1$-score, Precision, and Recall

# Baseline Model

➔ Used **forward stepwise subset selection** to choose features.

➔ The most informative feature is the proportion of authors on a paper that have been previously retracted

➔ Second most informative feature is a measure of how many retractions any institution associated to the paper has received.

➔ Results of subset selection correspond to our intuition on features that correlate to a risk of retraction
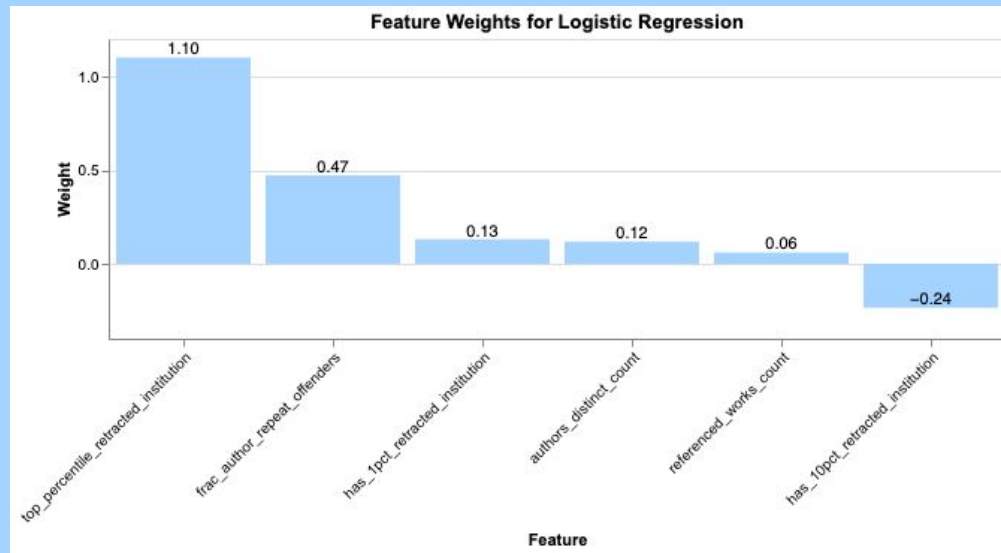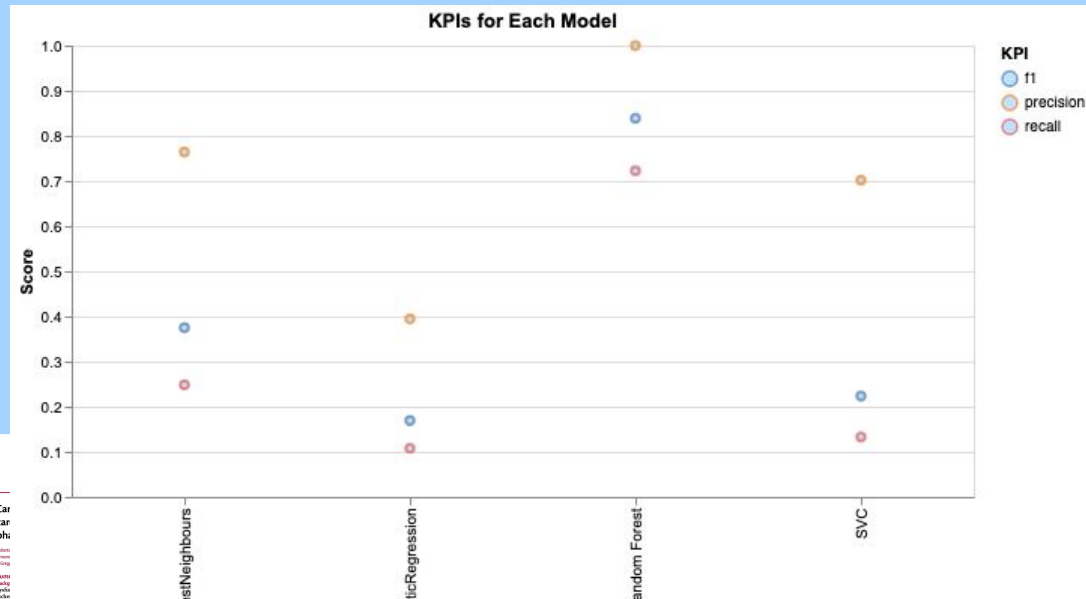
# Baseline Model

F₁-Score: 0.169

Precision: 0.394

Recall: 0.107



Feature Weights for Logistic Regression

# *Model Selection*

➔ Tested k-Nearest Neighbors (k = 1, 2,...., 10), random forests, and support vector classification.

➔ We used stratified 10-fold cross-validation to choose hyperparameters.

# Results

➔ Random Forest classification far outperformed our other models in training.

*Final Model: Random Forest with 500 estimators, max depth of 20*

*$F_1$:*            *0.289*

*Precision:*           *0.691*

*Recall:*            *0.182*

*Accuracy:*           *0.998*

# *Future Directions*

➔ Expand dataset to include other journals and features relating to funding sources, journal publishers, etc.

➔ Investigate papers published in the post ChatGPT-era. Can we use 'Contains AI generated text' as a feature?

# *Acknowledgements*

➔ Many thanks to our mentor Greg Edwards for his help and expertise throughout this project

➔ Thank you to the Erdos Institute for providing us with this excellent bootcamp.