

# Detecting Fake Users on Social Media with Neo4j and Random Forest Classifier

March 4, 2020

Presented by **Yichun Zhao**,  
Department of Human & Social Development

Supervised by **Dr. Jens Weber**,  
Department of Computer Science

## BACKGROUND

Fake users on social media are perceived as popular [1], and they spread false information or fake news by making it look real [2], manipulating real users in making a decision.

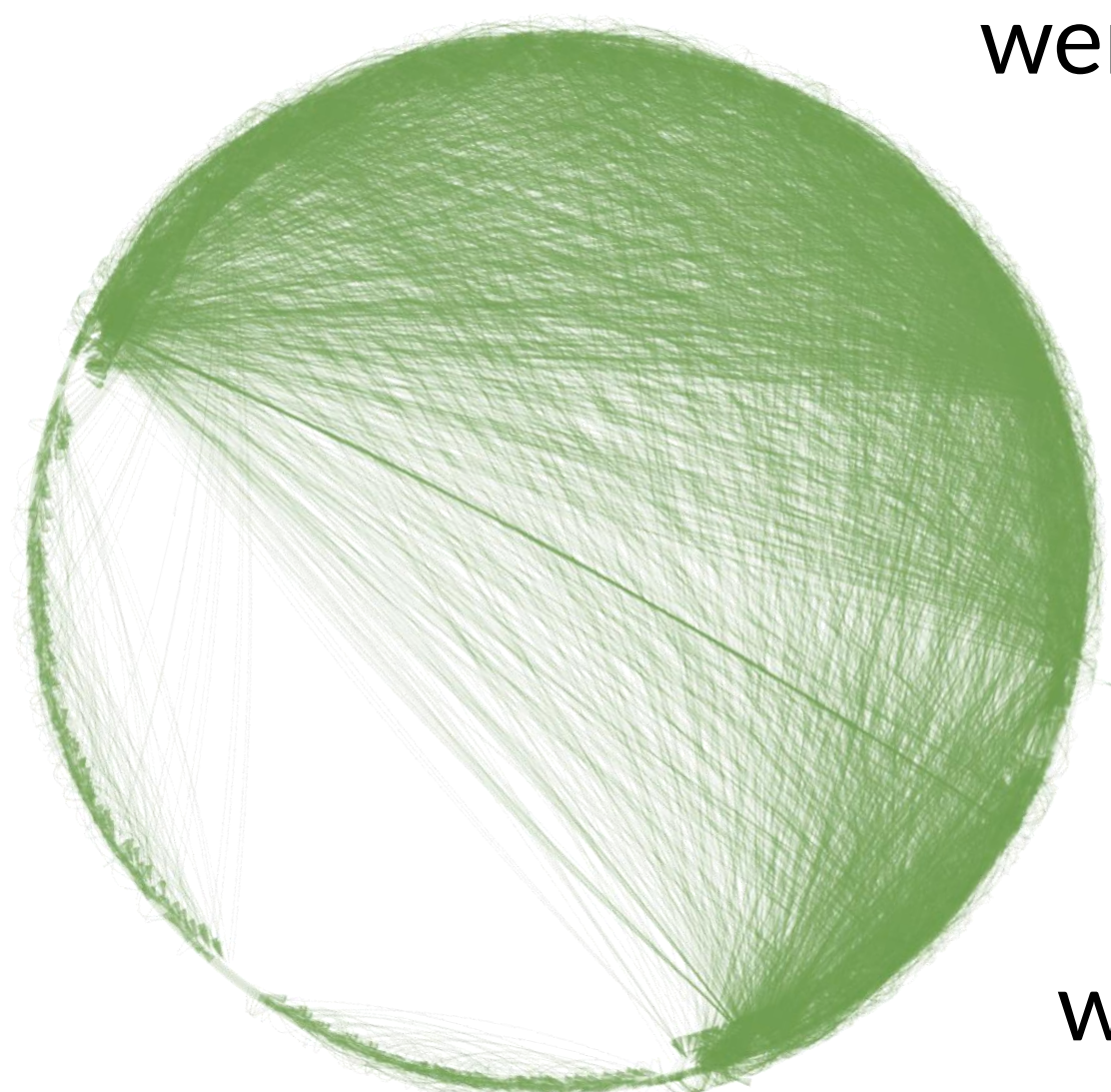
## OBJECTIVE

To improve the accuracy of detecting fake users in the previous study by A. Mehrotra, M. Sarreddy and S. Singh [1], by using different centrality measures supported by the Neo4j graph database.

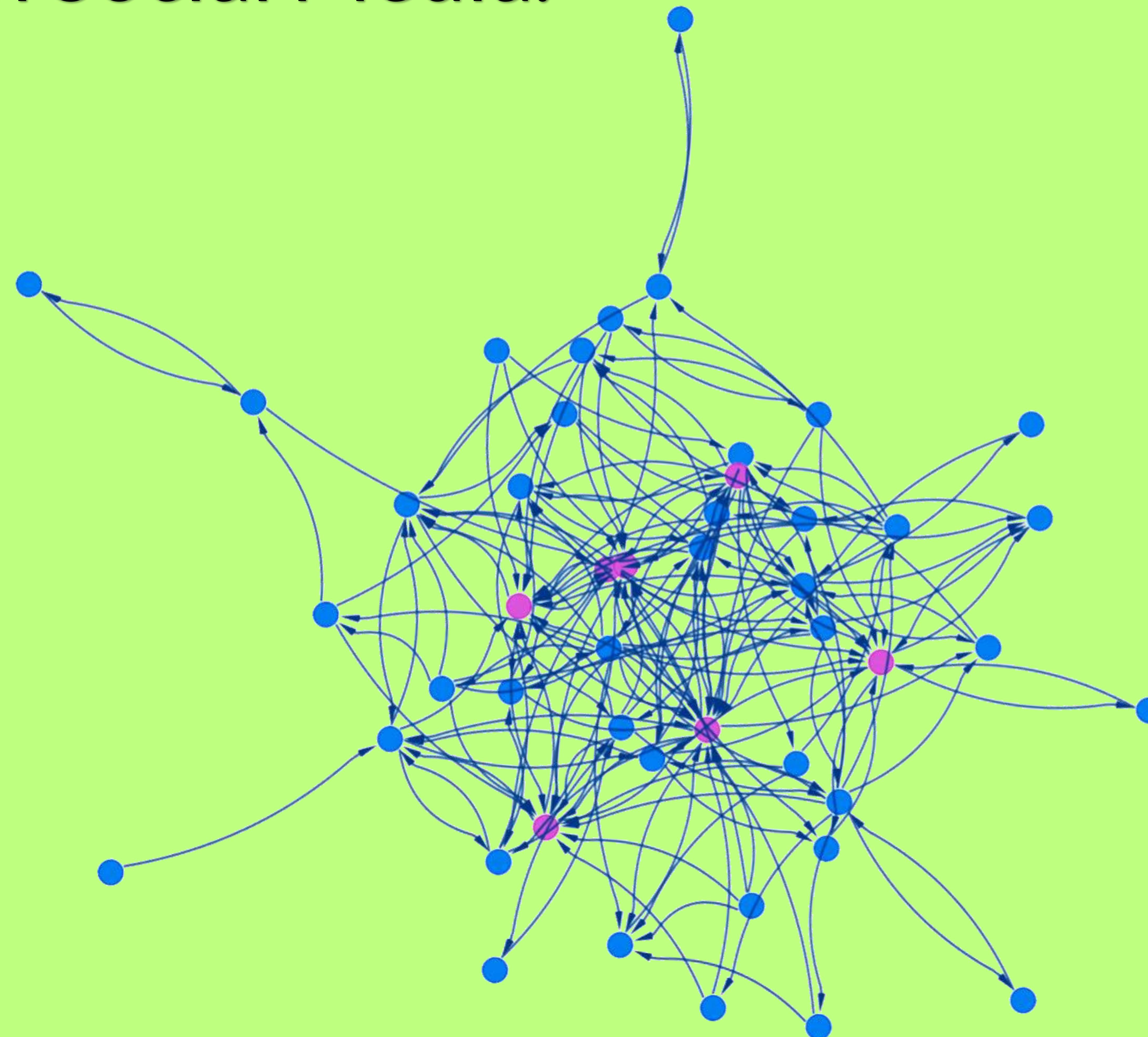
## METHODS

1. Used Neo4j graph database to import the same dataset (Fake Project) used in the previous study [1], and filtered the graph to remove isolated nodes which do not have relationships.
2. Calculated the centrality measures using Neo4j: Betweenness, Eigenvector, In-degree, Out-degree, Closeness, and Page-rank.
3. Used Random Forest Classifier with 10-fold cross-validation for training and testing, and centrality measures as features.
4. Evaluated the features to identify the most influential one.
5. Two more datasets (Class and Twitter)

were used for evaluating the methods. The procedures were applied to them.  
6. Due to the large size of the Twitter dataset, it was randomly sampled.



# Machine Learning Classifier using Centrality Measures as Features Detects Fake Users on Social Media.



## RESULTS

Dataset	Precision	Recall	Accuracy
Fake Project	99.5%	99.5%	99.5%
Class	90.9%	91.5%	91.5%
Twitter (sampled)	87.6%	87.7%	87.7%

Most influential feature: Closeness measure

## DISCUSSION

Comparing with the results from the previous study by A. Mehrotra, M. Sarreddy and S. Singh [1] using the Fake Project dataset (Precision = 89.0%; Recall = 100%; Accuracy = 95%), a significant increase in precision and accuracy is observed. This might be due to the inclusion of the new feature, Closeness measure, which has the highest correlation with the output. The slight decrease in recall might be due to the exclusion of Katz and Load measures, which were used in the previous study [1].

The lower results achieved from using the Class dataset might be due to its small size (47 nodes and 177 relationships). The randomly sampled Twitter dataset might not represent the same state of the original graph, and not preserve certain graph properties. In addition, both datasets refer to their “fake” nodes as “anomalous”, meaning these user might not be fake. This might also explain the lower results.

## CONCLUSION

The use of different centrality measures as features for the Random Forest classifier detects fake users with reasonable results.

Centrality measures are node properties. We can include edge properties such as link prediction measures as new features to possibly achieve even better results.

## REFERENCES

[1] A. Mehrotra, M. Sarreddy and S. Singh, "Detection of fake Twitter followers using graph centrality measures", 2016 2nd International Conference on Contemporary Computing and Informatics (IC3I), 2016. Available: 10.1109/ic3i.2016.7918016 [Accessed 12 February 2020].

[2] "08. Fake News, Accounts, and Bots — The Democratic Engagement Exchange", The Democratic Engagement Exchange, 2020. [Online]. Available: <https://www.engagedemocracy.ca/fake-news>. [Accessed: 12- Feb- 2020].

This research was supported by the Jamie Cassels Undergraduate Research Awards, University of Victoria.

