

## INTISARI

Teknologi serta internet banyak membawa perkembangan ke penggunaannya, baik positif maupun negatif. Hampir semua aktivitas dapat dilakukan di internet. Namun, pengguna dapat menggunakan data diri palsu di dunia maya, sehingga tindakan mereka akan sangat sulit dilacak. Semua orang memiliki potensi melakukan pemalsuan data, bahkan di Indonesia sekalipun. Hal ini menjadi salah satu yang mendorong peningkatan angka kejahatan online, tak terkecuali *cyberbullying*.

Penulis melakukan penelitian terhadap *cyberbullying* pada data *tweet* dalam Bahasa Indonesia. Fokus penelitian ini adalah *data-centric*, yang berusaha meningkatkan kualitas data sebelum dilakukan klasifikasi, dalam rangka meningkatkan performa model pendeteksi perilaku *cyberbullying* pada komunitas online, khususnya dalam Bahasa Indonesia. Aplikasi *data-centric* yang dilakukan pada penelitian ini adalah dengan memanfaatkan *data augmentation* untuk menanggulangi masalah keterbatasan data. Tujuan dilakukan *data augmentation* adalah untuk meningkatkan jumlah *data point* pada *dataset* sebelum digunakan untuk pelatihan mesin. Metode *data augmentation* yang digunakan adalah *character augmentation*, *word augmentation* dan *backtranslation*. Hasil dari masing-masing *data augmentation* kemudian diklasifikasi menggunakan *deep learning* LSTM dan *machine learning* SVM.

Secara keseluruhan, hasil dari penelitian ini menunjukkan *data augmentation* memberikan peningkatan performa baik untuk LSTM maupun SVM. Akurasi tertinggi pada pemodelan LSTM menghasilkan nilai akurasi 89,41% dengan peningkatan 6,66% dan skor F1 88,98% dengan peningkatan 7,01%. Sedangkan untuk pemodelan SVM, hasil nilai akurasi 92,16% dengan peningkatan 2,75% dan skor F1 92,42% dengan peningkatan 3,44%.

**Kata kunci** : *data augmentation, data-centric, prediction, cyberbullying detection, Twitter, long short-term memory, support vector machine, machine learning*

## ***ABSTRACT***

Technology and the internet have brought many developments to its users, both positive and negative. Almost all activities can be done on the internet. However, users can use fake personal data on the internet, thus their actions will be very difficult to track. Everyone has the potential to falsify data, even in Indonesia. This matter of fact has become one of the reasons for the increase in the number of online crimes, including cyberbullying.

The author conducted this research on cyberbullying tweets in Indonesian. The focus of this research is data-centric, which seeks to improve the quality of the data prior to classification, in order to improve the performance of the cyberbullying detection model in online communities, especially in Indonesian. The data-centric application carried out in this study is data augmentation to overcome the problem of data limitations. The purpose of data augmentation is to increase the number of data points in the dataset before being used for machine training. Data augmentation methods used are character augmentation, word augmentation and backtranslation. The results of each augmentation data are then classified using deep learning LSTM and machine learning SVM.

Overall, the results of this study show that data augmentation provides improvement for both LSTM and SVM performances. The result in LSTM modeling produces an accuracy value of 89,41% with an increase of 6.66% and a F1-score of 88.98% with an increase of 7.01%. As for the SVM modeling, the accuracy value is 92.16% with an increase of 2.75% and a F1-score of 92.42% with an increase of 3.44%.

***Keywords*** : *data augmentation, data-centric, prediction, cyberbullying detection, Twitter, long short-term memory, support vector machine, machine learning*