



2nd International Conference on Computer Science and Computational Intelligence 2017, ICCSCI
2017, 13-14 October 2017, Bali, Indonesia

Twitter Pornography Multilingual Content Identification Based on Machine Learning

Edo Barfian^{a*}, Bambang Heru Iswanto^b, Sani Muhamad Isa^a

^aBina Nusantara University, Jl. K. H. Syahdan No. 9 Kemanggisan, Jakarta 11480, Indonesia

^bJakarta State University, Jl. Rawamangun Muka, RT.11/RW.14, Rawamangun, Jakarta Timur, 13220, Indonesia

Abstract

Pornography on social media raises a lot of negative impact and affect the moral of children and teenagers. Social media used to spread pornography can have a negative impact. Thus, the spread of pornography on social media must be prevented. One of the social media which is often used as a medium pornography is Twitter. Pornography used on Twitter in the form of text and image. Among the two types of media, the text is very interesting to study because of the use of a variety of languages. In this study, the classification process will be conducted in Indonesian and English tweet and a combination of both languages. This classification uses three methods of machine learning, Decision Tree, Naive Bayes and Support Vector Machines for the purpose of comparing which method is the best in the classification process. In this study also conducted additional experiment was carried out with the aim of improving the performance in classification. The results showed that the level of accuracy is quite high. However, different grammar is a constraint that affects the accuracy of the results in the classification.

© 2017 The Authors. Published by Elsevier B.V.

Peer-review under responsibility of the scientific committee of the 2nd International Conference on Computer Science and Computational Intelligence 2017.

Keywords: Pornography; Twitter; Machine Learning

* Corresponding author.

E-mail address: bedhaw@gmail.com

1. Introduction

Negative content (pornography) that appear on the Internet have a negative impact and damage the morale of children and adolescents, according to the survey results of the Indonesia National Commission for Child Protection (KPAI) in 2012 to 4500 levels of high school students in 12 major cities in Indonesia. The survey found that 97.2% of respondents had access to a website that has a negative content. KPAI also added that 91% of respondents had been doing "kissing" and "petting". While 62.1% of female students have sex without marriage and 22% of female students had an abortion. According to the survey, sexual development in children and adolescents due to easy access to negative content on the internet ¹.

Prevention of pornography, especially on porn site or social media containing negative content with a technological approach is regarded as an effective way because it can prevent the spread of pornography directly ². Polpinij et al. also stated that the prevention of the spread of negative sites can be done through such negative content filtering URL Blocking, Keyword Filtering, and Rating System ³. URL blocking is done by entering the URL address to filter system. Keyword Filtering is done by determining whether a negative account or a website based on the text or keywords.

In this study, a Twitter account will be classified according to pornographic content multilingual. Then do the classification into two types of output: positive (non-pornographic) and negative (pornography) using TF-IDF algorithm to perform feature extraction and classified using the Decision Tree, Naive Bayesian Classifier, and Support Vector Machine.

2. Related Work

This previous research discusses about the sentiment analysis as a positive opinion and negative opinion on English data and Indonesian data by using Naïve Bayes Classifier (NBC) method and Support Vector Machine (SVM). Both the Naïve Bayes Classifier method and the Support Vector Machine method performed a good result. The results of the experiments show that the Support Vector Machine method performed better results than the Naïve Bayes Classifier method to classify English opinion and positive opinion in Indonesian language. While Naïve Bayes Classifier performed better results in classifying the negative opinion test data in Indonesian language. The result of accuracy is better on the classification of English opinion than the Indonesian language opinion due to the nature of the vocabulary ⁴.

The research goal is to create an application that is used to censor pornographic content based image and text using machine learning method. The author forms a corpus that is used to match pornographic words with data. The accuracy of the classification is 79.2% ³.

In this research, the authors performed a test to develop a multilingual-based sentiment analysis system. The authors evaluate Spanish and English-speaking methods and data. Then the author also performs additional testing in order to improve the performance of both language systems. The method used by the author is SMO SVM with unigram and bigram features. The resulting accuracy is 68.23% ⁵.

3. Method

In general, the system will be made in this study is a system that can analyze sentiment regarding pornographic content at Twitter. The system will provide sentiment analysis results opinions at pornographic content that is classified as pornographic or non-pornographic sentiment. The method for classification are Decision Trees, Naive Bayes and Support Vector Machines.

Decision tree is a flow chart like a tree structure, where each internal node shows a test in an attribute, each branch represents the result of the test, and the leaf node represents the classes or distribution classes ⁶.

Bayes's theorem was first proposed by a British Presbyterian priest in 1763 named Thomas Bayes who was used to calculate the probability of occurring on an event based on the observed falling of observations. Naïve Bayes Classifier is a simple opportunity classification based on Bayes theorem application with assumptions among independent variables. Then only the variance of a variable in a class is needed to determine the classification, not the whole of the covariance matrix. This research using probabilistic Naive Bayes.

Download English Version:

<https://daneshyari.com/en/article/4960435>

Download Persian Version:

<https://daneshyari.com/article/4960435>

[Daneshyari.com](https://daneshyari.com)