# Unsupervised Online Anomaly Detection of Internet Trolls in Streaming Data

**Albert Haque** [*]
Department of Computer Science
University of Texas at Austin
`akhaque@cs.utexas.edu`

## Abstract

With the increasing importance of online communities, discussion forums, and customer reviews, Internet "trolls" have proliferated thereby making it difficult for information seekers to find relevant and correct information. In this paper, we consider the problem of detecting and identifying Internet trolls, almost all of which are human agents. Identifying a human agent among a human population presents significant challenges compared to detecting automated spam or computerized robots. To learn a troll's behavior, we use contextual anomaly detection to profile each chat user. Using density-based clustering methods, we use contextual data such as the group's current goal, the current time, and the user to classify each point as an anomaly. A user with consistent anomalies will be labeled as a troll. We have successfully trained all three algorithms on a dataset consisting of 50 million data points from the viral Internet social experiment, Twitch Plays Pokemon. Using MapReduce techniques for preprocessing and feature extraction, we are able to achieve a prediction accuracy of Y% on streaming data as humans communicate with one another live.

## 1 Introduction

In Internet slang, a "troll" is a person who sows discord on the Internet by starting arguments or upsetting people, by posting inflammatory, extraneous, or off-topic messages in an online community (such as a forum, chat room, or blog), either accidentally or with the deliberate intent of provoking readers into an emotional response or of otherwise disrupting normal on-topic discussion [3, 4].

Application of the term troll is subjective. Some readers may characterize a post as trolling, while others may regard the same post as a legitimate contribution to the discussion, even if controversial [1]. Regardless of the circumstances, controversial posts may attract a particularly strong response from those unfamiliar with the dialogue found in some online, rather than physical, communities [1]. Experienced participants in online forums know that the most effective way to discourage a troll is usually to ignore it, because responding tends to encourage trolls to continue disruptive posts hence the often-seen warning: "Please don't feed the trolls" [5].

The contributions of this project are as follows. First, we propose a set of features used for identifying trolls in the viral, croudsourced Internet game, Twitch Plays Pokemon [2]. We use context-based techniques to understand the scenario a human is faced when entering input into the chatroom. We then compare the effects of different distance measures to understand their strengths and weaknesses.

The second major contribution of this paper is an online classification algorithm. It is initially trained using unsupervised methods on offline data. When switched to online mode, this algorithm updates as new data points are received from a live stream. We apply this algorithm to the Twitch Plays

---

[*]This paper was created on April 19, 2014.

Pokemon data stream and visualize the results. The primary purpose of this project, through these two contributions, is to in realtime, distinguish between trolls and humans on the Internet. Future work can be extended to email classifiers, comment moderation, and anomaly detection for online forums, reviews, and other communities.

The remainder of this paper is organized as follows: In Section 2 we conduct a review of the current literature, summarize related efforts, and give an overview of Twitch Plays Pokemon. Section 3 outlines how we extract features from our dataset. These features are then used by algorithms outlined in Section 4 and are evaluated in Section 6. We then conclude in Section 7.

## 2 Background

### 2.1 Twitch Plays Pokemon

Twitch Plays Pokmon [2] is a "social experiment" and channel on the video streaming website Twitch, consisting of a crowdsourced attempt to play Game Freak and Nintendo's Pokmon video games by parsing commands sent by users through the channel's chat room [1]. Tens of thousands of users are online and active in the chatroom. Next to the chat, a live video displays the current state of the Pokemon game. Users can then input commands into the chat. The game will then execute these commands in a FIFO order. Examples of commands include: up, down, left, right, and start. Up, down, left, and right move the main character in the game and start brings up the main menu. The objective of Twitch Plays Pokemon is to collective receive and process commands from users and attempt to beat the game of Pokemon.

### 2.2 Dataset

We wrote an Internet Relay Chat robot that collects user messages, commands, time, and user information in realtime. The robot collected data for several months. It is currently 3 GB in size and contains over 50 million data points. A data point is defined as a user's chat message and metadata such as their username and timestamp at millisecond resolution. It is still collecting data at the time we wrote this project proposal.

Users can enter non-commands into the chat as well. When a major milestone is accomplished in the game, it is common to see an increase in positive and ALL CAPS messages. Since a human logs into the system with the same username, we are able to build profiles of each user and learn from their past messages and commands.

## 3 Feature Selection

### 3.1 Context

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi ac sapien interdum, dictum sem ut, pellentesque urna. Integer eu odio velit. Praesent id dolor at libero cursus porta id vel augue. Ut et volutpat sapien, nec pulvinar nibh. Ut egestas aliquet quam ac rhoncus. Nunc rhoncus lacus a risus porttitor mattis. In gravida et lorem eu euismod. Fusce egestas, velit nec vulputate luctus, velit odio cursus dolor, non dictum dolor sem vel eros. Ut consequat turpis in justo aliquam, a lacinia ante aliquet. Vivamus feugiat dui justo, ut vestibulum nunc aliquam vel. Cras id orci non dui suscipit consectetur. Cras euismod felis in egestas vestibulum. In dictum malesuada vehicula.

### 3.2 Consensus Agreement

## 4 Algorithms

### 4.1 $k$-means Clustering

Morbi rhoncus nisi neque, sed viverra dui luctus eget. Praesent at dictum ante. Nam sit amet eleifend nisl. Phasellus sit amet neque vel diam facilisis consectetur sed nec ipsum. Mauris posuere elit eu libero interdum blandit. Aliquam elementum magna ac tincidunt tempor. Vivamus non orci

sollicitudin, auctor urna ut, consequat neque. Phasellus et molestie orci, nec gravida est. Nullam arcu nisi, ornare adipiscing condimentum in, ultricies sit amet nunc. Donec dictum lacinia neque, nec elementum augue eleifend id. Duis ultrices mattis velit eu lacinia. Morbi adipiscing consequat dolor, eu semper tellus hendrerit ut. Nullam et risus faucibus, mattis sem vitae, mattis lorem. Aenean fermentum eu turpis sed aliquet.

## 4.2 Point Anomaly Reduction

# 5 Distance Measures

## 5.1 Euclidean Distance

## 5.2 Local Outlier Factor

## 5.3 Multi-Granularity Deviation Factor

# 6 Evaluation

## 6.1 Cross Validation

## 6.2 Online Performance

# 7 Conclusion

## References

[1] Wikipedia. `https://www.wikipedia.org/`.

[2] Twitch plays pokemon. `http://www.twitch.tv/twitchplayspokemon`, April 2014.

[3] T. Campbell. Internet trolls. `http://web.archive.org/web/20011026130853/http://members.aol.com/intwg/trolls.htm`, July 2001.

[4] H. Fosdick. Why people troll and how to stop them. `http://www.osnews.com/story/25540`, January 2012.

[5] C. Heilmann. De-trolling the web: dont post in anger. `http://christianheilmann.com/2012/06/04/de-trolling-the-web-dont-post-in-anger/`, June 2012.