# Preventing Poisoning Attacks On AI Based Threat Intelligence Systems

Submitted by grigby1 on Wed, 11/04/2020 - 2:13pm

| | |
|---|---|
| Abstract | As AI systems become more ubiquitous, securing them becomes an emerging challenge. Over the years, with the surge in online social media use and the data available for analysis, AI systems have been built to extract, represent and use this information. The credibility of this information extracted from open sources, however, can often be questionable. Malicious or incorrect information can cause a loss of money, reputation, and resources; and in certain situations, pose a threat to human life. In this paper, we use an ensembled semi-supervised approach to determine the credibility of Reddit posts by estimating their reputation score to ensure the validity of information ingested by AI systems. We demonstrate our approach in the cybersecurity domain, where security analysts utilize these systems to determine possible threats by analyzing the data scattered on social media websites, forums, blogs, etc. |

# Citation Key khurana_preventing_2019

AI Poisoning AI systems Artificial Intelligence computer security cybersecurity domain Engines ensembled semi-supervised approach Human behavior learning (artificial intelligence) malicious information online social media poisoning attacks prevention pubcrawl resilience Resiliency Scalability security analysts security of data social networking (online) Support vector machines threat intelligence systems Twitter Web sites