

1 Article

2 A Survey of Attacks Against Twitter Spam Detectors 3 in an Adversarial Environment

4 Niddal Imam ¹

5 ¹ University of York; ni571@york.ac.uk

6 * Correspondence: ni571@york.ac.uk; Tel.: +44-7447-831-899

7 Received: date; Accepted: date; Published: date

8 **Abstract:** Online Social Networks (OSNs), such as Facebook and Twitter, have become a very
9 important part of many people's daily lives. Unfortunately, the high popularity of these platforms
10 makes them very attractive to spammers. Machine-learning (ML) techniques have been widely used
11 as a tool to address many cybersecurity application problems (such as spam and malware detection).
12 However, most of the proposed approaches do not consider the presence of adversaries that target
13 the defense mechanism itself. Adversaries can launch sophisticated attacks to undermine deployed
14 spam detectors either during training or the prediction (test) phase. Not considering these
15 adversarial activities at the design stage makes OSNs' spam detectors prone to a range of adversarial
16 attacks. This paper thus surveys the attacks against Twitter spam detectors in an adversarial
17 environment. In addition, a general taxonomy of potential adversarial attacks is proposed by
18 applying common frameworks from the literature. Examples of adversarial activities on Twitter
19 were provided after observing Arabic trending hashtags. A new type of spam tweet (*Adversarial
20 spam tweet*), which can be used to undermine deployed classifier, were found. In addition, possible
21 countermeasures that could increase the robustness of Twitter spam detectors against such attacks
22 are investigated.

23 **Keywords:** Twitter spam detection; Adversarial machine learning; Online Social Networks; Survey.
24

25 1. Introduction

26 Online Social Networks (OSNs), such as Facebook, WhatsApp, and Twitter have become a very
27 important part of daily life. People use them to make friends, communicate with each other, read the
28 news, and share their stories. The amount of information shared in these OSNs has continued to
29 increase over the past few years. One study shows that the number of profiles on Facebook, Twitter,
30 and LinkedIn reached more than two billion in 2016 [3].

31 Unfortunately, the high popularity of these OSNs has made them very attractive to malicious
32 users, or spammers. Spammers spread false information, propaganda, rumors, fake news, or
33 unwanted messages [29]. Spam is referred to as an unsolicited message that is received from a
34 random sender who has no relationship with the receiver. These messages may contain malware,
35 advertisements, or URLs directing the recipients to malicious websites [6]. Spamming on the Internet
36 first appeared in the 1990s in the form of email spam [3]. Although spam is prevalent in all forms of
37 online communication (such as email and the web), researchers' and practitioners' attention has
38 increasingly shifted to spam in OSNs due to the growing amount of spammers and the possible
39 negative effects on users [45, 6].

40 The first appearance of spam on Facebook was in 2008, while the first Twitter spam attack, in
41 which a number of Twitter accounts were hacked to spread advertisements, was in 2009 [50, 57]. On
42 Twitter, spammers tweet for several reasons, such as to spread advertisements, disseminate
43 pornography, spread viruses, phishing, or simply just to compromise a system's reputation. [7].
44 Furthermore, [26] added that a tweet is considered spam if it is not composed purely of text. Instead,
45 it may contain a hashtag, a mention, a URL or an image. Various types of spam are found in OSNs,
46 including textual pattern spam [58], image spam [8, 10], URL-based spam [53], and phone number-

47 based spam [30]. Whilst most previous studies have focused on detecting the above types of spam,
48 few have attempted to detect advertisement spam. The authors in [44] categorized adversarial
49 advertisements as: counterfeit goods, misleading or inaccurate claims, phishing, arbitrage, and
50 malware. The diversity of spam in OSNs makes it very hard for any single existing method to detect
51 most spam [27]. Several reported incidents show the danger of spammers in OSNs. For example, a
52 number of NatWest bank customers were victims of a phishing attack on Twitter that used spam
53 tweets that looked very similar to those from the official NatWest customer support account [3]. A
54 recent study noted that the increase in the number of OSN spammers, who distribute unsolicited
55 spam and advertise untrustworthy products, has an effect on the public's perception of companies,
56 which can eventually lead to people's opinion becoming biased [24].

57 The issue of spamming over OSNs has become an area of interest for many researchers. Many
58 solutions have been proposed to detect spam using techniques such as blacklisting and whitelisting,
59 Machine Learning (ML) and others. ML techniques have been shown to be effective when deployed
60 to solve security issues in different domains, such as email spam filters, intrusion detection systems
61 (IDSs), and malware detectors [21]. ML techniques aim to automatically classify messages as either
62 spam or non-spam. Various OSN spam detectors have been developed using ML algorithms,
63 including Supervised Vector Machine (SVM) [7], Random Forests (RF) [39, 55] and, more recently,
64 Deep Neural Networks [6].

65 Despite the success of these algorithms in detecting spam, the presence of adversaries
66 undermines their performance. These algorithms are vulnerable to different adversarial attacks
67 because they were not designed for adversarial environments [48, 2, 12]. The traditional assumption
68 of stationarity of data distribution in ML is that the dataset used for training a classifier (such as SVM
69 or RF) and the testing data (the future data that will be classified) have a similar underlying
70 distribution. This assumption is violated in the adversarial environment, as adversaries are able to
71 manipulate data either during training or before testing [46, 2].

72 Studying the robustness of OSNs' spam detectors against adversarial attacks is crucial. The
73 security of ML techniques is a very active area of research. Whilst several studies have examined the
74 security of IDS, email filters, and malware detectors, few have investigated the security of OSNs'
75 spam detectors. To the best of the researcher's knowledge, there is no survey of adversarial attacks
76 against OSNs' spam detectors. Recent studies suggest that if a secure system is to be achieved,
77 predicting potential attacks before they occur would help to develop suitable countermeasures [12].
78 Thus, the main goal of this paper is to present a comprehensive overview of different possible attacks,
79 which is the first step towards evaluating the security of OSNs' spam detectors in an adversarial
80 environment. This paper provides a general survey of the possible adversarial attacks against OSNs'
81 spam detectors. In addition, potential defense mechanisms that could reduce the effect of such attacks
82 are investigated. Ideas proposed in the literature were generalized to identify potential adversarial
83 attacks and countermeasures. Twitter, which is one of the most popular OSN platforms, was used as
84 a case study, and all examples of attacks were taken from Twitter. Examples of spam tweets that can
85 be used by an adversary to attack Twitter spam detectors were provided. This kind of spam tweet is
86 called Adversarial spam tweet.

87 The remainder of this survey is structured as follows: Section 2 describes previous research on
88 Twitter spam detection. Section 3 provides an overview of adversarial machine learning. Section 4
89 surveys the adversarial attacks that could be used against Twitter spam detectors and presents a
90 proposed taxonomy of such attacks. Section 5 briefly discusses possible defense strategies and
91 countermeasures. The conclusion and future works are presented in Section 6.

92 2. Techniques for Twitter spam detection

93 Twitter and the research community have proposed a number of spam detectors to protect users.
94 Twitter spam detection approaches can be divided into automated approaches, including machine
95 learning, and non-automated approaches that require human interaction [55].

96 Researchers who use ML approaches build their models by employing some of the common
97 spam detection techniques. Based on surveys in [36, 34], Twitter spam detectors can be classified into

98 four categories: user-based, content-based, hybrid-based, and relation-based techniques. User-based
 99 techniques are also referred to as account-based and classify tweets based on an account's features
 100 and other attributes that provide useful information about users' behavior. Content-based techniques
 101 use the content of a tweet, such as the linguistic properties of the text or the number of hashtags in
 102 the tweet, for classification. Hybrid techniques use a combination of user-based and content-based
 103 features. The last category was proposed to detect spam in real-time, in contrast to user-based
 104 techniques, which can only detect spam after a message has been received. Relation-based techniques
 105 can detect a tweet immediately if it is received from an unknown sender. The features used in these
 106 techniques are distance and connectivity (see Table 1).

107
 108 **Table 1:** Feature categories and description [34, 37].

Feature Category	Feature Name	Description
Account-based features	account_age	The number of days since the creation of an account.
	no_followers	The number of followers of an account.
	no_friends	The number of friends an account has.
	no_favorites	The number of favorites an account received.
	no_lists	The number of lists an account is a member of.
	no_reputation	The ratio of the number of followers and the total followers and friends of an account.
	no_statuses	The number of tweets an account has.
Tweet content-based features	no_words	The number of words in a tweet.
	no_chars	The number of characters in a tweet.
	no_hashes	The number of hashtags in a tweet.
	no_urls	The number of URLs in a tweet.
	no_phone	The number of phone numbers in a tweet.
	no_mentions	The number of mentions in a tweet.
Relation-based features	Distance	The length of the distance between accounts.
	Connectivity	The strength of the relationship between accounts.

109
 110 Additionally, the authors in [56] categorized the methodologies used for detecting Twitter spam
 111 into three groups: syntax-based, feature-based, and blacklist-based detection (see Figure 1). Syntax-
 112 based detectors analyze the content of tweets, including linguistic features and shortened URLs, to
 113 determine whether the tweet is spam or non-spam. The second group, feature-based detectors,
 114 extract a set of statistical features from tweets to help utilized classifier to determine whether the
 115 tweet is spam or non-spam. This group uses a combination of techniques: account-based features,
 116 tweet-based features, and social graph features. Account-based features include account age and
 117 number of followers, while tweet-based features are the number of characters and the number of
 118 URLs. To overcome some of the weaknesses of account-based and tweet-based features, some recent
 119 studies have found that adopting a social graph to detect spam by analyzing mathematical features,
 120 such as social distance and connectivity between followers, is more robust. In the last group, blacklist-
 121 based detectors, accounts and tweets are blocked based on users' feedback or the URL's reputation.
 122 [28] presented the first study of the effectiveness of some of the techniques that have been used in the
 123 past to detect Twitter spam. Examples include spam behavior, clickthrough, and blacklists. The
 124 authors found that the blacklist methods (for example Google SafeBrowsing) are too slow at detecting
 125 new threats. They found that although 90% of victims visit spam URLs within the first two days, it
 126 would take four to 20 days for the URLs in spam tweets to be blacklisted. In another study, it was
 127 determined that blacklists can protect only a few users, and asserted that studying the regional
 128 response rate could improve spam detection [20]. Furthermore, to overcome the limitations of the
 129 blacklist, some preliminary studies have used heuristic rules to filter Twitter spam [20].
 130

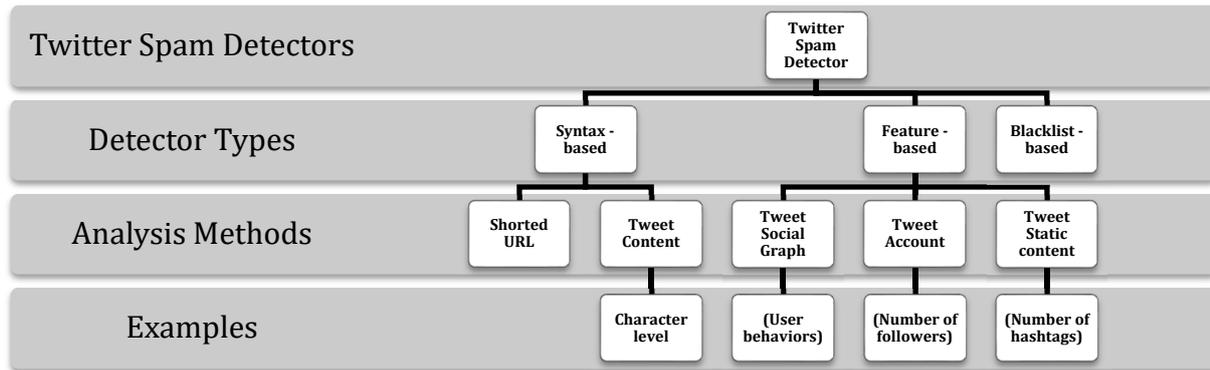


Figure 1 Types of Twitter Spam Detectors [56]

131 **Process of Detecting Spam through ML.** This process involves several steps. The first step
 132 involves collecting data from Twitter using its Streaming Application Programming Interference
 133 (API). This is followed by data pre-processing, which includes feature extraction, data labelling, and
 134 dataset splitting. However, for textual spam detectors, the pre-processing step may include more
 135 functions, such as tokenizing, removing stop words, and steaming. Extracting and selecting features
 136 from tweets or Twitter accounts helps the chosen ML classifier to distinguish between spam and non-
 137 spam, for example based on account age, the number of followers or friends, and the number of
 138 characters. Data labelling or ground truth is the process in which the collected data are labelled either
 139 manually or using a crowdsourcing site. The dataset then needs to be split into a training set and a
 140 test set. The last step involves training the chosen classification algorithm by using the labelled data,
 141 followed by performance evaluation, where the trained machine learning classifier can be used for
 142 spam detection [20, 1] (see Figure 2).
 143

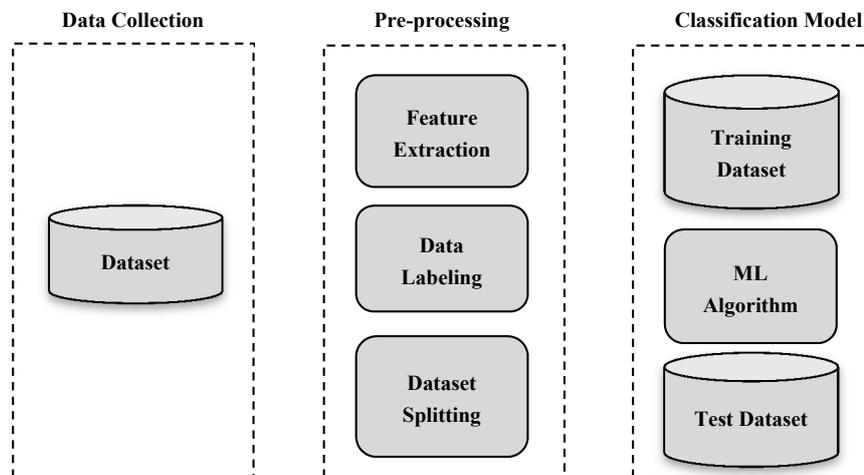


Figure 2: Learning Process of ML-based Spam Detection.

144 **Detecting Spam Tweets using ML Techniques.** As well as using blacklists for Twitter spam
 145 detection, other studies have used different ML techniques to detect spam tweets. As mentioned
 146 earlier, a number of steps are involved in the use of ML techniques to detect spam in Twitter; some
 147 of the important steps are discussed here. Several studies of Twitter spam detection developed their
 148 models by employing different ML algorithms, such as SVM, NB, and RF, of which RF has shown the
 149 best results in terms of detection accuracy. In [37], the authors compared and evaluated the detection
 150 accuracy, stability, and scalability of nine machine learning algorithms. The results showed that RF
 151 and C5.0 outperformed the other algorithms in terms of their superior detection accuracy. Similarly,
 152 a framework for detecting spam based on random forests was proposed in [39]. In addition, RF was
 153 chosen from five other algorithms in [55], as it has shown the best results in terms of evaluation
 154 metrics. Selecting features to help classify samples is as important a step as choosing the most suitable

155 algorithm for the required task. In [20], the authors collected a large dataset and labelled
 156 approximately 6.5 million spam tweets. They found that when using an imbalanced data set that
 157 simulated a real-world scenario, the classifiers' ability to detect spam tweets is reduced. On the other
 158 hand, when features are discretized, the performance of classifiers improves.
 159

160 **Detecting Spam Campaigns in Twitter.** Unlike the above spam detection models developed
 161 to detect a single spammer, some approaches can be used to detect campaign spam (spambots).
 162 According to [24], social spambots are a growing phenomenon and current spam detectors designed
 163 to detect a single spam account are not capable of capturing spambots. Although their study shows
 164 that neither humans nor existing machine learning models can detect spambots accurately, the result
 165 of an emerging technique deploying digital DNA has achieved a very promising detection
 166 performance. Similarly, [49] stated that methods designed to detect spam using account-based
 167 features cannot detect *crowdturfing accounts* (accounts created by crowdsourcing sites that have
 168 crowdsourcing and astroturfing characteristics). Another study [55] noted that spammers tend to
 169 create account bots to quickly reach their goals by systematically posting a large amount of spam in
 170 a short period of time. Consequently, they proposed an approach that uses the time property (for
 171 example the account creation date and tweet posting time), which cannot be modified by spammers,
 172 to reduce the creation of bots. [27] proposed an approach called Tangram, which uses a template-
 173 based model to detect spam in OSNs. After analyzing the textual pattern of a large collection of spam,
 174 they found that the largest proportion of spam is generated with an underlying template compared
 175 to other spam categorizes (for example paraphrase, no-content, and others).
 176

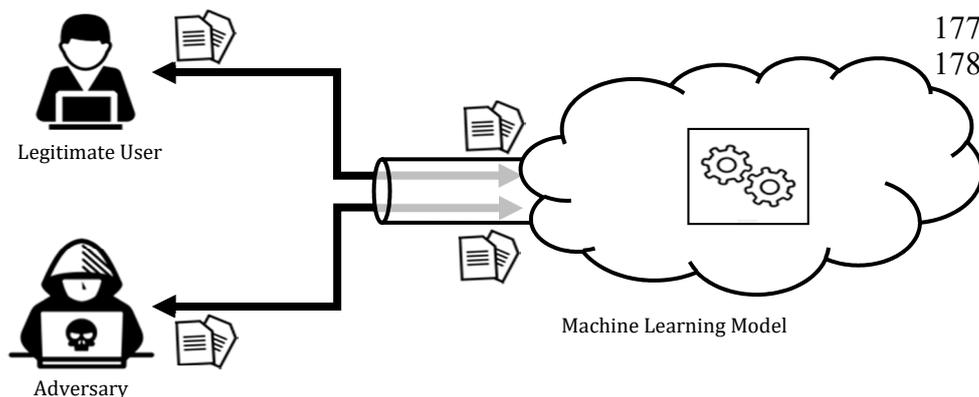


Figure 3: An adversary uses the same channel as legitimate users to exploit knowledge about the system.

179 **Security of Twitter spam detectors.** Despite the success and high level of accuracy of the models
 180 described here in detecting Twitter spam, they are nevertheless vulnerable as they were not
 181 developed for adversarial settings. A popular framework for evaluating secure learning was
 182 proposed in [4], and extended in [5, 31, 12]; it enables different attack scenarios to be envisaged
 183 against machine learning algorithms. The framework suggests the following steps: (1) identifying
 184 potential attacks against machine learning models based on the popular taxonomy (see Section 3.1);
 185 (2) simulating these attacks to evaluate the resilience of ML models; assumption that the adversary's
 186 attacks can be formed based on the adversary's goals, knowledge, and capabilities/resources. (3)
 187 investigating some possible defense strategies against these attacks. Defense against adversarial
 188 attacks is challenging as these attacks are non-intrusive in nature and an adversary launches their
 189 attacks using the same channel as legitimate users. Thus, defense strategies against these attacks
 190 cannot employ traditional encryption/security techniques [47]. Figure 3 demonstrates how an
 191 adversary can use the same channel as legitimate users to access an ML model and learn some of its
 192 characteristics. Designing *proactive* models rather than traditional *reactive* models is a necessity in the
 193 adversarial environment. Whereas reacting to detected attacks will never prevent future attacks,
 194 proactively anticipating adversaries' activities enables suitable defense methods to be developed

195 before an attack occurs [12]. This has motivated researchers to develop different attack scenarios
 196 against machine learning algorithms and classification models and propose some countermeasures.
 197 Table 2 shows an outline of recent spam detectors proposed in the literature.
 198

199 **Adversarial attacks against Twitter spam detectors.** In [54] the authors evaluated the security
 200 of a ML detector that is designed to detect spam generated by malicious crowdsourcing users of
 201 Weibo (the Chinese version of Twitter) against evasion and poisoning attacks. Their focus was on
 202 adversaries that use crowdsourcing sites to launch attacks. To study evasion attacks, two attacks were
 203 simulated: basic evasion, where an adversary has limited knowledge, and optimal evasion, where
 204 the adversary has perfect knowledge. The results show that an optimal evasion attack has a much
 205 higher impact than the basic one. However, in the real world, it is very difficult for adversaries to
 206 have perfect knowledge about the system. Thus, the less knowledge adversaries have about the
 207 system, the harder it is for them to evade detection. In causative attacks, two mechanisms for
 208 launching poisoning attacks are used. The aim of the first poisoning attack is to mislead the system
 209 by using crowdturfing admins to inject misleading samples directly into the training data. In the
 210 second poisoning attack, adversaries pollute training data by crafting samples that mimic benign
 211 users' behavior. After analyzing both attacks, it was found that injecting misleading samples causes
 212 the system to produce more errors than the second poisoning attack.

213 Another study by [42] analyzed the robustness of a Twitter spam detector called POISED against
 214 evasion and poisoning attacks. POISED is designed to distinguish between spam and non-spam
 215 messages based on the propagation of messages in each campaign. The authors suggested a
 216 poisoning attack, where the goal of an adversary is to contaminate training data by joining
 217 communities to alter their network and structure. The adversary posts manipulated messages in these
 218 compromised communities to mislead the system. Similarly, in an evasion attack, the adversary joins
 219 communities and resembles the propagation of non-spam messages to evade detection. The results
 220 show that the performance of POISED decreases when the percentage of compromised communities
 221 increases in both attacks. Thus, the authors suggest that if the adversary is to successfully attack
 222 systems, he or she needs to have a perfect knowledge about the structure and network of the targeted
 223 community.

224 The above studies suggest that the adversary's level of knowledge about the deployed system
 225 plays a very important role in determining the success of the attack.
 226

Table 2: Outline of some recent techniques used for detecting spam in Twitter: some of these works are discussed in Section 2.

Title	Methodology	Type of Spam	Type of Detector	Learning Approach	Results/Accuracy
6 Million Spam Tweets - A Large Ground Truth for Timely Twitter Spam Detection [62]	Different ML algorithms were used; balanced and imbalanced datasets were tested.	Spam tweet	Feature-based	Supervised	RF outperforms other algorithms.
A Hybrid Approach for Spam Detection for Twitter [63]	J48, Decorate and Naive-Bayes (NB).	Spam tweet	Feature-, user-, and graph-based	Supervised	J48 outperforms other algorithms.
Leveraging Time for Spammers Detection on Twitter [55]	Time-based features were used, and different ML algorithms were tested.	Spam tweet	Feature-based	Supervised	RF outperforms other algorithms.
Twitter spam detection based on Deep Learning [56]	Deferent ML algorithms with Word2Vector technique were used.	Spam tweet	Syntax-based	Supervised	RF with Worrdd2Vec outperforms other algorithms.
Semi-supervised spam detection (S3D) [64]	Utilizes four lightweight detectors (supervised and unsupervised) to detect spam tweets and updates the models periodically in batch mode.	Spam tweet	Feature-based and Blacklist	Semi-supervised	Confidential labeling process, which uses blacklisted, near-duplicated, and reliable non-spam tweets, makes the deployed classifier

						more efficient when detecting new spam tweets.
CrowdTarget: Target-based Detection of Crowdturfing in Online Social Networks [49]	Detects spam tweets that received retweets from malicious crowdsourcing customers. It uses four new retweet-based features and KNN as a classifier.	Spam tweet	Feature-based	Supervised		CrowdTarget detects malicious retweets created by crowdturfing users with a True Positive Rate of 0.98 and False Positive Rate of 0.01.
Beating the Artificial Chaos - Fighting OSN Spam using Its Own Templates [58]	Detects template-based spam, paraphrase spam, and URL-based spam.	Spam campaign	Syntax-based	Supervised		Template-detection outperforms URL blacklist-detection.
POISED - Spotting Twitter Spam Off the Beaten Paths [42]	Detects spam campaign based on community and topic of interest.	Spam campaign	Syntax-based	Supervised and unsupervised		NB, SVM, and RF all achieve about 90% detection accuracy.
Collective Classification of Spam Campaigners on Twitter: A Hierarchical Meta-Path Based Approach (HMPS) [65]	Builds heterogeneous networks and detects nodes connected by the same phone number or URL.	Spam campaign	Social graph-based	Supervised		HMPS outperforms feature- and content-based approaches. Prediction accuracy improves when using HMPS with feature- and user-based approaches.
Social Fingerprinting - Detection of Spambot Groups Through DNA-Inspired Behavioral Modeling [61]	Models users' behaviour using DNA fingerprinting technique to detect spambots.	Spam campaign	Social graph-based	Supervised and unsupervised		The results show that the proposed approach achieves better detection accuracy when using a supervised learning approach.

227 3. Adversarial Machine Learning

228 Adversarial ML is a research field that investigates the vulnerability of ML to *adversarial examples*,
 229 along with the design of suitable countermeasures [13]. Adversarial examples are inputs to ML that
 230 are designed to cause incorrect output [15]. The term was first introduced in [51] and used for
 231 computer vision, but in the context of spam and malware detection, the term *evasion attacks* is used in
 232 [12]. This section is going to discuss about different adversarial attacks and countermeasures. Table
 233 3 and 4 outline recent works in adversarial machine learning.
 234

Table 3: Outline of different adversarial attacks. These studies are discussed in Section 3.

Type of Influence	Title	Name of Attack	Attack Target	Attack Method
Causative	Manipulating Machine Learning: Poisoning Attacks and Countermeasures for Regression Learning [32]	Poisoning	Regression Learning	Optimization-based poisoning attack, in which different optimization approaches were used. Statistical-based poisoning attack (StatP) that queries a deployed model to find an estimate of the mean and covariance of training data.
	Support vector machines under adversarial label noise [59]	Label Flipping	SVM	Two different label flipping attacks were used: random and adversarial label flips.
	Curie - A method for protecting SVM Classifier from Poisoning Attack [35]	Label Flipping	SVM	Two label flipping attack were used. In the first, the loss maximization framework was used to select points that needed their label to be flipped. In second attack, the selected data points are moved to other points in the feature space.
	Adversarial Machine Learning [31]	Dictionary	Spam filter	An adversary builds a dictionary of tokens learned from the targeted model, and then sends attack messages to cause misclassification.

	Thwarting Signature Learning by Training Maliciously [41]	Red Herring	Polymorphic worm signature generation algorithms	An adversary sends messages with fake features to trick the deployed model.
	Man vs. Machine: Practical Adversarial Detection of Malicious Crowdsourcing Workers [54]	Poisoning	NB, BN, SVM, J48, RF	Two types of poisoning attacks were performed: Injecting misleading samples and altering training data.
Exploratory	Data Driven Exploratory Attacks on Black Box Classifiers in Adversarial Domains [47]	Anchor Points (AP) and Reverse Engineering attacks (RE)	SVM, KNN, DT, RF	AP attack is not affected by the chosen model (linear or non-linear), unlike RE, which is affected when a defender uses DT or RF.
	Evasion Attacks against Machine Learning at Test Time [10]	Evasion	SVM, Neural Network	A gradient-descent evasion attack was proposed.
	Good Word Attacks on Statistical Spam Filters [38]	Good Word	NB, Maximum entropy filter	Active and passive good word attacks against email spam filters were evaluated.
	Adding Robustness to Support Vector Machines Against Adversarial Reverse Engineering [2]	Reverse Engineering	SVM	Three different query selection methods, which help learn the decision boundary of deployed classifier, were used. Random, selective, and uncertainty sampling.
	Man vs. Machine: Practical Adversarial Detection of Malicious Crowdsourcing Workers [54]	Evasion	NB, BN, SVM, J48, RF	Two evasion attack were launched: Basic evasion attack and Optimal evasion attack, where an adversary knows features Needs to be altered.

235

236

3.1 Taxonomy of attacks against ML

237

A popular taxonomy proposed in [4, 5, 12] categorized attacks against ML systems along the three following axes:

238

239

The Attack INFLUENCE

240

- **Causative:** the attack influences the training data to cause misclassification.

241

- **Exploratory:** the attack exploits knowledge about the deployed classifier to cause misclassifications without influencing training data.

242

243

The Type of SECURITY VIOLATION

244

- **Integrity violation:** an adversary evades detection without compromising normal system operations.

245

246

- **Availability violation:** an adversary compromises the normal system functionalities available to legitimate users.

247

248

- **Privacy violation:** an adversary obtains private information about the system (such as its users, data, or characteristics) by reverse-engineering the learning algorithm.

249

250

The Attack SPECIFICITY

251

- **Targeted** attacks focus on a particular instance.
- **Indiscriminate** attacks encompass a wide range of instances.

252

253

254

The first axis, which is the attack influence, divides an adversary's capability to influence a classifier's learning systems into causative and exploratory. The influence is causative if an adversary misleads the deployed classifier by contaminating (poisoning) the training data by injecting carefully crafted samples into it. In contrast, the influence is exploratory if an adversary gains knowledge about the deployed classifier to cause misclassification at the testing phase without influencing training data.

257

258

259

260

The second axis describes the type of security violation committed by an adversary. The security violation can be an integrity violation if it enables an adversary to bypass the deployed classifier as a false negative. In addition, the attack can violate the model's availability if it creates denial of service, misclassifying non-spam samples as false positives, or if it prevents legitimate users from accessing

261

262

263

264 the system. The security violation can be a privacy violation if it allows an adversary to exploit
 265 confidential information from the deployed classifier.

266 The third axis of the taxonomy refers to the specificity of an attack. In other words, it indicates
 267 how specific an adversary's goal is. The attack specificity can be either targeted or indiscriminate,
 268 depending on whether the attack causes the classifier to misclassify a single or few instances, or
 269 undermines the classifier's performance on a larger set of instances.
 270

Table 4: Outline of techniques used for mitigating adversarial attacks: all of these works are discussed in Section 3

Type of Influence	Title	Name of Attack	Type of Classifier	Defense Category	Defense Method
Causative	Mitigating Poisoning Attacks on Machine Learning Models: A Data Provenance Based Approach [60]	Poisoning	SVM	Data Sanitization	Filtering out poisoned data from the training dataset using a provenance framework that records the lineage of data points.
	Curie- A method for protecting SVM Classifier from Poisoning Attack [35]	Poisoning	SVM	Data Sanitization	The data are clustered in the feature space, and the average distance of each point from the other points in the same cluster is calculated, with the class label considered as a feature with proper weight. The data points with less than 95% confidence are removed from the training data.
	Bagging Classifiers for Fighting Poisoning Attacks in Adversarial Classification Tasks [9]	Poisoning	Bagging and weighted bagging ensembles	Data Sanitization	Using an ensemble construction method (bagging) to remove outliers (adversarial samples) from training dataset.
	Data sanitization against adversarial label contamination based on data complexity [18]	Label Flipping	SVM	Data Sanitization	Data complexity, which measures the level of difficulty of classification problems, was used to distinguish adversarial samples in the training data.
	Support vector machines under adversarial label noise [59]	Label Flipping	SVM	Robust learning	Adjusting the kernel matrix of SVM depending on noise (adversarial) samples' parameters increases the robustness of the classifier.
Exploratory	Manipulating Machine Learning: Poisoning Attacks and Countermeasures for Regression Learning [32]	Poisoning	Regression Learning	Robust learning	The TRIM algorithm, which regularized linear regression by applying trimmed optimization techniques, was proposed
	Robust support vector machines against evasion attacks by random generated malicious samples	Evasion	SVM	Robust learning	Trains the SVM classifier with random malicious samples to enclose the decision function.
	Adding Robustness to Support Vector Machines Against Adversarial Reverse Engineering [2]	Reverse Engineering	SVM	Randomization	Learning a distribution of classifiers and picking a decision boundary randomly makes reverse engineering attacks harder to launch.
	Handling adversarial concept drift in streaming data [48]	Evasion	SVM	Disinformation	Hiding the importance of features and using an ensemble of classifiers.
	Adversarial Pattern Classification using Multiple Classifiers and Randomization [11]	Evasion	Spam Filter, SVM, NB	Multiple Classifiers and Randomization	Multiple Classifiers Strategy MCS, where different classifiers are trained by different features to randomize a model's decision boundary.

271

272 **3.2 Common Types of Threat Models**

273 In this regard, after presenting the taxonomy of attacks against ML systems, the next step
 274 towards identifying potential attack scenarios is threat modelling, which involves defining an
 275 adversary's goal, knowledge, and capability [4, 5, 12]. According to the above taxonomy, the
 276 attacker's goal may be based on the type of security violation (integrity, availability, or privacy), and
 277 on the attack specificity (targeted or indiscriminate). For instance, the adversary's goal could be to
 278 violate the system's integrity by manipulating either a specific instance or different instances. An
 279 attacker's level of knowledge about the classifier varies, and may include perfect-knowledge (white
 280 box attack), limited-knowledge (grey box attack), or zero-knowledge (black box attack). Attacker
 281 capability can involve either influencing training data (causative attack) or testing data (exploratory
 282 attack).

283

284 **3.3 Adversarial Attacks and Defense Strategies**

285 The existing literature on adversarial ML provides different attack examples and defense
 286 methods for both adversarial attack types (causative and exploratory). This section reviews common
 287 attack examples and some defense strategies against these attacks (see Table 5).

288

289

Table 5: Common Adversarial Attacks and Defenses

	Causative Attack	Exploratory Attack
Attack	Poisoning	Probing
	Red Herring	Evasion
	Label-Flipping	Reverse Engineering
		Good Words Attack
Defense	RONI	Randomization
	Game Theory based	Disinformation
	Multiple Learners	

290

291 **3.3.1 Causative Attacks**

292 One of the most common types of causative attack is a *poisoning attack*, in which an adversary
 293 contaminates the training dataset to cause misclassification [5]. An adversary can poison training
 294 data by either directly injecting malicious samples or sending a large number of malicious samples
 295 to be used by the defender when retraining the model [54]. A *label-flipping attack* is another example
 296 of a causative attack. Here, an adversary flips the label of some samples and then injects these
 297 manipulated samples into the training data. Different methods are used to perform this
 298 attack. Adversaries can either select samples that are nearest to or farthest from a classifier's decision
 299 boundary and flip their label [35]. The easiest method is to randomly flip the label of some samples
 300 that might be used for retraining. In [59], it was shown that randomly flipping about 40% of the
 301 training data's labels decreased the prediction accuracy of deployed classifier. A *red herring* attack is
 302 a type of causative attack in which the adversary adds irrelevant patterns or features into the training
 303 data to mislead the classifier to focus on these irrelevant patterns [41, 2]. Defense against causative
 304 attacks is challenging because ML classifiers need to be retrained periodically to adapt to new
 305 changes. Retraining the classifier makes it vulnerable as the data used for retraining is collected from
 306 an adversarial environment [35].

307

308 **3.3.2 Causative Defense Methods**

309 Although preventing these attacks is difficult, there are some defense methods proposed in the
 310 literature that can reduce the effect of these attacks. Defense methods against causative attacks may
 311 rely on *Game Theory*, where the defense problem is modeled as a game between the adversary and
 312 the classifier [2, 23, 25, 14]. *Data sanitization* methods focus on removing contaminated samples that
 313 have been injected by an adversary from a training dataset before training a classifier, while *robust*
 314 *learning* focuses on increasing the robustness of a learning algorithm to reduce the influence of

315 contaminated samples [18]. *Reject-on-negative-impact* (RONI) is one of the simplest and most effective
316 defense methods against causative attacks, and is considered to be a data sanitization method. In
317 RONI, all the training data go through preliminary screening to find and reject samples that have a
318 negative impact on the classification system. To distinguish between contaminated and untainted
319 samples, a classifier is trained using base training data before adding suspicious samples to the base
320 training data and train another classifier. The prediction accuracy for both classifiers compared to
321 over labelled test data is evaluated. If adding suspicious samples to the training data reduces the
322 prediction accuracy, these samples must be removed [31]. Another defense method involves using
323 *multiple classifiers*, which has been shown to reduce the influence of poisoned samples in training data
324 [9].

325 326 3.3.3 Exploratory Attacks

327 The most popular types of exploratory attacks are *evasion* and *reverse engineering*. Both attacks
328 start with a *probing attack*, in which an adversary sends messages to reveal some information about
329 the targeted classifier. Once the adversary gains some knowledge about the system, he or she can
330 either carefully craft samples that can evade the system (an evasion attack), or use that information
331 to build a substitute system (a reverse engineering attack) [4]. Furthermore, a *Good Word Attack* is a
332 type of exploratory attack in which the adversary either adds or appends words to spam messages
333 to evade detection. Good Word attacks can be passive or active. In a passive attack, the adversary
334 constructs spam messages by guessing which words are more likely to be bad or good (for example
335 a dictionary attack). In an active attack, the adversary has access to a targeted system that enables
336 him or her to discover bad and good words [38].

337 338 3.3.4 Exploratory Defense Methods

339 As with causative attacks, it is difficult to prevent exploratory attacks from happening because
340 in most cases systems cannot differentiate between messages sent for a legitimate purpose and those
341 sent to exploit the system. However, there are currently two common defense methods: *disinformation*
342 and *randomization*. In disinformation methods, the defender's goal is to hide some of the system's
343 functions (for example classification algorithms or features used by the classifier) from an adversary.
344 In contrast, in randomization methods, the defender's aim is to randomize the system's feedback to
345 mislead an adversary. [4]

346 Although most of these attack strategies and defense methods were proposed for domains such
347 as email spam filtering, IDS, and malware detection, the underlying approach can be applied in
348 Twitter spam detectors. The following section examines some of these techniques in the context of
349 Twitter spam detectors.

350 3. Taxonomy of Attacks Against Twitter Spam Detectors

351 This section surveys attacks against Twitter spam detectors in an adversarial environment.
352 Examples of adversarial spam tweets that can be used by adversaries to attack Twitter are also
353 provided.

354 355 4.1 Methodology

356 Different hypothetical attack scenarios against Twitter spam detectors are proposed. Attack
357 tactics were formularized based on the framework of the popular attack taxonomy presented in [4]
358 [5] that categorizes attacks along three axes: influence, security violations, and specificity. This
359 framework was extended in [12] to derive the corresponding optimal attack strategy by modelling
360 an adversary's goal, knowledge and capability. The adversary's goals considered in this study are
361 either to influence training or test data or to violate the system's integrity, availability, or privacy.
362 The adversary's knowledge is considered as: perfect knowledge (white-box attack) and zero-
363 knowledge (black-box attack). This ensures that both the worst-case and best-case scenarios are
364 considered for the adversary when attacking spam detectors. The adversary's capability is based on
365 desired goals. For example, if the goal is to influence the training data, the adversary must be capable

366 of doing so. Examples of adversarial spam tweet were extracted from Arabic trending hashtags. The
367 number of spam tweets using Arabic trending hashtags was found to be high, the reasons for which
368 are beyond the scope of this study. However, it was found that there are very active spam campaigns
369 spreading advertisements for untrustworthy drugs, for example weight loss drugs, Viagra, and hair
370 treatment drugs, targeting Arabic-speaking users. The attack scenarios can be modelled as follows:

- 371 1. Categorizing attacks based on their influence and type of violation (such as causative integrity
372 attacks).
- 373 2. Identifying the attack's settings, which includes an adversary's goal, knowledge and capability.
- 374 3. Defining the attack strategy that provides potential attack steps.

375

376 4.2 Potential Attack Scenarios

377 Here, attacks against Twitter spam detectors were categorized into four groups: causative
378 integrity, causative availability, exploratory integrity, and exploratory availability attacks. Four
379 hypothetical attack scenarios are provided, and different examples for each category are presented.
380 Some spam tweets were extracted from Arabic hashtags to show how an adversary can manipulate
381 tweets.

382

383 4.2.1 Causative Integrity Attacks

384

385 **Example 1: Poisoning Attack**

386 In this attack scenario, an adversary attempts to influence training data to cause new spam to
387 bypass the classifier as false negatives. The settings of the attack scenario are as follows: The
388 **adversary's goal** is to compromise the integrity of Twitter spam detectors and the attack specificity
389 can be either targeted or indiscriminate. The **adversary's knowledge** is assumed to be perfect (white-
390 box attack). In terms of the **adversary's capability**, it is assumed that the adversary is capable of
391 influencing the training data. After defining the attack scenario's setting, the next step is the attack
392 strategy. A potential attack strategy is as follows:

- 393 • As the adversary's knowledge of the system is considered to be perfect, it is not necessary to send
394 probing tweets to gain knowledge.
- 395 • The adversary would carefully craft a large number of malicious tweets.
- 396 • The crafted tweets must resemble non-spam tweets and include both spam components, such as
397 malicious URLs, and non-spam components or words (see Figure 4).
- 398 • The adversary would then post these tweets randomly using different trending hashtags and
399 hope that these malicious tweets are used by Twitter when retraining their system.

400 Figure 4 shows an example of a spam tweet that has been carefully crafted and can be used to
401 poison training data. The spam tweet mimics non-spam tweets by avoiding the inclusion of any spam
402 words, telephone numbers, or hashtags. In addition, the account resembles a legitimate user's
403 account in terms of having a decent number of followers and friends, and has a profile photo and
404 description. This spam tweet bypasses Twitter's spam detector and could be used for retraining the
405 classifier.



Figure 4: A spam tweet resembling a non-spam tweet to poison training data.

406
407
408

Example 2: Probing and Red Herring Attack

409
410
411
412
413
414
415

As in [41], in this attack scenario, the adversary's aim is to mislead Twitter's spam detectors by influencing training data. **The adversary's goal** is to compromise the integrity and privacy of Twitter's spam detectors, and the attack specificity can be either targeted or indiscriminate. **The adversary's capability** is similar to the previous example. However, the **adversary's knowledge** about Twitter's spam detectors is assumed to be zero (black-box attack). Based on the scenario's settings, a potential attack strategy is as follows:

416
417
418
419
420
421
422

- As the adversary has zero-knowledge about the system, sending probing tweets to gain knowledge is required (privacy violation).
- A probing attack is an exploratory type of attack, and will be discussed in the next section.
- The adversary would craft samples with spurious or fake features and post these samples on trending hashtags to trick Twitter's spam detectors into using these samples for retraining.
- If Twitter spam detectors are trained on these samples, the adversary will discard these spurious features in future tweets to bypass the classifier.

423
424
425
426
427
428

Figure 5 (a) shows an example of a spam tweet that has a spurious feature (phone number). As the number of tweets that have a phone number have increased on Twitter, some proposed spam detectors suggest using a phone number as an indicator of spam tweets [1, 30]. However, Figure 5(b) shows how the adversary can trick Twitter into using a phone number as a feature, and avoid including phone numbers in his spam tweets. Instead, the adversary includes a phone number inside an image to evade detection.



429
430

Figure 5: (a) A spam tweet containing spurious feature (a mobile number).



Figure 5: (b) A spam tweet with a mobile number inside an image to evade detection.

431

432

433

434 Example 3: Probing and Label-Flipping Attack

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

454 Example 1: Poisoning Attack

455

456

457

458

459

460

461

462

463

464

465

466

467

468

4.2.2 Causative Availability Attack

In this type of attack, an adversary tends to influence training data to either subvert the entire classification process or to make future attacks (such as evasion attacks) easier. The settings of the attack scenario are as follows: the **adversary's goal** is to violate the availability of Twitter and the attack specificity can be either targeted or indiscriminate. The **adversary's knowledge** is assumed to be perfect (white-box attack). In terms of the **adversary's capability**, it is assumed that the adversary is capable of influencing the training data. After defining the attack scenario's setting, the next step is the attack strategy. A potential attack strategy is as follows.

- As the adversary's knowledge about the system is considered to be perfect, sending probing tweets to gain knowledge is not required.
- The adversary would carefully craft a large number of misleading tweets that consist of a combination of spam and non-spam components (see Figure 4).
- The adversary needs to contaminate a very large proportion of training data for this attack to be successful. Using crowdsourcing sites or spambots to generate contaminated tweets helps the adversary to launch such an attack.

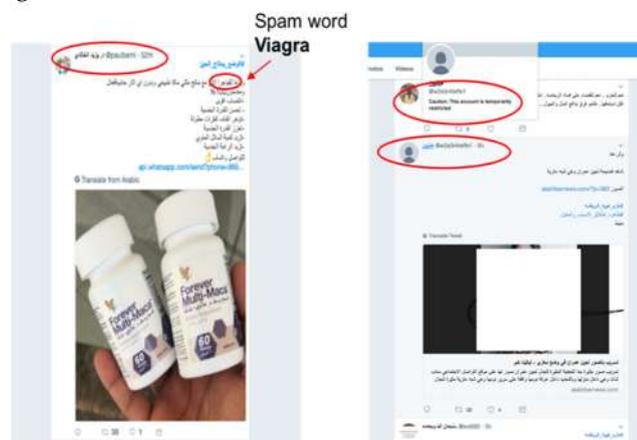
- 469 • The last step would be to post these tweets randomly using different trending hashtags to quickly
 470 spread these tweets in the hope that Twitter will use them when retraining their system.
 471

472 **Example 2: Dictionary Attack**

473 In this attack, as in [31], an adversary aims to corrupt the classification process by influencing
 474 training data and lead future legitimate tweets to be misclassified. The settings of the attack scenario
 475 are as follows: **an adversary's goal** is to violate the availability and integrity of Twitter spam
 476 detectors, and the attack specificity can be either targeted or indiscriminate. The **adversary's**
 477 **knowledge** is assumed to be perfect (white-box attack). In terms of the **adversary's capability**, it is
 478 assumed that the adversary is capable of influencing the training data. After defining the attack
 479 scenario's setting, the next step is the attack strategy. A potential attack strategy is as follows.

- 480 • As the adversary's knowledge about the system is considered to be perfect, sending probing
 481 tweets to gain knowledge is not required.
 482 • Based on the adversary's knowledge, he or she would build a dictionary of words or phrases
 483 frequently used by legitimate users and use this to craft malicious tweets.
 484 • The adversary would post tweets that contain a large set of tokens (non-spam words, phrases, or
 485 tweet structure) from the dictionary in trending hashtags.
 486 • If these tweets were used to train the system, non-spam tweets will be more likely to be classified
 487 as spam because the system will give a higher spam score for tokens used in the attack.

488 Figure 6 shows how a causative availability attack can affect Twitter spam detectors. The two
 489 spam tweets remain undetected for a long period of time because of the attack. As mentioned earlier,
 490 availability attacks overwhelm the system, which leads to difficulty in detecting spam tweets. A spam
 491 tweet on the left-hand side of the image below contains a very common spam word and should be
 492 very easily detected by the classifier, yet due to the attack, the tweet stays in place for longer than 52
 493 minutes. In addition, a spam tweet on the right-hand side remains undetected for longer than five
 494 hours, which is a very long time.



495
 496 **Figure 6:** Spam tweets bypass the detection system due to the availability attack.
 497

498 **4.2.3 Exploratory Integrity Attack**

500 **Example 1: Probing Attack**

501 In this attack scenario, the aim is to learn or expose some of the deployed classifier's
 502 functionalities without any direct influence over the training data. The settings of the attack scenario
 503 are as follows: **an adversary's goal** is to compromise the privacy of Twitter's spam detectors and the
 504 attack specificity can be either targeted or indiscriminate. The **adversary's knowledge** is assumed to
 505 be zero (black-box attack). As in [12], in terms of the **adversary's capability**, it is assumed that the
 506 adversary is only capable of influencing the testing data. After defining the attack scenario's setting,
 507 the next step is the attack strategy. A potential attack strategy is as follows.

- 508 • As the adversary does not have sufficient knowledge of how the Twitter spam detector works,
 509 sending probing tweets to gain knowledge is required.

- 510 • The adversary would send a large number of tweets, each with different features, to learn about
 511 the system (see Figure 7).
 512 • Based on the information that is learned, the adversary would carefully craft tweets to evade
 513 detection.
 514

515 Figure 7 shows an example of three spam tweets advertising the same weight-loss products.
 516 However, the adversary uses different features in each tweet. The first tweet consists of text, a URL,
 517 and an image, and the second has text and an image. The last one contains text only. The goal here is
 518 to learn how the classifier works. For example, if the first tweet is detected, the adversary will learn
 519 that a blacklist of URLs could be one of the features used by the classifier.



520
 521 **Figure 7:** An example of a probing attack.
 522

523 **Example 2:** Evasion Attack – Good Word Attack

524 In this attack scenario, the aim is to evade being detected by the deployed classifier without any
 525 direct influence over the training data. The settings of the attack scenario are as follows: **an**
 526 **adversary's goal** is to compromise the integrity of the Twitter spam detector and the attack specificity
 527 can be either targeted or indiscriminate. The **adversary's knowledge** is assumed to be perfect (white-
 528 box attack). In terms of the **adversary's capability**, as in [12], it is assumed that the adversary is only
 529 capable of influencing the testing data. After defining the attack scenario's setting, the next step is the
 530 attack strategy. A potential attack strategy is as follows.

- 531 • As the adversary's knowledge of the system is considered to be perfect, sending probing tweets
 532 to gain knowledge is not required.
 533 • Based on the adversary's knowledge, he or she would carefully craft tweets by modifying and
 534 obfuscating spam words (such as Viagra) or the tweet's features to evade detection (such as
 535 number of followers) (see Figure 8).

536 Figure 8 shows a spam tweet that has been carefully crafted to evade detection. The adversary
 537 avoids including any spam words in the text. Instead, the tweet contains a description of the drug
 538 (Viagra), and the spam word was inserted inside an image.
 539



Figure 8: Spam image tweet crafted to evade detection.

540

541

542

543 Example 3: Probing and Reverse Engineering Attacks

544 Evading the classifier without influencing the training data is the aim considered in this attack
 545 scenario. The scenario's settings are: **the adversary's goal** is to violate the integrity and privacy of
 546 Twitter's spam detectors, and the attack specificity can be either targeted or indiscriminate. The
 547 **adversary's capability** is similar to the previous example, but the **adversary's knowledge** about
 548 Twitter's spam detectors is assumed to be zero (black-box attack). Based on scenario's settings, a
 549 potential attack strategy is as follows:

- 550 • As the adversary has zero knowledge about the system, the first step would be to send probing
 551 tweets to learn how the system works (privacy violation).
- 552 • Based on the exploited knowledge, the adversary would build a substitute model that could be
 553 used for launching different exploratory attacks [2].
- 554 • Once the substitute model is built, the adversary would craft different spam tweets to evade
 555 detection, and spam tweets that successfully evade the model will be used against the Twitter
 556 spam detector.

557

558 3.2.4 Exploratory Availability Attack

559

560 Example 1: Denial of Service and Evasion Attack

561 In this attack scenario, the main aim is to evade the classifier by sending a large number of
 562 adversarial spam tweets to overwhelm the classifier without any direct influence over the training
 563 data. The settings of the attack scenario are as follows: **an adversary's goal** is to violate the availability
 564 and integrity of the Twitter spam detector and the attack specificity can be either targeted or
 565 indiscriminate. The **adversary's knowledge** is assumed to be perfect (white-box attack). In terms of
 566 the **adversary's capability**, as in [12], it is assumed that the adversary is only capable of influencing
 567 the testing data. After defining the attack scenario's setting, the next step is the attack strategy. A
 568 potential attack strategy is as follows.

- 569 • As the adversary has perfect knowledge about the system, sending probing tweets to gain
 570 knowledge is not required.
- 571 • Based on the gained knowledge, the adversary would carefully craft spam tweets. As the
 572 adversary cannot influence training data, the adversary would craft tweets that require more
 573 time to be processed by the classifier, such as image-based tweets [5].
- 574 • The adversary would then flood the system (for example a particular trending hashtag) with
 575 spam tweets to prevent users from reading non-spam tweets and cause difficulty in detecting
 576 spam tweets.

577 Figure 9 shows an example of an availability attack, where the adversary post a large number of
 578 spam tweets from a different account that only contains an image. As mentioned earlier, image

579 processing overwhelms the deployed classifier and causes a denial of service. In this kind of attack,
580 the adversary may use crowdsourcing sites or spambots to generate spam tweets.



Figure 9: An adversary floods the hashtag with spam tweets.

581
582
583

584 Example 2: Probing and Denial of Service Attacks

585 The aim of this attack scenario is similar to the previous example, but the scenario's settings are
586 slightly different. **The adversary's goal** is to violate the integrity, availability and privacy of Twitter's
587 spam detectors, and the attack specificity can be either targeted or indiscriminate. The **adversary's**
588 **capability** is similar to the previous example, but the **adversary's knowledge** about Twitter's spam
589 detectors is assumed to be zero (black-box attack). Based on scenario's settings, a potential attack
590 strategy is as follows:

- 591 • As the adversary has zero-knowledge about the system, the first step would be to probe the
592 classifier with some tweets to learn how it works.
- 593 • Based on the exploited knowledge, the adversary would craft a large number of spam tweets and
594 post them in a specific hashtag to cause denial of service and make future attacks easier. [2].

595
596
597
598
599
600
601
602

All attack examples can be either targeted if an adversary focuses on a specific spam tweet (for
example URL-based spam, or weight-loss ads), or indiscriminate, if an adversary targets multiple
types of spam tweet (such as URL-based and advertisements). Although presented adversarial spam
tweets look very similar to spam tweets that targeted users, this special type of spam tweets need to
be studied more as it aims to subvert Twitter spam detectors. Table 6 summarizes the taxonomy of
potential attacks.

Table 6: Taxonomy of potential attacks against Twitter spam detectors

Type of Influence	Potential Attack	Security Violation	Specificity
Causative	Poisoning Attack	Integrity	Targeted/
	Probing and Red Herring Attack	Integrity & Privacy	Indiscriminate
	Probing and Label-Flipping Attack	Integrity & Privacy	
	Poisoning Attack	Availability	
	Dictionary Attack	Availability & Integrity	
Exploratory	Probing Attack	Privacy	
	Good Word Attack	Integrity	
	Probing and Reverse Engineering Attacks	Integrity & Privacy	
	Denial of Service and Evasion Attack	Availability & Integrity	
	Probing and Denial of Service Attacks	Availability, Integrity & Privacy	

603

604 5 Potential Defense Strategies

605 This section discusses some possible defense strategies against adversarial attacks that can be
606 considered when designing a spam detector for Twitter. Some of the popular defense methods
607 proposed in the literature are discussed in the context of Twitter spam detection.

608

609 5.1 Defenses Against Causative Attacks

610 Existing approaches to defend against causative attacks focus on filtering or screening all the
611 training data before using them to update a deployed classifier, such as RONI, data sanitization
612 techniques, and bagging of classifiers. Although these methods have been shown to reduce the
613 influence of contaminated samples on training data, in some cases in which contaminated samples
614 overlap with untainted samples, discriminating between the two becomes very difficult [18]. Some
615 recent studies have suggested using a data collection oracle to retrain a deployed classifier [48, 33].
616 However, trusting an oracle to label training data could be problematic. The authors in [40] stated
617 that using crowdsourcing sites to label data might produce noisy data, thus increasing complexity.
618 Furthermore, Song et al. added that adversaries can increase the popularity of malicious tweets by
619 using artificial retweets generated by crowdsourcing workers [49]. Thus, developing a fully
620 automated model that can filter these poisoned samples is important. Nowadays, the trend is towards
621 fully automated systems to eliminate human errors. However, the above defense methods require
622 human interventions.

623

624 5.2 Defenses Against Exploratory Attacks

625 As mentioned in Section 3, the common defense methods against exploratory attacks are
626 disinformation and randomization. The goal in disinformation methods is to hide some of the
627 important information about the system from an adversary. Although determining the features used
628 by the classifier is not difficult, manipulating or mimicking all of these features may be impossible
629 for an adversary. Some features can be neither manipulated nor mimicked. In [31] and [55], the
630 authors found that time-based features (such as account age) are unmodifiable. Furthermore, the
631 authors in [54] discussed how altering some features comes at a cost, while others cannot even be
632 altered. For example, the number of tweets, and the number of followers and following, are features
633 that can easily be mimicked, and they might cause the adversary to create a large number of accounts
634 and buy lots of friends. On the other hand, profile and interaction features are much harder to alter.
635 Consequently, considering the robustness of selected features and applying the disinformation
636 method when designing a spam detector would help reduce the effect of adversaries' activities.
637 However, this cannot stop determined adversaries from trying every way possible to accomplish
638 their goals [46]. Furthermore, as stated in [2], relying on obscurity in an adversarial environment is
639 not a good security practice, as one should always overestimate rather than underestimate the
640 adversary's capabilities. In randomization, the defender's aim is to mislead the adversary by
641 randomizing system's feedback. Unlike the disinformation method, this strategy cannot prevent
642 adversaries from exploiting some information about the system, but makes it harder for them to gain
643 any information [4], especially in Twitter, where the adversary uses the same channel as benign users
644 to discover the system. This makes randomization methods less effective against exploratory attacks
645 in Twitter.

646 However, some recent studies have proposed an approach that can detect adversarial samples
647 using the deployed classifier's uncertainty in predicting samples' labels. In [48], the authors use
648 multiple classifiers (predict and detect) for detecting adversarial activities. Each classifier detects
649 samples that lie within a classifier's region of uncertainty (blind posts), where the classifier needs to
650 use its best guess. Then, if there is a disagreement between the two classifiers' output, the sample will
651 be tested with labelled samples for confirmation.

652 6 Conclusion and Future Work

653 The use of machine learning techniques in security applications has become very common. As
654 spam on OSNs is considered to be an adversarial problem, investigating the security of the machine

655 learning models used to detect spam is very important. Adversaries tend to launch different types of
656 attacks to evade detection by influencing the deployed model either at the training or test phase.
657 Recent studies have shown an increased interest in studying the security of machine learning in
658 domains such as IDS, malware detection, and email spam filters. However, the security of OSNs'
659 spam detectors has not been evaluated sufficiently.

660 The main contribution of this paper is to provide a general taxonomy of potential adversarial attacks
661 against Twitter spam detectors and a discussion on possible defense strategies that can reduce the
662 effect of such attacks. Examples of adversarial spam tweets that can be used by an adversary were
663 provided. This study is the first step towards evaluating the robustness of Twitter spam detectors, as
664 it identifies potential attacks against them. Hypothetical examples of possible attacks against Twitter
665 spam detectors were based on common frameworks proposed in [4, 5, 12]. In addition, defense
666 methods commonly proposed in the literature and ways of deploying these methods in the context
667 of Twitter spam detection were discussed.

668 Throughout the paper, a number of challenging issues were mentioned; future research needs
669 to focus on addressing them. Detecting image-based spam is an ongoing problem, as the processing
670 of images overwhelms classifiers and affects detection performance. Adversaries take advantage of
671 this issue, and the amount of image-based spam is increasing. Furthermore, spam detectors designed
672 for spam campaigns may fail to detect single spam and vice-versa. This issue can also be exploited
673 by adversaries when attacking spam detectors. Most proposed defense strategies can make attacks
674 against Twitter spam detectors very hard for adversaries, but, as most adversarial attacks are non-
675 intrusive [47], they cannot completely prevent attacks from happening.

676 In terms of a future direction, after identifying potential attacks against Twitter spam detectors,
677 the next step is to simulate some of these attacks to evaluate the robustness of Twitter spam detectors.
678 Evaluating the security of Twitter spam detectors experimentally will help design adversarial-aware
679 spam detectors that are more robust against adversarial activities.

680 **Acknowledgments:** I would like to warmly thank my supervisor Dr. Vasileios Vasilakis for his feedback and
681 suggestions.

682 References

- 683
- 684 1. Al Twaresh, N., Al Tuwajri, M., Al Moammar, A., Al Humoud, S., 2016, "Arabic Spam Detection in
685 Twitter", *The 2nd Workshop on Arabic Corpora and Processing Tools 2016 Theme: Social Media*, 2016
 - 686 2. Alabdulmohsin, I.M., Gao, X. & Zhang, X., 2014. Adding Robustness to Support Vector Machines Against
687 Adversarial Reverse Engineering. *Proceedings of the 23rd ACM International Conference on Conference on*
688 *Information and Knowledge Management*, pp.231–240.
 - 689 3. Al-Zoubi A., Alqatawna J. & Faris H. 2017, "Spam profile detection in social networks based on public
690 features", *2017 8th International Conference on Information and Communication Systems (ICICS)*, pp. 130.
 - 691 4. Barreno, M. et al., 2006. Can Machine Learning Be Secure? In *Proceedings of the 2006 ACM Symposium on*
692 *Information, Computer and Communications Security*. ASIACCS '06. New York, NY, USA: ACM, pp. 16–25.
 - 693 5. Barreno, M. et al., 2010. The security of machine learning. *Machine learning*, 81(2), pp.121–148.
 - 694 6. Barushka, A. & Hajek, P., 2018. Spam filtering using integrated distribution-based balancing approach and
695 regularized deep neural networks. *Applied Intelligence*, 48(10), pp.3538–3556.
 - 696 7. Benevenuto, F., Magno, G., Rodrigues, T. and Almeida, V., 2010, July. Detecting spammers on twitter. In
697 *Collaboration, electronic messaging, anti-abuse and spam conference (CEAS)* (Vol. 6, No. 2010, p. 12).
 - 698 8. Biggio, B., Fumera, G., et al., 2011a. A survey and experimental evaluation of image spam filtering
699 techniques. *Pattern recognition letters*, 32(10), pp.1436–1446.
 - 700 9. Biggio, B., Corona, I., et al., 2011b. Bagging Classifiers for Fighting Poisoning Attacks in Adversarial
701 Classification Tasks. In *Lecture Notes in Computer Science*. pp. 350–359.
 - 702 10. Biggio, B. et al., 2017. Evasion Attacks against Machine Learning at Test Time. *arXiv [cs.CR]*. Available at:
703 <http://arxiv.org/abs/1708.06131>.
 - 704 11. Biggio, B., Fumera, G. & Roli, F., 2008. Adversarial Pattern Classification using Multiple Classifiers and
705 Randomisation. *12th Joint IAPR International Workshop on Structural and Syntactic Pattern Recognition (SSPR*
706 *2008)*, 5342, pp.500–509.

- 707 12. Biggio, B., Fumera, G. & Roli, F., 2014. Security Evaluation of Pattern Classifiers under Attack. *Knowledge*
708 *and Data Engineering*, \dots, 26(4), pp.984–996.
- 709 13. Biggio, B. & Roli, F., 2017. Wild Patterns: Ten Years After the Rise of Adversarial Machine Learning. *arXiv*
710 [cs.CV]. Available at: <http://arxiv.org/abs/1712.03141>.
- 711 14. Bruckner, M., 2012. Static Prediction Games for Adversarial Learning Problems. *Journal of machine learning*
712 *research: JMLR*, 13, pp.2617–2654..
- 713 15. Buckman, J. et al., 2018. Thermometer Encoding: One Hot Way To Resist Adversarial Examples.
- 714 16. Cai, Z. et al., 2018. Collective Data-Sanitization for Preventing Sensitive Information Inference Attacks in
715 Social Networks. *IEEE Transactions on Dependable and Secure Computing*, 15(4), pp.577–590.
- 716 17. Chakraborty, M. et al., 2016. Recent developments in social spam detection and combating techniques: A
717 survey. *Information processing & management*, 52(6), pp.1053–1073.
- 718 18. Chan, P.P.K. et al., 2018. Data sanitization against adversarial label contamination based on data
719 complexity. *International Journal of Machine Learning and Cybernetics*, 9(6), pp.1039–1052.
- 720 19. Chan, P.P.K. et al., 2015. Spam filtering for short messages in adversarial environment. *Neurocomputing*,
721 155, pp.167–176.
- 722 20. Chen, C. et al., 2015. Asymmetric self-learning for tackling Twitter Spam Drift. In *2015 IEEE Conference on*
723 *Computer Communications Workshops (INFOCOM WKSHPs)*. pp. 208–213.
- 724 21. Chen, L., Ye, Y. & Bourlai, T., 2017. Adversarial Machine Learning in Malware Detection: Arms Race
725 between Evasion Attack and Defense. In *2017 European Intelligence and Security Informatics Conference*
726 *(EISIC)*. pp. 99–106.
- 727 22. Chen, S. et al., 2018. Automated poisoning attacks and defenses in malware detection systems: An
728 adversarial machine learning approach. *Computers & Security*, 73, pp.326–344.
- 729 23. Corona, I., Giacinto, G. & Roli, F., 2013. Adversarial attacks against intrusion detection systems: Taxonomy,
730 solutions and open issues. *Information sciences*, 239, pp.201–225.
- 731 24. Cresci, S. et al., 2017. The Paradigm-Shift of Social Spambots: Evidence, Theories, and Tools for the Arms
732 Race. In *Proceedings of the 26th International Conference on World Wide Web Companion*. International World
733 Wide Web Conferences Steering Committee, pp. 963–972.
- 734 25. Dalvi, N. et al., 2004. Adversarial Classification. In *Proceedings of the Tenth ACM SIGKDD International*
735 *Conference on Knowledge Discovery and Data Mining*. KDD '04. New York, NY, USA: ACM, pp. 99–108.
- 736 26. El-Mawass N. & Alaboodi S. 2016, "Detecting Arabic spammers and content polluters on Twitter", *2016*
737 *Sixth International Conference on Digital Information Processing and Communications (ICDIPC)*, pp. 53.
- 738 27. Gao, H. et al., 2014. Spam ain't as diverse as it seems: throttling OSN spam with templates underneath. In
739 *Proceedings of the 30th Annual Computer Security Applications Conference*. ACM, pp. 76–85.
- 740 28. Grier, C., Thomas, K., Paxson, V. and Zhang, M., 2010, October. @ spam: the underground on 140 characters
741 or less. In *Proceedings of the 17th ACM conference on Computer and communications security* (pp. 27-37). ACM
- 742 29. Gupta A.& Kaushal R. 2015, "Improving spam detection in Online Social Networks", *2015 International*
743 *Conference on Cognitive Computing and Information Processing(CCIP)*, pp. 1.
- 744 30. Gupta, P., Perdisci, R. & Ahamad, M., 2018. Towards Measuring the Role of Phone Numbers in Twitter-
745 Advertisised Spam. In *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*.
746 ACM, pp. 285–296.
- 747 31. Huang, L. et al., 2011. Adversarial Machine Learning. In *Proceedings of the 4th ACM Workshop on Security*
748 *and Artificial Intelligence*. AISec '11. New York, NY, USA: ACM, pp. 43–58.
- 749 32. Jagielski, M. et al., 2018. Manipulating Machine Learning: Poisoning Attacks and Countermeasures for
750 Regression Learning. , (1). Available at: <http://arxiv.org/abs/1804.00308>.
- 751 33. Kantchelian, A. et al., 2013. Approaches to adversarial drift. *Proceedings of the 2013 ACM workshop on*
752 *Artificial intelligence and security - AISec '13*, pp.99–110.
- 753 34. Kaur, P., Singhal, A. & Kaur, J., 2016. Spam detection on Twitter: A survey. In *2016 3rd International*
754 *Conference on Computing for Sustainable Global Development (INDIACom)*. pp. 2570–2573.
- 755 35. Laishram, R. & Phoha, V.V., 2016. Curie: A method for protecting SVM Classifier from Poisoning Attack.
756 *arXiv [cs.CR]*. Available at: <http://arxiv.org/abs/1606.01584>.
- 757 36. Lalitha, L.A., Hulipalled, V.R. & Venugopal, K.R., 2017. Spamming the mainstream: A survey on trending
758 Twitter spam detection techniques. In *2017 International Conference On Smart Technologies For Smart Nation*
759 *(SmartTechCon)*. pp. 444–448.

- 760 37. Lin G., Sun N., Nepal S., Zhang J., Xiang Y. & Hassan H. 2017, Statistical Twitter Spam Detection
761 Demystified: Performance, Stability and Scalability
- 762 38. Lowd, D. & Meek, C., 2005. Good Word Attacks on Statistical Spam Filters. *Conference on email and anti-*
763 *spam*.
- 764 39. Meda, C. et al., 2016. Spam detection of Twitter traffic: A framework based on random forests and non-
765 uniform feature sampling. In *2016 IEEE/ACM International Conference on Advances in Social Networks Analysis*
766 *and Mining (ASONAM)*. pp. 811–817.
- 767 40. Miller, B. et al., 2014. Adversarial Active Learning. Proceedings of the 2014 Workshop on Artificial
768 Intelligent and Security Workshop - AISec '14, (August), pp.3–14.
- 769 41. Newsome, J., Karp, B. & Song, D., 2006. Paragraph: Thwarting Signature Learning by Training Maliciously.
770 In *Lecture Notes in Computer Science*. pp. 81–105.
- 771 42. Nilizadeh, S. et al., 2017. POISED: Spotting Twitter Spam Off the Beaten Paths. In *Proceedings of the 2017*
772 *ACM SIGSAC Conference on Computer and Communications Security*. ACM, pp. 1159–1174.
- 773 43. Pawar, K. & Patil, M., 2015. Pattern classification under attack on spam filtering. In 2015 IEEE International
774 Conference on Research in Computational Intelligence and Communication Networks (ICRCICN). pp.
775 197–201.
- 776 44. Sculley, D. et al., 2011. Detecting adversarial advertisements in the wild. Proceedings of the 17th ACM
777 SIGKDD international conference on Knowledge discovery and data mining - KDD '11, p.274.
- 778 45. Sedhai, S. & Sun, A., 2015. HSpam14: A Collection of 14 Million Tweets for Hashtag-Oriented Spam
779 Research. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in*
780 *Information Retrieval*. ACM, pp. 223–232.
- 781 46. Sethi, T.S., Kantardzic, M., Lyua, L., et al., 2018a. A Dynamic-Adversarial Mining Approach to the Security
782 of Machine Learning. Available at: <http://dx.doi.org/10.1002/widm.1245>.
- 783 47. Sethi, T.S. & Kantardzic, M., 2018b. Data driven exploratory attacks on black box classifiers in adversarial
784 domains. *Neurocomputing*, 289, pp.129–143.
- 785 48. Sethi, T.S. & Kantardzic, M., 2018c. Handling adversarial concept drift in streaming data. *Expert systems*
786 *with applications*, 97, pp.18–40.
- 787 49. Song, J., Lee, S. & Kim, J., 2015. CrowdTarget. Proceedings of the 22nd ACM SIGSAC Conference on
788 Computer and Communications Security - CCS '15, (i), pp.793–804.
- 789 50. Stringhini, G., Kruegel, C. & Vigna, G., 2010. Detecting Spammers on Social Networks. In *Proceedings of the*
790 *26th Annual Computer Security Applications Conference*. ACSAC '10. New York, NY, USA: ACM, pp. 1–9.
- 791 51. Szegedy, C. et al., 2013. Intriguing properties of neural networks. *arXiv [cs.CV]*. Available at:
792 <http://arxiv.org/abs/1312.6199>.
- 793 52. Verma, M., Divya, D. & Sofat, S., 2014. Techniques to Detect Spammers in Twitter- A Survey. *International*
794 *Journal of Computer Applications in Technology*, 85(10), pp.27–32.
- 795 53. Wang, D. et al., 2013. Click Traffic Analysis of Short URL Spam on Twitter. In ICST. Available at:
796 <http://dx.doi.org/10.4108/icst.collaboratecom.2013.254084>
- 797 54. Wang, G. et al., 2014. Man vs. Machine : Practical Adversarial Detection of Malicious Crowdsourcing
798 Workers. *the 23rd USENIX Security Symposium*, pp.239–254.
- 799 55. Washha, M., Qaroush, A. & Sedes, F., 2016. Leveraging time for spammers detection on Twitter. In
800 *Proceedings of the 8th International Conference on Management of Digital EcoSystems*. ACM, pp. 109–116.
- 801 56. Wu, T. et al., 2017. Twitter Spam Detection Based on Deep Learning. In *Proceedings of the Australasian*
802 *Computer Science Week Multiconference*. ACSW '17. New York, NY, USA: ACM, pp. 3:1–3:8.
- 803 57. Yang, C., Harkreader, R. & Gu, G., 2013. Empirical Evaluation and New Design for Fighting Evolving
804 Twitter Spammers. *IEEE Transactions on Information Forensics and Security*, 8(8), pp.1280–1293.
- 805 58. Zhu, T. et al., 2016. Beating the Artificial Chaos: Fighting OSN Spam Using Its Own Templates. *IEEE/ACM*
806 *Transactions on Networking*, 24(6), pp.3856–3869.
- 807 59. Biggio, B., 2011. Support Vector Machines Under Adversarial Label Noise. *JMLR workshop and conference*
808 *proceedings*, 20, pp.97–112.
- 809 60. Baracaldo, N. & Chen, B., Mitigating Poisoning Attacks on Machine Learning Models: A Data Provenance
810 Based Approach. Available at: <https://dl.acm.org/ccs.cfm>.
- 811 61. Cresci, S. et al., 2018. Social Fingerprinting: Detection of Spambot Groups Through DNA-Inspired
812 Behavioral Modeling. *IEEE Transactions on Dependable and Secure Computing*, 15(4), pp.561–576.

- 813 62. Chen, C. et al., 2015. 6 million spam tweets: A large ground truth for timely Twitter spam detection. In *2015*
814 *IEEE International Conference on Communications (ICC)*. pp. 7065–7070.
- 815 63. Mateen, M. et al., 2017. A hybrid approach for spam detection for Twitter. In *2017 14th International Bhurban*
816 *Conference on Applied Sciences and Technology (IBCAST)*. pp. 466–471.
- 817 64. Sedhai, S. & Sun, A., 2018. Semi-Supervised Spam Detection in Twitter Stream. *IEEE Transactions on*
818 *Computational Social Systems*, 5(1), pp.169–175.
- 819 65. Gupta, S. et al., 2018. Collective Classification of Spam Campaigners on Twitter: A Hierarchical Meta-Path
820 Based Approach. In *Proceedings of the 2018 World Wide Web Conference*. International World Wide Web
821 Conferences Steering Committee, pp. 529–538.
- 822 66.

823