

# A Machine Learning-Based Approach For Detecting Malicious PyPI Packages

Haya Samaana<sup>1</sup>, Diego Elias Costa<sup>2</sup>, Emad Shihab<sup>2</sup>, Ahmad Abdellatif<sup>3</sup>

<sup>1</sup>An Najah National University, Nablus, Palestine; hayasam@najah.edu

<sup>2</sup>Concordia University, Montreal, Quebec, Canada; diego.costa@concordia.ca, emad.shihab@concordia.ca

<sup>3</sup>University of Calgary, Calgary, Alberta, Canada; ahmad.abdellatif@ucalgary.ca

## Abstract

**Background.** In modern software development, the use of external libraries and packages is increasingly prevalent, streamlining the software development process and enabling developers to deploy feature-rich systems with little coding. While this reliance on reusing code offers substantial benefits, it also introduces serious risks for deployed software in the form of malicious packages - harmful and vulnerable code disguised as useful libraries. **Aims.** Popular ecosystems, such as PyPI, receive thousands of new package contributions every week, and distinguishing safe contributions from harmful ones presents a significant challenge. There is a dire need for reliable methods to detect and address the presence of malicious packages in these environments. **Method.** To address these challenges, we propose a data-driven approach that uses machine learning and static analysis to examine the package's metadata, code, files, and textual characteristics to identify malicious packages. **Results.** In evaluations conducted within the PyPI ecosystem, we achieved an F1-measure of 0.94 for identifying malicious packages using a stacking ensemble classifier. **Conclusions.** This tool can be seamlessly integrated into package vetting pipelines and has the capability to flag entire packages, not just malicious function calls. This enhancement strengthens security measures and reduces the manual workload for developers and registry maintainers, thereby contributing to the overall integrity of the ecosystem.

## CCS Concepts

• **Security and privacy** → **Intrusion/anomaly detection and malware mitigation**;

## Keywords

Malicious packages, Supply chain security, Vocabulary, PyPI

## 1 Introduction

Code reuse drives technological innovation by enhancing developer productivity and enabling the creation of feature-rich, maintainable systems [10, 25, 36]. Package managers like npm and PyPI facilitate this innovation, with PyPI allowing developers to contribute freely to a vast repository of Python packages [29, 58]. PyPI meets diverse needs, including those in artificial intelligence, with millions of packages available [32]. Its popularity has made Python the most favored programming language as of April 2024, according to the TIOBE index.<sup>1</sup>

Developers are increasingly reusing code, which inadvertently heightens the risk of integrating malicious code into applications

[31]. Malicious actions can include inserting backdoors and stealing sensitive data. The PyPI registry faces various attacks, including compromised accounts, where attackers gain full control of a maintainer's account to publish malware [29, 59]. Typosquatting [1, 2] exploits typographical errors in package names, while combosquatting [52, 59] leverages the order of nouns in names. A notable example is the malicious package "jellyfish," which replaced a character in the benign "jellyfish" package name to steal SSH and GPG keys, remaining undetected for a year [58].

Several code-scanning methods have been proposed to identify malicious packages in popular ecosystems like NPM, PyPI, and Maven. These approaches encompass both traditional methods (e.g., anomaly detection [32], dynamic and static analysis [20, 21], and machine learning-based methods (e.g., unsupervised k-means clustering [23], supervised learning [48])). However, managers of large package registries, such as PyPI, struggle to handle the daily influx of packages, making it challenging for manual oversight [26, 64]. According to the libraries.io database [3], which is an open source repository and dependency database that catalogues libraries of the most popular ecosystems, and it has been used by previous work as a source of library metadata [5, 18], analysis reveals that over the course of one week, PyPI developers publish around 1,800 public package versions, including both new packages and updated versions of existing packages. To this end, prior works show that many existing scanning tools face limitations such as scalability issues [34, 48], a narrow focus on specific ecosystem aspects like package updates [23, 48], high resource costs [34], and high false alerts [8, 34]. In this context, we have detailed code scanners that work well if we analyse a few packages, however, there is a necessity for an approach that can extend its scalability to encompass the entire ecosystem. Our approach is designed to fill this gap, addressing issues of false positives and negatives. Its strength lies in the integration of newly crafted features operating at the package level, distinguishing it from most existing tools that operate at the function call level. Our approach considers the interactions of multiple factors across code, function, file levels, and package metadata, to provide a scalable classification of the likelihood of a malicious package.

To evaluate the effectiveness of our approach, we conduct a study to answer the following three research questions.

**RQ1: How accurate is our approach in classifying malicious packages?** We developed six machine learning classifiers—Random Forest, Decision Tree, Support Vector Machine, Multilayer Perceptron, Naive Bayes (Bernoulli version), and stacking—using eight features: 2 metadata-related, 2 file-related, 3 code-related, and 1 text feature. The classifiers were evaluated on a dataset of 5,193 benign and 138 malicious packages. Additionally, we tested the method's

<sup>1</sup><https://www.tiobe.com/tiobe-index/>

generalizability on an unseen dataset of 397 typical and 143 new malicious packages. The stacking ensemble classifier achieved an F1-score of 94.2% in predicting malicious packages, particularly when incorporating the text-related feature. Our findings indicate strong performance and generalizability of the approach.

**RQ2: What features are the best indicators of malicious packages?** The study investigates features that differentiate malicious packages from benign ones, identifying metadata-related features as the best indicators. Key tokens like `GETATTR`, `CONNECT`, `READ`, `OPEN` are crucial for malicious package identification. The research also examines the effects of keyword removal, stop word removal, and stemming on classifier performance. It concludes that stemming does not influence performance, while removing keywords and stop words has minimal impact.

**RQ3: Is our approach useful?** We assess the viability and efficiency of our approach by comparing it with two state of the art tools, namely `bandit` and `packj`. This evaluation is conducted on a random subset comprising 50 benign packages and an additional 50 samples of malicious packages obtained from an external dataset and is performed on two scenarios: the entire package and the `setup.py` file. The results suggest that our approach effectively detects a considerable number of malicious packages in real-world scenarios, demonstrating a low rate of false alerts when compared to the tested tools. As an illustration, our tool successfully detected 4 out of 5 recently released malicious packages, outperforming the tested tools that proved unsuccessful in this regard. Furthermore, by selecting features related to the Python ecosystem, our method can detect a broader spectrum of malicious Python packages, in contrast to the state-of-the-art approach [37].

Our study contributes to the research and practice on four fronts.

- Unlike existing strategies that focus on detecting suspicious function calls, our approach operates at the package level, applicable to the entire ecosystem.
- This is the first study to employ a vocabulary-based method to automatically detect malicious packages within the PyPI ecosystem.
- We introduce significant new features previously unexplored in predicting malicious functions and packages.
- We publicly share our dataset [7] for further research on enhancing the security of package managers against supply chain attacks, comprising 5,331 packages (138 malicious and 5,193 popular) with various features.

## 2 Background

Malicious packages are software components intentionally designed to harm systems, often containing malware or code for unauthorized activities like stealing sensitive information or installing backdoors [11, 50]. In contrast, benign packages are safe and intended to perform their designated functions without compromising security.

Numerous studies indicate that malware is continuously evolving, employing diverse techniques to evade detection tools [21, 61, 62]. Attackers target all stakeholders in the software supply chain, including end-users, developers, package maintainers, and registry maintainers. One prevalent tactic is typosquatting and combosquatting [52, 59], as many registries lack security policies [21].

In typosquatting, malicious packages are published with names similar to popular packages to deceive developers into downloading them. Combosquatting manipulates the order of words in package names, such as changing `'python-nmap'` to `'nmap-python'`. An illustrative incident occurred in May 2022 when attackers published a malicious package named `pymafka`, which mimicked the legitimate `PyKafka` package, leading to 325 downloads before its removal [43]. Additionally, attackers publish new malicious packages directly, often employing code obfuscation methods to conceal harmful code from analysis. Techniques like encoding and encryption are commonly used, as seen in the `colourama@0.1.6` typosquatting variant of `colorama`, which utilized `base64` encoding to evade detection as shown in listing 1 (Line 1). Malicious packages like `hipid` and `hpid` [22] have been reported using uncommon `base32` encoding, while the `botaa3` typosquatting package employs bitwise XOR encryption and `base64` encoding to obscure its malicious payloads [13].

Many tools are developed to detect malicious attacks against popular ecosystems (e.g., NPM and PyPI) such as `Malware-check` [35], `OSSGadget` [40], `Maloss` [34], `Packj` [41], `Bandit4mal` [9], and `Bandit` [8]. To the best of our knowledge, the proposed tools work in detecting risky factors in Python packages, always reporting alerts at the level of Python functions. To put our results into perspective, we opt to select both `Packj` and `Bandit` as baselines in our evaluation because 1) they have been selected as benchmarks in previous research [48, 56, 58] and 2) they are mature tools that provide detailed reports that facilitate our manual analysis. The **Bandit** tool [8] is a widely-used static analysis tool designed for identifying security vulnerabilities in Python files. It employs pre-defined rules and the Abstract Syntax Tree (AST) representation of source code to enhance its analysis. `Bandit` rates issues based on severity and confidence levels (low, medium, high), providing insights into potential impact and reliability [46]. It has also been used as a benchmark in previous research [55, 56]. On the other hand, the **Packj** tool [41] employs a comprehensive analysis approach with three steps: static code analysis, metadata analysis, and optional dynamic analysis. It examines package code for filesystem, network, and process API usage, validates metadata attributes, and can perform dynamic analysis. `Packj` is built upon the `MalOSS` project [34] and is recognized in research as a benchmarking framework [48, 56, 58]. Given that attackers employ diverse tactics to obfuscate the detection of malicious packages, and existing registries lack a robust review process for package publication [21], and most existing tools suffer from considerable limitations, our emphasis is on integrating various features. These include meta-related, code-related, and text-related features. These features aim to effectively capture the range of tactics employed by attackers, enhancing the ability to identify malicious packages.

### Listing 1: `colourama` package [17] uses code obfuscation in `setup.py` file to defeat analysis (snippet code).

```
def run(self):
    exec("b3MxID0WNv...cmIudA==" .decode('base64')) \\Line 1
    os = platform.system()
    req = urllib2.Request('https://grabify.link/...',
        texto = urllib2.urlopen( req ).read()
```

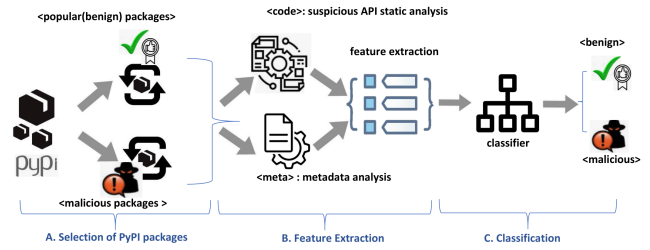
**Table 1: Existing techniques for analyzing PyPI suspicious/malicious packages.**

Tool	Input	Technique
Malware Checks [35]	setup.py file	static (Regular Expression)
OSSGadget [40]	package+artifacts	static (Regular Expression)
MalOSS [21]	package	hybrid (metadata, static, dynamic)
Packj [41]	package	hybrid (metadata, static, dynamic)
bandit4mal [9]	package	static (Abstract Syntax Tree)
bandit[8]	package	static (Abstract Syntax Tree)

### 3 Related work

**Traditional approaches.** Several methods have been proposed for identifying malicious packages. Liang et al. [32] used anomaly detection, combining abstract syntax tree (AST) and regular expression techniques, achieving a 97.51% reduction in review workload, though it struggles with small malicious artifacts. Duan et al. [20] reported 339 malicious packages through dynamic and static analysis but lacked false positive analysis. Their MALOSS framework, while comprehensive, requires significant resources. Vu et al. [55] identified differences between published packages and source repositories but missed packages without repositories. Ohm et al. [39] relied on dynamic analysis with heavy manual intervention, while Rieck et al. [45] monitored behavior in a sandbox but lacked detailed analysis. Various methods targeting typosquatting often suffer from false positives and negatives due to their focus on Levenshtein distance alone [49, 52, 59].

**Machine learning approaches.** Related work includes Garrett et al. [23], who used unsupervised learning to identify suspicious NPM package updates based on features like resource access and API usage, flagging 539 updates but not investigating malicious packages in depth. Sejfia and Schäfer [48] proposed an automated method for detecting malicious NPM packages that examined version changes and employed classifiers, successfully identifying 95 out of 96,287 unknown malicious packages with acceptable false positives. However, their approach lacks scalability and is limited to package updates, unable to address single-release packages. Halder et al. [27] developed MeMPtec, which utilizes features from package metadata. However, our approach goes deeper by analyzing metadata and file-related features. We parse the contents of license and configuration files, rather than simply reporting their existence. Recent research by Ohm et al. [37] focuses on detecting malicious NPM packages using a combination of classifiers, achieving true positive rates over 70% by evaluating 25,210 models. Their optimal combination involved Support Vector Machine, Multi Layer Perceptron, and Random Forest, successfully identifying 13 previously unknown malicious packages. This work is the most closely related to ours, as it also seeks to classify malicious packages within the NPM ecosystem. However, our methodology differs significantly. We adopt a more holistic approach by incorporating features related to code, licenses, file configurations, and author information. Unlike Ohm et al. [37], which treats suspicious APIs as boolean features, we analyze entire lines of code as text features. To our knowledge, no previous research has integrated linter outputs with metadata to classify software packages as malicious. Ohm et al. [37] developed a methodology tailored for the NPM ecosystem, specifically focusing on identifying malicious JavaScript packages with features relevant to that environment. In contrast, our work is explicitly designed for the Python ecosystem. We have crafted features that incorporate python-specific setup configuration files,



**Figure 1: The workflow for identifying malicious packages.**

allowing for a more accurate analysis of potential threats within Python packages. This distinction highlights the adaptability of security measures to different programming ecosystems.

**Python malware detection tools.** In the literature, a multitude of methods such as malware-check [35], OSSGadget [40], Maloss [34], Packj [41], Bandit4mal [9], and bandit [8] have been suggested for detecting suspicious package. The majority of these methods examine various facets of a package through the utilization of metadata, static analysis, or dynamic analysis as shown in Table 1. Analysis methods for metadata (such as package name and author information) involve scrutinizing these metadata elements to detect potentially problematic packages. One instance of this is the utilization of package names and their popularity to identify suspicious packages that might involve tactics like typosquatting or combosquatting.

However, existing tools struggle with unproductive time use, overlooked security issues, and high resource demands for dynamic analysis. Many, like OSSGadget and malware-check, rely on limited rule-based detection, potentially missing harmful code and hindering efficient detection of malicious actions.

Unlike previous research efforts that created comparative frameworks for different ecosystems or aimed at identifying suspicious packages at the function calls level, our research stands out by focusing on identifying malicious packages at the package level. Notably, it surpasses the scope of proposed works by evaluating six classifiers on a comprehensive set of newly crafted features and features inherited from existing literature, ensuring the successful detection of malicious packages throughout the entire ecosystem.

### 4 Study Setup

The main goal of our study is to identify malicious packages from a set of suspicious packages. Many techniques are proposed for this purpose as shown in Table 1, however, they do not work at the package level, and requiring maintainers to inspect multiple alerts before identifying malicious packages. Therefore, in order to attain our objective and address the limitations of current methods, we develop a machine learning (ML) approach. This approach is centered around combining specific static features, as demonstrated in the upcoming Section 4.2. We resort to use ML technique as it is a popular technique in the area of information security [21, 28, 37, 48].

For a comprehensive view, Figure 1 shows the full workflow of our approach, which consists of three phases: collecting benign and malicious PyPI packages (Section 4.1), feature extraction including metadata, file, code, and text related features (Section 4.2), and classification phase which details the machine learning classifiers and presents the evaluation process (Section 4.3).

**Table 2: An overview of feature set used to identify malicious packages in PyPI, including reused features from Ohm et al. [37].**

Dimensions	Features	Definition	Reused Features From Ohm et al. [37]
<b>Metadata-related</b>	Has invalid or no homepage	The package includes invalid or no homepage.	✗
	Has invalid or no author email	The package includes invalid or no author email.	✗
<b>Text-related</b>	Suspicious LOC	Lines of code including suspicious APIs.	✓ (Consider each API as a boolean feature)
<b>Code-related</b>	Has (post) install command	The package includes install script.	✓
	Has suspicious URL	The package includes suspicious URL or IP address.	✓
	Has long string	The package includes very long string (obfuscation).	✓
<b>File-related</b>	Has minimum setup configuration	The developer does not specify the package details in the setup.cfg file.	✗
	Has mismatch license	The package is missing license type uniformity in the three core positions.	✗

## 4.1 Collecting Training and Testing Datasets

**Training dataset.** One of the main challenges in building a classifier model to identify malicious packages is the quality and abundance of data. As thousands of packages are published weekly, it is well known that only a small fraction of malicious packages are immediately flagged as security threats; hence, we cannot assume that other recently published packages are benign (non-malicious). To solve this problem, we use popular packages that have been used by projects for years to compose the training/validation set of non-malicious packages, similarly as done in prior work [37, 49, 63]. Also, widely used and trusted popular packages are often copied by malicious packages (e.g., typosquatting), aiming to camouflage their activities and reach a larger number of users and developers [26, 38]. As such, our approach must be able to distinguish between popular packages and their malicious copies. It is important, however, to note that 1) we refrain from using any popularity metric as a feature of our classifier, our classifier should use only code and metadata-related features, and 2) we also test our model’s performance on a set of benign packages that are also not popular, to simulate a real-case scenario of using our approach.

To build our training and validation dataset, we collect **5,193 benign packages**. We first collect the top 5,000 most downloaded PyPI packages, from the PyPI registry [54] and select the top 5,000 most dependent upon packages from PyPI, as recorded in the libraries.io database [3]. We then merge both datasets, removing all duplicates (the vast majority), leading to a total of 5,193 popular Python packages in the PyPI ecosystem.

To build our set of malicious packages for model training, we use the dataset of 252 malicious packages collected by Ohm et al. [38]. As this dataset contains multiple versions of the same malicious packages, we only kept the latest version of malicious packages, similarly as done in previous work [32], to avoid overrepresenting a single package in our training set. We also remove (18) packages that were deemed not complete (e.g., package that included only the setup.py file or just the malicious payload). Finally, the malicious dataset contains (138) packages. This dataset of malicious packages spans 2015 to 2023, and the majority of attack vectors target install time rather than runtime. The dataset includes a variety of different injection techniques, such as TypoSquatting (52%) and Trojan Horse (27%), with varying infiltration objectives, e.g., data exfiltration (44%), droppers (15%), backdoors (8%). More information can be found in the study of Ohm et al. [38].

**Test dataset.** While using popular packages as a proxy for benign packages is a sensible choice, it is expected that distinguishing between popular packages and malicious packages is easier than

finding malicious packages in a batch of newly published Python packages [57]. Thus, we craft a test dataset that, by definition, does not contain any package (malicious or not) seen by the model during training but also better represents the average packages in the PyPI ecosystem. To achieve this, we 1) randomly select 397 packages from the PyPI registry, hereby named as **typical packages**. For the malicious packages, we incorporated a new collection of **143 malicious packages**, distinct from our training dataset selected from the dataset introduced by Ohm et al. [38]. During the implementation of the study, the dataset maintainer supplemented the collected training dataset with these 143 malicious packages, which we subsequently designated as the test dataset.

## 4.2 Feature Extraction

We need to collect features that capture different characteristics of malicious packages. We observed that majority of related work targeted NPM ecosystem, thus a set of features are adapted from NPM ecosystem [23, 48], called Code-related features in our case. Other features were chosen based on the grey literature [33] and expert knowledge of the differences between benign and malicious packages (in our case called File-related features). Moreover, a few of metadata-related features (e.g., README length, dependency analysis, author information, homepage, number of versions) are explored in the context of machine learning based solution [32, 48, 55]. In general, it’s important to note that not all metadata features hold significance; certain attributes might introduce noise and subsequently impact the model’s performance. Hence, we expand the feature set to include metadata-related features. Finally, a text-related feature was derived from various sensitive APIs and permissions included in the source code. These APIs have the following behaviours: produce new code during runtime ("getattr"), create forks or terminate operating system processes ("exit", "Thread", "system"), access obfuscated (hidden) code, modify system or environment variables ("clear"), access files and directories ("open"), establish connections with external networks ("open\_connection"), or readwrite user input ("input", "mkdtemp").

Our approach involves integrating eight different sets of features summarized in Table 2, to capitalize on the synergies that can arise from the interactions among these distinct features.

**Metadata-related features.** Every PyPI package has a set of metadata that contains different information such as homepage, project-url, authors and maintainer-email. For the computation of these features, we analyze the PKG-INFO file associated with each package. The validity of the email address is assessed by confirming its domain, while the homepage URL’s validity is verified by ensuring its security and association with a well-known domain. The

compilation of popular domains is derived from the list provided in [19]. We consider two features of this dimension in our approach. *Has invalid or no homepage (repository)*: Previous work showed that the differences in source code between build artifacts of a package and the respective source code repository are a strong indicator for its maliciousness [55, 63]. Our hypothesis suggests that this feature could aid in differentiating between malicious and benign packages, as supported by previous studies that emphasize the importance of having a homepage or/and source code repository for a given package. Moreover, based on a quantitative analysis of our dataset, we find that only 20% of malicious packages have a valid homepage or repository, while 73% of benign packages own a valid homepage. *Has invalid or no author email*: A recent investigations highlight that packages lacking a valid author email address, often a result of improper packaging guidelines, may raise suspicion and serve as an indicator of potential malicious intent [57, 65]. Among the inspected packages of the before mentioned dataset, we find about 87% of benign packages have a valid author’s email address in comparison to only 46% valid email address in malicious packages.

**Listing 2: Snapshot of a generated report from packj tool**

```
"performs a process operation": [{
  "filepath": "10Cent10-999.0.4\setup.py",
  "api_name": "spawn",
  "lineno": "17"
}]
```

**Text-related feature**: Most known attacks had malicious code injected into setup.py file [38, 58]. In this particular situation, our hypothesis revolves around the possibility that the malicious attacker might disperse the harmful payload among various files. To better understand and address this scenario, we extract lines of code corresponding to the suspicious APIs. To avoid reinventing the wheel, we rely on the generated static analysis report of the packj tool [41] to construct a portion of text feature. These APIs are linked to file paths and line numbers, as demonstrated in listing 2. To formulate a complete text feature, we enrich these suspicious lines of code by including the source code that originates from the setup.py file, in addition to any other .py source code that involve suspicious URLs. Moreover, based on a quantitative analysis of our dataset, we observe that setup.py file was the most targeted file by attackers (75%). Our observation is inline with prior works [9, 38, 55] showed that setup.py file is the most likely file to be manipulated by attackers.

**Code-related features**. Many characteristics have been extensively discussed in prior research studies like [23, 32, 37, 48], highlighting their significance in proficiently recognizing malicious software packages. These characteristics include installation command, suspicious URL, and long string. In this context, we leverage the previously discussed generated text feature to extract the following features.

*Has (post) install command*: Prior works [37, 48] show that the install command initiates an external operation, which is a common characteristic observed in malicious packages. Thus, we use regular expressions to search for the install command in the generated text feature.

*Has suspicious URL*: Different studies showed that attackers often inject their IP address or URL address in malicious code [32, 37]. We use regular expression to extract URLs and IP addresses for a particular package, then we record the suspicion of each URL/IP based on different criteria such if the URL is insecure, and if it belongs to unpopular domain. We rely on a list of approved URL domains provided by Amazon to verify the domain of a URL. This list contains the top 1 million most popular domains from Alexa, and if the domain of the URL is not present in the list we mark the URL as suspicious [19].

*Has long string*: The obfuscated code is often very long [30, 32]. Attackers commonly apply specialized encoding methods such as base64 to obscure harmful payloads. A string is categorized as lengthy when its length surpasses a specific threshold (40 characters) [14]. Consequently, we adopt this characteristic with string length larger than 40, following prior works [14, 32, 37]. Our dataset indicates that this attribute seldom appears in harmless packages.

**File-related features**. Research conducted earlier in the NPM ecosystem [47] demonstrated that during typosquatting and com-bosquatting attacks, the malicious package imitated a popular package by utilizing the README file. Therefore, we resort to examine files other than .py files, such as PKG-INFO and setup.cfg files. Subsequently, we formulate the hypothesis that these files might play a role in discerning between malicious and benign packages. From this perspective, we extract two distinct features.

*Has minimum setup configuration*: In the last few years, package distribution guidelines are constantly evolving [44, 53]. In Python development, both setup.cfg and setup.py are common for the purpose of packaging and distribution. The setup.cfg file includes the configuration of various aspects of package distribution, such as metadata details. This approach is considered cleaner, more contemporary, and modular, enabling efficient management of settings without cluttering the main script. We compute this feature by parsing the content of the setup.cfg file. We define the minimal configuration as equivalent to the default contents automatically generated by the utilized packaging tool. We hypothesize that attackers focus on the malicious payload rather than the package design. In some instances, attackers may leave the setup.cfg file untouched, adhering to a minimal configuration. Hence, we assume that the absence of a robust design is posited as a potential indicator of a malicious package, particularly when the setup.cfg contents resemble the minimum configuration. To support our assumption, we examine our dataset, discovering that 79% of the malicious packages exhibit inadequate design, while only 45% of benign packages deviate from the best practices.

*Has mismatch license*: A grey literature [33] highlighted that around 10% of PyPI packages lack clear usage licenses, posing a potential risk for malicious attacks. Python packages employ three methods to express licenses: license classifier, field, and file [33, 60]. We calculate this feature by examining the agreement among the types of licenses employed in the aforementioned three methods. Trusted packages adhere to proper license usage, as indicated by a recent study[15]. However, our dataset analysis reveals that approximately 66% of malicious packages exhibit license discrepancies, compared to just 1% in benign packages.

Until this point, we have prepared all the features for both malicious and benign packages. Next, we perform pairwise Pearson correlation [12] between the features to check whether two independent variables have a linear relationship. We found no correlation.

### 4.3 Classifiers and Performance Evaluation

To perform our predictions, we leverage six machine learning classifiers from Scikit-learn python library [42]: Random Forest, Support-Vector Machine, Decision Tree, Multilayer Perceptron, Naive Bayes (Bernoulli version), and Stacking. These classifiers have been used in prior works [37, 48], as well as other software engineering works [4, 24, 51]. To measure the performance of each classifier, we compute the precision, recall, and F1-score.

## 5 Case Study Results

### 5.1 RQ1: How accurate is our approach in classifying malicious packages?

**Motivation:** Most registries have little to no review process for publishing packages [21], which can be exploited by attackers to publish different types of malware to harm all downstream stakeholders. To maintain the ecosystem health from malicious actions, we need to identify the malicious packages among millions of published packages. Furthermore, different scanning methods have been developed to detect packages that raise suspicion rather than those that are undoubtedly malicious. In this RQ, we aim to evaluate the effectiveness of our approach in detecting malicious packages. Our approach helps registry maintainers to have a timely and accurate prediction of the malicious packages even before publishing them, and keeps the registry clean.

**Approach:** To answer this question, we evaluated the classifiers in predicting malicious packages in the two datasets discussed in Section 4.2. During training and validation, we performed a stratified 10-fold cross-validation, which trained six classifiers on 90% of the dataset and measured precision, recall, and F1-score on the remaining 10%. To ensure that the results are not skewed by a particular random initialization or split, this process is repeated ten times for each classifier, and then an average performance of these runs is computed to present the overall performance for that classifier. Moreover, to understand the sort of vocabulary used in distinguishing malicious packages from benign ones, we extracted all tokens from a text feature and converted them to lowercase. We removed stop words from the set of tokens, and we evaluated the impact of these choices on the performance of classifiers. Then we used a combination of standard NLP techniques, including keywords and stop-word removal to convert the text into tokens. Then we fed these tokens as inputs for the text classification algorithms with/without the feature set extracted from the metadata-related, code-related, and file-related features. We relied on **Scikit-learn** python library [42] for all classifiers, with the best configurations as shown in Table 4, and **Keras** python library [16] for text pre-processing (Tokenizer). Next, to understand the generalizability of our model, we tested our trained model on the test set discussed in Section 4.1, applying the same process shown in Figure 1 to extract the corresponding features. For this evaluation, we reported only results from our best model, the Stacking classifier.

**Fitness results:** Table 4 presents the performance of experimented classifiers in terms of precision, recall, and F1-score values. Generally, the classifiers achieve high performance (F1-score > 87%) in identifying malicious packages in all classifiers except the Naive Bayes (52%). Ensemble stacking classifier outperforms all other classifiers achieving the highest F1-score (94.2%) considering text vocabulary and metadata-related, file-related, and code-related features. Upon closer examination of the impact of using text vocabulary for the best performed classifier (i.e., stacking), the results show that the text vocabulary has a significant contribution to the prediction. When comparing the performance of stacking classifier only on metadata-related, file-related, and code-related as input features, we find that while the combination of most extracted features (i.e., metadata-related, file-related, and code-related features) achieve F1-score (84%), text vocabulary feature achieves higher F1-score (89%) as shown in Table 3. The evaluation assures our conjecture that text vocabulary could be used as a strong candidate feature to distinguish malicious package from benign one.

**Test results:** To test the generalizability of our trained model, we report the classification performance when inferring only the test set (the model was not trained further). The results shown in Table 5 display both the instances of false positives and false negatives generated by our approach when employed on the test set. We observe that our approach performs consistently well, reaching a F1-score of 90%, suggesting that our trained model is capable of identifying malicious packages even when regular (non-popular) Python packages are in the mix. We can observe from Table 5 that all typical packages are correctly classified as benign packages. On the other hand, when looking at the malicious packages, we observe that 105 out of 143 are classified correctly as malicious packages. However, the 38 malicious packages that were wrongly classified as benign were clones of popular packages, to inject their malicious payload. Upon manual analysis, we note some of these package’s malicious payloads seem insufficient for our model to distinguish them from their benign counterparts. Further exploration of other features may help reduce these false negative mistakes.

**The stacking ensemble classifier performs best in detecting malicious packages, with a fitting F1-score of 94.2% and a test F1-score of 90% on unseen data. The text vocabulary features contribute significantly to distinguishing malicious from benign packages.**

### 5.2 RQ2: What features are the best indicators of malicious packages?

**Motivation:** Given that stacking classifier achieves a high fitting and test performance in identifying the malicious packages, we want to better understand the most important features that contribute to the prediction. By knowing these important features, we can set the characteristics that distinguish malicious packages from benign ones.

**Approach:** We rely on the permutation feature importance technique [6] to find the most useful features in the stacking classifier. This technique randomly permutes the values of one feature while preserving the values of the remaining features. This process is

**Table 3: Performance of stacking classifier on a train dataset.**  
**Best configuration (bold):** num-folds = [10, 15] num-trees = [12,100, 200] num-words =[200, 500,1000, 2000]  
tokenizer modes = [binary, tfidf, count, freq] lower-states = [ False, **True**] (M/B) = (Malicious/Benign)

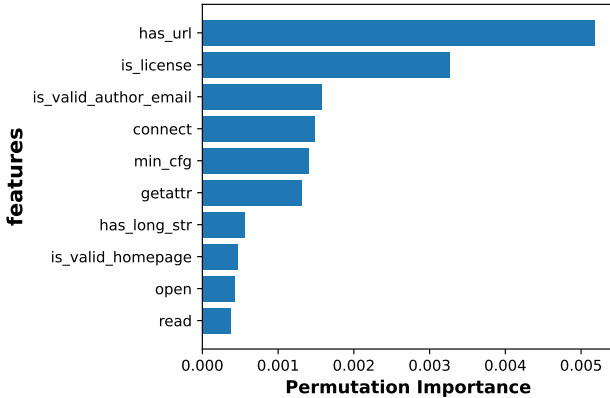
Input Features	Overall			Per-class (M/B)		
	Precision (%)	Recall (%)	F1-score (%)	Precision (%)	Recall (%)	F1-score (%)
Metadata-Code-File	94	77	84	(89/99)	(55/100)	(68/99)
Text	97	83	89	(96/99)	(67/100)	(79/100)
Text + preprocessing	98	83	89	(97/99)	(67/100)	(79/100)
Text + Metadata-Code-File	98	91	94	(96/100)	(82/100)	(88/100)
<b>Text + Metadata-Code-File + preprocessing</b>	<b>98</b>	<b>91</b>	<b>94</b>	<b>(97/100)</b>	<b>(81/100)</b>	<b>(89/100)</b>

**Table 4: Performance of different classifiers under precision (P), recall (R), and F1 score.**

Classifier	P (%)	R (%)	F1 (%)
Random Forest (RF) [n_estimators=12]	97	83	89
Support Vector Machine (SVM) [kernel= linear]	99	80	87
Decision Tree (DT) [max_depth=15]	89	89	89
Multilayer Perceptron (MLP) [activation= logistic, solver= adam]	97	90	93
Naive Bayes-Bernoulli (NB)	54	81	52
<b>Stacking (RF+SVM+MLP+DT+NB)</b>	<b>98</b>	<b>91</b>	<b>94</b>

**Table 5: Performance of our approach on the training and test dataset.**

Dataset	Package type	Total	FP	FN	Precision (%)	Recall (%)	F1-score (%)
Training	popular	5,193	3	-	98	91	94
	malicious	138	-	23			
Test	typical	397	0	-	96	87	90
	malicious	143	-	38			



**Figure 2: The permutation feature importance.**

applied to all features discussed in Section 4.2. We use permutation\_importance function in the Scikit-learn library [42] to compute the feature importance values of the stacking classifier.

**Results:** We find that the most important features are related to the following features: has suspicious url, has a licence, has a valid author email, has minimum configuration, has long string, and has a set of suspicious API (GETATTR, CONNECT, OPEN, AND READ), as shown in Figure 2. Upon a deeper examination of the percentage of the most important feature in malicious and benign packages in our dataset, we find that while 28% of malicious packages contain a suspicious url, only 9% of the benign packages include it. Another important feature is the "Has a License." Based on our analysis, the result shows that 77% of the malicious packages have a mismatch license between the three known locations (i.e., file, classifier, and field), and only 1% of the benign packages have this property. This

confirms that these features have a significant impact on the target variable, at least in the given model’s context. Moreover, it is evident that features such as "Has long string" and "Has suspicious URL" possess attributes indicative of malicious packages. Moreover, our dataset suggests that malicious packages tend to have more invalid homepage (80% compared to 27% for benign packages). The same observation for the "Has minimum configuration" (79% compared to 45% for benign packages). Note that only the suspicious APIs are not enough to capture most malicious packages (low recall) as shown in Table 4, but when combined with metadata features, the model reaches excellent performances.

**Code-related, file-related, and metadata-related features are all contribute to identify the malicious packages. Moreover, suspicious APIs such as getattr, connect, open, and read, contribute to distinguishing malicious from benign packages.**

### 5.3 RQ3: Is our approach useful?

**Motivation:** In RQ1, we assessed the performance of our approach and found that the stacking classifier achieves the best results in identifying malicious packages with F1 score of 94.2%. Moreover, a recent study [57] has brought attention to the fact that popular packages are different from a typical Python package. The study emphasizes that these popular packages demonstrate enhanced engineering and a stronger alignment with standard Python programming conventions. Consequently, using only popular packages as the benign dataset might lead to unrealistic benchmark results since these packages might be relatively easy for detection tools to classify as benign.

**Approach:** To put our results into perspective, we conducted two experiments. We compared our approach with (1) two benchmarking tools: Bandit and Packj, as in prior works [48, 56, 58], (2) the method of Ohm et al. [37]. For the first experiment, we utilized the above mentioned tools for the reasons specified in Section 2. We selected these tools, despite their differing threat models, to explore their effectiveness in detecting malicious dependencies and to measure the manual effort needed to filter through all signals they generate. We run both tools and manually examine the returned alerts to identify malicious packages, simulating how practitioners use these tools to identify suspicious functions and conclude the presence of malicious packages in their applications. We avoid the computational overhead associated with Packj tool by focusing solely on the static code analysis capability for identifying API usage, which is then processed further to isolate the relevant phantom

**Table 6: Number of generated alerts for a random sample of external dataset for packj and bandit tools.**

Malicious Package	# Packj Alerts		# Bandit Alerts	
	whole pkg	setup	whole pkg	setup
2022-requests-3.0.0	69	16	555	21
typing-unions-3.10.0.1	38	0	61	1
rumihelling-0.0.1	15	0	4	0
dldsord-1.0.3	0	0	0	0

lines of code that constitute the text-related feature. In this experiment, we determine the number of alerts generated by applying the bandit and packj tools on a random subset containing 50 benign packages and another 50 samples of malicious packages sourced from the test dataset (Table 5). This evaluation is conducted in two situations: the whole package and exclusively the setup.py file.

In the second experiment, we exploited the method proposed by Ohm et al. [37], due to its close similarity to ours. Their dataset, drawn from the same source as ours, comprised nearly equivalent sizes, with 150 malicious npm packages compared to our 138 Python packages. We leveraged their method, which relied on the intersection of outcomes from three classifiers (SVM, RF, and MLP). It is important to highlight, however, that the approach of Ohm et al [37] was tailored to the JavaScript ecosystem, while our approach accounted for the specificities of Python packages. Thus, in this experiment, we are also evaluating how ecosystem-tailored features contribute to classifying malicious Python packages.

**Result:** Figure 3 presents the average number of alerts generated by bandit and packj tools in two scenarios, whole package and only setup.py file. The presented findings include both true positives and false positives. The result of Figure 3 (a) reveals that both the bandit and packj tools generated a considerable number of false alerts when assessing the entire benign package. Nonetheless, by considering only the setup.py file, Figure 3 (b), the average alert count peaked at a maximum of two alerts. Clearly, the number of alerts, for both tools has risen when considering the whole package both of benign and malicious packages, Figure 3 (a) and (c). Moreover, it can be noted from the same figure that when focusing solely on the setup.py file, we see that malicious packages tend to produce a higher number of alerts, averaging around six alerts, Figure 3 (d), in contrast to benign packages (two alerts as in Figure 3 (b)). This reconfirms the earlier conclusions presented in references such as [38, 58], which identified that the setup.py file was the most targeted file by attackers.

To provide further details, Table 6 shows the number of alerts produced by the packj and bandit tools for a subset of a selected sample of malicious packages. We observed a significant volume of alerts being generated by both tools, especially when examining the complete package scenario. This necessitated significant manual effort from developers to verify the status of these packages. The bandit tool, in particular, generated 555 alerts, and a manual investigation revealed that a number of these alerts were wrongly triggered. Our approach misclassified 16 samples out of 50 (32%) where the total misclassification was 38 out of 143 samples (about 27%) as shown in Table 5. This highlights the capability of our approach to reduce the manual workload for registry maintainers. In evaluating the two tools for package security, it is evident that both tools have limitations in identifying suspicious packages, as shown in instances like typing-unions-3.10.0.1, rumihelling-0.0.1,

and dldsord-1.0.3 (Table 6), particularly when focusing solely on the setup.py file. Notably, our approach successfully classified all benign packages, in contrast to packj and bandit tools, which generated excessive false alerts for benign packages. Nonetheless, it’s worth emphasizing that our method results in a relatively small number of false negatives. In the second experiment, we compared our approach, with features tailored to the Python ecosystem, against the approach of Ohm et al. [37]. We observed that the method proposed by Ohm et al. [37] failed to identify certain malicious packages. In contrast, our approach proved to be more effective, successfully detecting a wider range of these packages, as depicted in Figure 4. When employing a stacking classifier, we managed to detect 115 out of 138 in the training dataset and 105 out of 143 in the testing dataset. However, relying on the intersection method of the three classifiers only enabled us to detect 73 out of 138 and 81 out of 143 for the training and testing datasets, respectively. This indicates that our approach, which is tailored to the Python ecosystem, is more effective in detecting malicious Python packages compared to the approach proposed by Ohm et al. [37].

**Our approach is reasonably precise and does not produce an overwhelming number of false negatives and false positives, making manual investigation extremely feasible, when apply to the entire ecosystem.**

## 6 Model Misclassifications

This section analyzes the misclassification of packages by our model in both training and external datasets to understand the underlying reasons. In the training dataset, 23 out of 138 malicious and 3 out of 5,193 benign packages were misclassified. For the external dataset, 38 out of 143 malicious packages were misclassified, with no benign misclassifications among 397.

**Misclassification in popular packages.** Our investigation found that less than 0.06% of popular packages were misclassified by our model, primarily due to poorly structured code and deviation from standard Python conventions. Misclassification occurred because of a lack of essential metadata, such as homepage, author email, or license, as seen in packages like ordereddict-1.1, cookies-2.2.1, and blob-0.16. For instance, cookies-2.2.1 had an unknown license and used suspicious APIs like 'getattr', 'setattr', 'open', and 'call', which can indicate malicious behavior. We merged vocabulary from the setup.py file with lines from suspicious files, highlighting how missing metadata and the presence of suspicious APIs lead to misclassification. Ultimately, misclassifying benign packages as malicious is considered less risky than the reverse scenario.

**Misclassification in malicious packages.** Our model falsely classified some packages due to many reasons. (1) The package incorporates a reference to an external harmful function. For example, libpesh-0.1 package has been identified as malicious because it includes a reference to a harmful function called "rn" which can be found in the file named ENTRY\_POINT.TXT as "EGGSECUTABLE = LIBARI.PR:RN." Due to the absence of the text code, which are critical feature that our model relies on, this incident has the potential to deceive the model. (2) Using a group of suspicious APIs that are



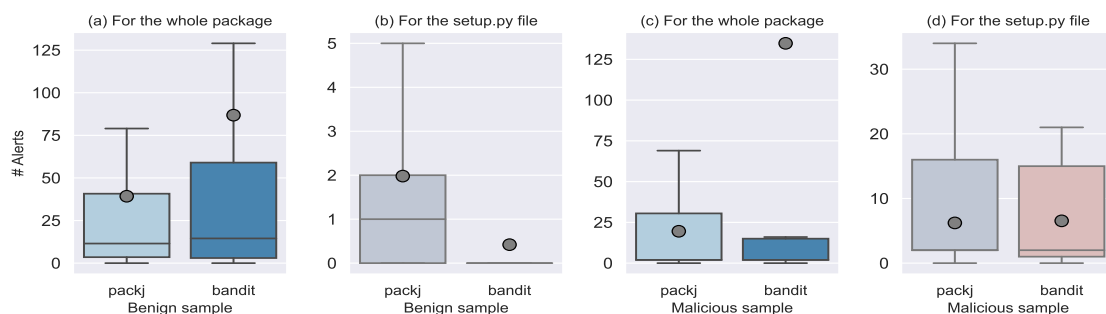


Figure 3: The average number of alerts generated by bandit and packj tools from the whole package and setup.py file.



Figure 4: The result of the intersection approach [37] of three classifiers on our train and test dataset.

commonly employed by benign packages. For example, the package `bzip-0.98` has been designed to appear as the legitimate `bz2file` package, and its installation script, `SETUP.PY` has been altered to contain a malicious code that is not particularly harmful.

## 7 Threats to validity

**Internal Validity.** The internal validity threats in this study include the potential for false positives and false negatives, although these are manageable. Manual review helps mitigate false positives. The dataset may also be biased due to clustering of similar malicious samples and assuming popular PyPI packages as benign. The selection criteria for popular packages could introduce variability, impacting the results. Despite these threats, we believe the dataset is sufficient for the experiments.

**Construct Validity.** One threat to construct validity is the reliance on the Packj tool [41] for static analysis, as we used it primarily to extract features for training our text classifier rather than drawing final conclusions. Another threat arises from the feature set construction, as we inherited some features from prior studies that effectively identify malicious packages [21, 23, 41, 48]. However, this reliance may limit our outcomes, highlighting the need for future research to investigate additional features and their impact. **External Validity.** External validity concerns the generalization of our findings. In our study, we assessed the performance of our approach using 5,193 benign packages and 138 malicious ones. Hence, our results may not generalize to other datasets, as malicious dataset may not accurately represent all malicious packages in the wild, and may there are malware with different characteristics than those in our dataset. However, we still believe that our dataset is comprehensive and serve the goal of this study. Moreover, to alleviate

this threat, we evaluate the model on an external (unseen) dataset, and we found that our approach performs very well, suggesting its generalizability.

## 8 Conclusion

We presented a machine learning-based method for detecting malware in PyPI packages using features from text, file, code, and metadata. Among the six classifiers evaluated, the stacking classifier performed best, while Naive Bayesian failed to detect known malicious packages. Our approach demonstrated practical effectiveness with low false negatives on external datasets. This method holds promise for automatic malware detection in PyPI. Future work includes expanding features and applying the technique to other ecosystems like NPM and RubyGems.

## References

- [1] (accessed: 23.07.2024). Bertusk. <https://bertusk.medium.com/discord-token-stealer-discovered-in-pypi-repository-e65ed9c3de06>
- [2] (accessed: 23.07.2024). dateutil. <https://snyk.io/blog/malicious-packages-found-to-be-typo-squatting-in-pypi/>
- [3] Accessed on 21/7/2024. Libraries- the open source discovery service. <https://libraries.io/>
- [4] Rabe Abdalkareem, Suhaib Mujahid, and Emad Shihab. 2020. A machine learning approach to improve the detection of ci skip commits. *IEEE Transactions on Software Engineering* (2020).
- [5] Mahmoud Alfadhel, Diego Elias Costa, and Emad Shihab. 2021. Empirical analysis of security vulnerabilities in python packages. In *2021 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 446–457.
- [6] André Altmann, Laura Tološi, Oliver Sander, and Thomas Lengauer. 2010. Permutation importance: a corrected feature importance measure. *Bioinformatics* 26, 10 (2010), 1340–1347.
- [7] Anonymous. 2024. A Machine Learning-Based Approach For Detecting Malicious PYPI Packages | Zenodo. <https://zenodo.org/records/13825064>.
- [8] bandit. (accessed: 27.07.2024). bandit. <https://github.com/PyCQA/bandit>
- [9] bandit4mal. (accessed: 20.07.2024). bandit4mal. <https://github.com/lyvd/bandit4mal>
- [10] Victor R Basili, Lionel C Briand, and Walcélio L Melo. 1996. How reuse influences productivity in object-oriented systems. *Commun. ACM* 39, 10 (1996), 104–116.
- [11] bleepingcomputer. (accessed: 12.08.2024). <https://www.bleepingcomputer.com/news/security/malicious-pypi-packages-hijack-dev-devices-to-mine-cryptocurrency/>
- [12] Kenneth A Bollen and Kenney H Barb. 1981. Pearson’s r and coarsely categorized measures. *American Sociological Review* (1981), 232–239.
- [13] botaa3. (accessed: 20.07.2024). botaa3. <https://blog.sonatype.com/another-day-of-malware-malicious-botaa3-pypi-package>
- [14] Davide Canali, Marco Cova, Giovanni Vigna, and Christopher Kruegel. 2011. Prophiler: a fast filter for the large-scale detection of malicious web pages. In *Proceedings of the 20th international conference on World wide web*. 197–206.
- [15] Bodin Chinthanet, Brittany Reid, Christoph Treude, Markus Wagner, Raula Gaikovina Kula, Takashi Ishio, and Kenichi Matsumoto. 2021. What makes a good Node. js package? Investigating Users, Contributors, and Runnability. *arXiv preprint arXiv:2106.12239* (2021).

- [16] François Chollet et al. 2018. Keras: The python deep learning library. *Astrophysics source code library* (2018), ascl-1806.
- [17] colourama. (accessed: 23.07.2024). colourama. <https://bertusk.medium.com/cryptocurrency-clipboard-hijacker-discovered-in-pypi-repository-b66b8a534a8>
- [18] Alexandre Decan, Tom Mens, and Eleni Constantinou. 2018. On the impact of security vulnerabilities in the npm package dependency network. In *Proceedings of the 15th international conference on mining software repositories*. 181–191.
- [19] domains. (accessed: 20.06.2024). domains. <http://s3.amazonaws.com/alexa-static/top-1m.csv.zip>
- [20] Ruian Duan, Omar Alrawi, Ranjita Pai Kasturi, Ryan Elder, Brendan Saltaformaggio, and Wenke Lee. 2020. Measuring and preventing supply chain attacks on package managers. *arXiv preprint arXiv:2002.01139* (2020), 18–52.
- [21] Ruian Duan, Omar Alrawi, Ranjita Pai Kasturi, Ryan Elder, Brendan Saltaformaggio, and Wenke Lee. 2020. Towards measuring supply chain attacks on package managers for interpreted languages. *arXiv preprint arXiv:2002.01139* (2020).
- [22] encode32. (accessed: 23.07.2024). encode32. <https://jfrog.com/blog/jfrog-discloses-3-remote-access-trojans-in-pypi/>
- [23] Kalil Garrett, Gabriel Ferreira, Limin Jia, Joshua Sunshine, and Christian Kästner. 2019. Detecting suspicious package updates. In *2019 IEEE/ACM 41st International Conference on Software Engineering: New Ideas and Emerging Results (ICSE-NIER)*. IEEE, 13–16.
- [24] Mehdi Golzadeh, Alexandre Decan, Damien Legay, and Tom Mens. 2021. A ground-truth dataset and classification model for detecting bots in GitHub issue and PR comments. *Journal of Systems and Software* 175 (2021), 110911.
- [25] Rebecca Elizabeth Grinter. 1996. *Understanding dependencies: A study of the coordination challenges in software development*. University of California, Irvine.
- [26] Yaocong Gu, Lingyun Ying, Yingyuan Pu, Xiao Hu, Huajun Chai, Ruimin Wang, Xing Gao, and Haixin Duan. 2023. Investigating package related security threats in software registries. In *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE, 1578–1595.
- [27] Sajal Halder, Michael Bewong, Arash Mahboubi, Yinhao Jiang, Md Rafiqul Islam, Md Zahid Islam, Ryan HL Ip, Muhammad Ejaz Ahmed, Gowri Sankar Ramachandran, and Muhammad Ali Babar. 2024. Malicious Package Detection using Metadata Information. In *Proceedings of the ACM on Web Conference 2024*. 1779–1789.
- [28] Yung-Tsung Hou, Yimeng Chang, Tsuhan Chen, Chi-Sung Lai, and Chia-Mei Chen. 2010. Malicious web content detection by machine learning. *expert systems with applications* 37, 1 (2010), 55–60.
- [29] Berkay Kaplan and Jingyu Qian. 2021. A survey on common threats in npm and PyPI registries. In *International Workshop on Deployable Machine Learning for Security Defense*. Springer, 132–156.
- [30] Byung-Ik Kim, Chae-Tae Im, and Hyun-Chul Jung. 2011. Suspicious malicious web site detection with strength analysis of a javascript obfuscation. *International Journal of Advanced Science and Technology* 26 (2011), 19–32.
- [31] Piergiorgio Ladisa, Henrik Plate, Matias Martinez, and Olivier Barais. 2023. Sok: Taxonomy of attacks on open-source software supply chains. In *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE, 1509–1526.
- [32] Genpei Liang, Xiangyu Zhou, Qingyu Wang, Yutong Du, and Cheng Huang. 2021. Malicious Packages Lurking in User-Friendly Python Package Index. In *2021 IEEE 20th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)*. IEEE, 606–613.
- [33] licenses. (accessed: 20.07.2024). licenses. <https://blog.inedo.com/python/python-package-licenses>
- [34] maloss. (accessed: 1.08.2024). maloss tool. <https://github.com/ossanitizer/maloss>
- [35] malwarecheck. (accessed: 20.07.2024). malwarecheck. <https://warehouse.pypa.io/development/malware-checks.html>
- [36] Parastoo Mohagheghi, Reidar Conradi, Ole M Killi, and Henrik Schwarz. 2004. An empirical study of software reuse vs. defect-density and stability. In *Proceedings. 26th International Conference on Software Engineering*. IEEE, 282–291.
- [37] Marc Ohm, Felix Boes, Christian Bungartz, and Michael Meier. 2022. On the Feasibility of Supervised Machine Learning for the Detection of Malicious Software Packages. In *Proceedings of the 17th International Conference on Availability, Reliability and Security*. 1–10.
- [38] Marc Ohm, Henrik Plate, Arnold Sykosch, and Michael Meier. 2020. Backstabber’s knife collection: A review of open source software supply chain attacks. In *International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment*. Springer, 23–43.
- [39] Marc Ohm, Arnold Sykosch, and Michael Meier. 2020. Towards detection of software supply chain attacks by forensic artifacts. In *Proceedings of the 15th international conference on availability, reliability and security*. 1–6.
- [40] OSSGadget. (accessed: 20.07.2024). OSSGadget. <https://github.com/microsoft/OSSGadget>
- [41] packj. (accessed: 20.07.2024). packj. <https://github.com/ossilate-inc/packj>
- [42] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research* 12 (2011), 2825–2830.
- [43] pymafka. (accessed: 20.08.2024). pymafka. <https://www.bleepingcomputer.com/news/security/malicious-pypi-package-opens-backdoors-on-windows-linux-and-macs/>
- [44] Kenneth Reitz and Tanya Schlusser. 2016. *The Hitchhiker’s guide to Python: best practices for development*. " O’Reilly Media, Inc."
- [45] Konrad Rieck, Philipp Trinius, Carsten Willems, and Thorsten Holz. 2011. Automatic analysis of malware behavior using machine learning. *Journal of computer security* 19, 4 (2011), 639–668.
- [46] Jukka Ruohonen, Kalle Hjerpe, and Kalle Rindell. 2021. A large-scale security-oriented static analysis of python packages in PyPI. In *2021 18th International Conference on Privacy, Security and Trust (PST)*. IEEE, 1–10.
- [47] Simone Scalco, Ranindya Paramitha, Duc-Ly Vu, and Fabio Massacci. 2022. On the feasibility of detecting injections in malicious npm packages. In *Proceedings of the 17th International Conference on Availability, Reliability and Security*. 1–8.
- [48] Adriana Sejfa and Max Schäfer. 2022. Practical Automated Detection of Malicious npm Packages. *arXiv preprint arXiv:2202.13953* (2022).
- [49] Matthew Taylor, Raturaj K Vaidya, Drew Davidson, Lorenzo De Carli, and Vaibhav Rastogi. 2020. Spellbound: Defending against package typosquatting. *arXiv preprint arXiv:2003.03471* (2020).
- [50] thehackernews. (accessed: 12.07.2024). stealing. <https://thehackernews.com/2022/08/10-credential-stealing-python-libraries.html>
- [51] Christoph Treude and Martin P Robillard. 2016. Augmenting API documentation with insights from stack overflow. In *2016 IEEE/ACM 38th International Conference on Software Engineering (ICSE)*. IEEE, 392–403.
- [52] Nikolai Philipp Tschacher. 2016. *Typosquatting in programming language package managers*. Ph.D. Dissertation. Universität Hamburg, Fachbereich Informatik.
- [53] Raturaj K Vaidya, Lorenzo De Carli, Drew Davidson, and Vaibhav Rastogi. 2019. Security issues in language-based software ecosystems. *arXiv preprint arXiv:1903.02613* (2019).
- [54] Hugo van Kemenade and Richard Si. 2022. hugovk/top-pypi-packages: Release 2022.08. <https://doi.org/10.5281/zenodo.6947954>
- [55] Duc-Ly Vu, Fabio Massacci, Ivan Pashchenko, Henrik Plate, and Antonino Sabetta. 2021. Lastpymile: identifying the discrepancy between sources and packages. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 780–792.
- [56] Duc-Ly Vu, Zachary Newman, and John Speed Meyers. 2022. A Benchmark Comparison of Python Malware Detection Approaches. *arXiv preprint arXiv:2209.13288* (2022).
- [57] Duc-Ly Vu, Zachary Newman, and John Speed Meyers. 2023. Bad Snakes: Understanding and Improving Python Package Index Malware Scanning. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 499–511.
- [58] Duc Ly Vu, Ivan Pashchenko, Fabio Massacci, Henrik Plate, and Antonino Sabetta. 2020. Towards using source code repositories to identify software supply chain attacks. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*. 2093–2095.
- [59] Duc-Ly Vu, Ivan Pashchenko, Fabio Massacci, Henrik Plate, and Antonino Sabetta. 2020. Typosquatting and combosquatting attacks on the python ecosystem. In *2020 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*. IEEE, 509–514.
- [60] Weiwei Xu, Hao He, Kai Gao, and Minghui Zhou. 2023. Understanding and Remediating Open-Source License Incompatibilities in the PyPI Ecosystem. *arXiv preprint arXiv:2308.05942* (2023).
- [61] Wei Xu, Fangfang Zhang, and Sencun Zhu. 2012. The power of obfuscation techniques in malicious JavaScript code: A measurement study. In *2012 7th International Conference on Malicious and Unwanted Software*. IEEE, 9–16.
- [62] Wei Xu, Fangfang Zhang, and Sencun Zhu. 2013. Jstill: mostly static detection of obfuscated malicious javascript code. In *Proceedings of the third ACM conference on Data and application security and privacy*. 117–128.
- [63] Nusrat Zahan, Thomas Zimmermann, Patrice Godefroid, Brendan Murphy, Chandra Maddila, and Laurie Williams. 2022. What are weak links in the NPM supply chain?. In *2022 IEEE/ACM 44th International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE, 331–340.
- [64] Junan Zhang, Kaifeng Huang, Bihuan Chen, Chong Wang, Zhenhao Tian, and Xin Peng. 2023. Malicious Package Detection in NPM and PyPI using a Single Model of Malicious Behavior Sequence. *arXiv preprint arXiv:2309.02637* (2023).
- [65] Markus Zimmermann, Cristian-Alexandru Staicu, Cam Tenny, and Michael Pradel. 2019. Small world with high risks: A study of security threats in the npm ecosystem. In *28th USENIX Security Symposium (USENIX Security 19)*. 995–1010.