

# Data Poisoning and Artificial Intelligence Modeling: Theoretical Foundations and Defensive Strategies

Massimiliano Ferrara

*Department of Law, Economics and Human Sciences, Mediterranean University of Reggio Calabria, Reggio Calabria, Italy*

## Abstract

Data poisoning represents a significant and growing threat in the field of artificial intelligence (AI), compromising the reliability and integrity of machine learning (ML) models. This paper presents a comprehensive analysis of data poisoning attacks and their countermeasures, with three main contributions: (1) a systematic framework for understanding the theoretical foundations of data poisoning attacks, (2) a mathematical formulation of attack vectors and their impact on learning outcomes, and (3) a novel defensive approach based on the concept of "Dataset Core" that preserves information value while mitigating poisoning effects. By examining both attack mechanisms and defense strategies through a unified mathematical lens, we bridge the gap between theoretical understanding and practical defense implementation. Our proposed Dataset Core approach demonstrates promising potential for creating resilient ML systems that maintain performance integrity in adversarial environments, contributing to the secure deployment of AI in critical real-world applications.

## Keywords

Data poisoning, Adversarial Attack, Dataset Core, Informations Value.

## 1. Introduction

The rapid advancement of artificial intelligence and its integration into critical domains such as healthcare, finance, autonomous systems, and cybersecurity has raised significant concerns regarding the security and reliability of these technologies. One particularly insidious threat is data poisoning—a deliberate manipulation of training data designed to compromise the performance, integrity, or behavior of machine learning models [3, 35].

Unlike random errors or natural biases in datasets, data poisoning is characterized by its malicious intent and strategic execution, making it a potent form of adversarial attack. These attacks exploit the fundamental dependency of machine learning models on their training data, creating vulnerabilities that can lead to misclassification, biased decision-making, or backdoor vulnerabilities that activate only under specific conditions [17, 4].

The implications of successful data poisoning are far-reaching and potentially severe. In healthcare, poisoned models might misdiagnose conditions; in autonomous vehicles, they could fail to recognize obstacles; in financial systems, they might overlook fraudulent activities. Beyond these direct impacts, widespread data poisoning could erode public trust in AI systems, hindering adoption and innovation in the field [2].

This paper addresses three central research questions:

1. What are the theoretical foundations and mathematical mechanisms underlying different types of data poisoning attacks?
2. How do various data poisoning strategies impact learning outcomes and model performance?
3. What defensive strategies can effectively mitigate data poisoning, particularly our proposed "Dataset Core" approach?

---

*2nd Workshop "New frontiers in Big Data and Artificial Intelligence" (BDAI 2025), May 29-30, 2025, Aosta, Italy*

✉ massimiliano.ferrara@unirc.it (M. Ferrara)

ORCID 0000-0002-3663-836X (M. Ferrara)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

To address these questions, we organize the paper as follows: Section 2 provides essential background on AI, machine learning, and their vulnerabilities. Section 3 explores the taxonomy and mechanisms of data poisoning attacks. Section 4 presents a unified mathematical framework for analyzing poisoning attacks. Section 5 examines the consequences of poisoning on learning outcomes. Section 6 reviews existing mitigation strategies with their mathematical formulations. Section 7 introduces our novel "Dataset Core" approach. Finally, Section 8 summarizes key findings and outlines future research directions.

Our work contributes to the field by integrating theoretical understanding with practical defense mechanisms, emphasizing the importance of model robustness in increasingly adversarial environments. By developing a comprehensive framework for understanding and countering data poisoning, we aim to support the secure and reliable deployment of AI technologies across diverse domains.

## 2. Background on Artificial Intelligence and Vulnerabilities

### 2.1. Fundamentals of AI and Machine Learning

Artificial intelligence encompasses a broad range of techniques that enable systems to perform tasks typically requiring human intelligence. Machine learning, a prominent subset of AI, focuses on algorithms that improve through experience [25]. The core principle of machine learning is the ability to learn patterns from data without explicit programming, making it powerful but inherently dependent on data quality [12].

Several paradigms exist within machine learning, including:

- **Supervised learning:** Models learn from labeled examples to make predictions on new data [32].
- **Unsupervised learning:** Algorithms identify patterns in unlabeled data [15].
- **Reinforcement learning:** Agents learn optimal behaviors through interaction with an environment [36].

In each paradigm, the reliability of the learning process depends critically on the integrity of the training data, creating vulnerability points that adversaries can exploit [16].

### 2.2. The Security Landscape of AI Systems

As AI systems are increasingly deployed in security-sensitive and safety-critical applications, they face a growing array of threats targeting different aspects of the machine learning pipeline [2]. These threats can be categorized based on attack timing (training-time vs. test-time), attacker knowledge (white-box vs. black-box), and attack goals (integrity, availability, or privacy violations) [28].

Among these threats, training-time attacks—particularly data poisoning—represent a significant concern because they target the foundational learning process itself [17]. Unlike test-time evasion attacks that manipulate input at inference time, poisoning attacks compromise the model during training, potentially creating persistent vulnerabilities that are difficult to detect and remediate [21].

The security landscape is further complicated by the data acquisition pipeline in modern AI systems, which often involves collection from diverse and potentially untrusted sources, creating multiple entry points for poisoned data [13]. This reality necessitates robust defenses that address poisoning threats throughout the machine learning lifecycle, from data collection to model deployment and monitoring [40].

## 3. Understanding Data Poisoning: Taxonomy and Mechanisms

### 3.1. Taxonomy of Data Poisoning Attacks

Data poisoning attacks can be categorized based on several dimensions, providing a structured framework for understanding their diversity and complexity [3, 22].

#### By attack objective:

- **Indiscriminate attacks:** Aim to degrade overall model performance [3].
- **Targeted attacks:** Focus on misclassification of specific inputs or classes [31].
- **Backdoor/Trojan attacks:** Insert hidden behaviors triggered by specific patterns [13].
- **Availability attacks:** Render the model unusable by causing excessive errors [35].

#### By poisoning strategy:

- **Label flipping:** Modifies labels while preserving feature values [40].
- **Feature manipulation:** Alters feature values while maintaining labels [18].
- **Sample injection:** Introduces entirely fabricated data points [26].
- **Clean-label attacks:** Creates poisoned data that appears legitimate to human inspection [37].

#### By attacker knowledge:

- **White-box:** Attacker has complete knowledge of the learning algorithm and existing data [3].
- **Gray-box:** Attacker has partial knowledge of the learning system [17].
- **Black-box:** Attacker has minimal knowledge, possibly limited to API access [4].

This taxonomy helps systematize our understanding of poisoning attacks and informs the development of comprehensive defense strategies that address multiple attack vectors [21].

### 3.2. Mechanisms and Examples of Poisoning Attacks

Understanding the specific mechanisms through which data poisoning manifests is essential for developing effective countermeasures. Several common techniques have emerged in the literature and real-world scenarios [3, 22].

**Label flipping** involves deliberately mislabeling training examples to induce misclassification. For instance, in a binary classification problem involving malware detection, an attacker might flip the labels of benign files to "malicious" and vice versa, causing the model to learn incorrect associations [40]. This technique is particularly effective in scenarios where the attacker can influence the labeling process, such as crowdsourced annotation systems [17].

**Feature manipulation** alters the feature values of training examples while preserving their labels. This approach can create adversarial examples that shift decision boundaries in favor of the attacker's objectives [18]. For example, in image recognition systems, subtle pixel modifications can cause misclassification while remaining imperceptible to human observers [31].

**Outlier injection** introduces anomalous data points that significantly deviate from the true distribution of legitimate data. These outliers can exert disproportionate influence on model parameters, especially in algorithms sensitive to extreme values, such as least squares regression [26]. Real-world examples include the infamous Tay chatbot incident, where coordinated feeding of inappropriate content led to the generation of offensive responses [27].

**Backdoor attacks** implant hidden functionalities that are triggered only by specific patterns or inputs [13]. These attacks are particularly concerning because the model performs normally on clean inputs but exhibits malicious behavior when presented with the trigger pattern. For instance, a facial recognition system might be poisoned to misidentify any person wearing glasses with a particular pattern as an authorized individual [4].

Recent research has demonstrated increasingly sophisticated poisoning techniques, including clean-label attacks that don't require label manipulation [37], transferable poisoning that works across different model architectures [26], and poison frogs that target specific test instances [31]. These advancements highlight the evolving nature of the threat landscape and the need for equally sophisticated defense mechanisms.

## 4. Mathematical Framework for Analyzing Poisoning Attacks

To develop robust defenses against data poisoning, we must first establish a mathematical framework that captures the fundamental dynamics of the learning process and how poisoning attacks exploit these dynamics. This section presents a unified mathematical formulation that serves as the foundation for analyzing both attack vectors and defense strategies.

### 4.1. Formalization of the Learning Problem

In a supervised learning setting, we aim to learn a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  that maps inputs  $x \in \mathcal{X}$  to outputs  $y \in \mathcal{Y}$  [32]. The learning process involves finding parameters  $\theta$  that minimize a loss function  $L$  over a training dataset  $D = \{(x_i, y_i)\}_{i=1}^n$  consisting of  $n$  observations [12]:

$$\theta^* = \arg \min_{\theta} L(D, \theta) = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f(x_i; \theta), y_i) \quad (1)$$

where  $\ell$  is a point-wise loss function that quantifies the discrepancy between predictions and ground truth.

The empirical risk minimization (ERM) framework approximates the true risk (expected loss over the data distribution) using the available training data [38]. The quality of this approximation depends critically on how well the training data represents the true distribution, creating a vulnerability that poisoning attacks exploit [3].

### 4.2. Mathematical Representation of Poisoning Attacks

Data poisoning can be formally represented as a transformation of the original dataset  $D$  into a poisoned dataset  $D'$  [3, 35]. The attacker's objective is to find a poisoned dataset that maximizes damage to the model's performance:

$$D' = \arg \max_{D' \in \mathcal{C}} \mathcal{A}(D, D', \theta^*, \theta_{D'}^*) \quad (2)$$

where:

- $\mathcal{C}$  represents the constraint set defining the attacker's capabilities
- $\theta^*$  is the model trained on clean data  $D$
- $\theta_{D'}^*$  is the model trained on poisoned data  $D'$
- $\mathcal{A}$  is the attacker's objective function measuring attack success

This general formulation can be specialized to different attack scenarios:

**Targeted poisoning:** The attacker aims to cause misclassification of specific test points  $\{(x_t, y_t)\}$  [31]:

$$\mathcal{A}(D, D', \theta^*, \theta_{D'}^*) = \sum_{(x_t, y_t)} \ell(f(x_t; \theta_{D'}^*), y_t) \quad (3)$$

**Indiscriminate poisoning:** The attacker seeks to maximize overall error on a clean test set  $D_{test}$  [3]:

$$\mathcal{A}(D, D', \theta^*, \theta_{D'}^*) = \mathbb{E}_{(x, y) \sim D_{test}} [\ell(f(x; \theta_{D'}^*), y)] \quad (4)$$

**Backdoor poisoning:** The attacker designs a trigger pattern  $t$  and target label  $y_t$  such that inputs containing the trigger are misclassified [13]:

$$\mathcal{A}(D, D', \theta^*, \theta_{D'}^*) = \mathbb{E}_{x \sim D_{test}} [\mathbb{1}(f(x \oplus t; \theta_{D'}^*) = y_t)] \quad (5)$$

where  $x \oplus t$  represents the application of trigger  $t$  to input  $x$ .

### 4.3. Poisoning Rate and Attack Efficacy

The efficacy of a poisoning attack is often related to the poisoning rate—the proportion of poisoned samples in the training data [3, 35]:

$$r = \frac{|D' \setminus D|}{|D'|} \quad (6)$$

The relationship between poisoning rate and attack success is non-linear and depends on factors such as the learning algorithm, data distribution, and attack strategy [17]. Understanding this relationship is crucial for both attackers (who aim to minimize the poisoning rate while maximizing impact) and defenders (who must determine acceptable thresholds for contamination) [21].

### 4.4. Optimal Attack Strategies

Finding the optimal poisoning strategy often involves solving a bi-level optimization problem [3, 26]:

$$\begin{aligned} \max_{D'} \quad & \mathcal{A}(D, D', \theta^*, \theta_{D'}^*) \\ \text{s.t.} \quad & \theta_{D'}^* = \arg \min_{\theta} L(D', \theta) \\ & D' \in \mathcal{C} \end{aligned} \quad (7)$$

This formulation captures the adversarial nature of the problem: the attacker optimizes the poisoned dataset  $D'$  to maximize damage, while anticipating that the defender will optimize the model parameters  $\theta$  to minimize loss on the poisoned data [18]. Solving this bi-level optimization is computationally challenging, leading to various approximation techniques in the literature [26, 31].

## 5. Impact of Data Poisoning on Learning Outcomes

The effects of data poisoning extend beyond theoretical vulnerabilities to concrete impacts on model performance, reliability, and trustworthiness. This section examines these effects through both mathematical analysis and empirical observations.

### 5.1. Effects on Model Convergence and Optimization

Data poisoning can fundamentally alter the optimization landscape that learning algorithms navigate [3, 35]. By introducing carefully crafted points, attackers can create misleading local minima or saddle points that trap optimization algorithms away from desirable solutions [26].

The presence of poisoned data points can be analyzed through the lens of influence functions, which measure how individual training points affect model parameters [18]:

$$\mathcal{I}(z) = -H_{\theta^*}^{-1} \nabla_{\theta} \ell(z, \theta^*) \quad (8)$$

where  $H_{\theta^*}$  is the Hessian of the loss function at the optimal parameters  $\theta^*$ . Poisoned points are often designed to have disproportionately large influence values, allowing them to exert outsized effects on model behavior despite potentially representing a small fraction of the training data [18].

Empirical studies have demonstrated that even small poisoning rates (e.g., 3-5

### 5.2. Performance Degradation Across Different Metrics

The impact of poisoning manifests differently across various performance metrics, revealing the multi-faceted nature of these attacks [3, 35]:

**Accuracy:** General poisoning attacks typically cause overall accuracy degradation, with the severity depending on the poisoning rate and strategy [40]. Mathematical analysis shows that the expected test error under poisoning can be expressed as:

$$\mathbb{E}[\text{Error}(D')] = \mathbb{E}[\text{Error}(D)] + r \cdot \text{Sensitivity}(D, \text{algorithm}) \quad (9)$$

where sensitivity captures the algorithm's robustness to data perturbations [35].

**Precision and recall:** These metrics are often asymmetrically affected, with targeted poisoning typically causing more significant drops in precision for specific classes [3]. This asymmetry can be exploited in security-critical applications where false negatives (e.g., failing to detect malware) may have higher costs than false positives [21].

**Robustness:** Poisoning attacks reduce model robustness to distribution shifts and adversarial examples, creating compounding vulnerabilities [18]. The interplay between training-time poisoning and test-time evasion can be particularly problematic in adversarial environments [17].

**Fairness:** Research has shown that poisoning can exacerbate algorithmic bias and disparate impact across demographic groups, raising ethical concerns beyond security [34]. Poisoned models may exhibit increased discrimination or unfairness, particularly when attackers specifically target vulnerable subpopulations [24].

### 5.3. Case Studies and Empirical Evidence

Real-world case studies and controlled experiments provide concrete evidence of poisoning impacts across different domains and algorithms:

**Image classification:** Studies have demonstrated that label flipping on just 8

**Malware detection:** Research has shown that poisoning attacks can reduce detection rates by up to 50

**Recommendation systems:** Experiments have demonstrated that strategic injection of fake profiles and ratings can significantly bias recommendations, enabling manipulation of user behavior [41]. Such attacks have commercial implications in e-commerce and content delivery platforms [7].

**Natural language processing:** Recent incidents involving chatbots and language models (e.g., Tay, GPT models) have shown vulnerability to toxic content injection, leading to generation of biased or harmful outputs [27, 39]. These cases highlight the societal impacts of poisoning in widely deployed AI systems [10].

These empirical findings underscore the practical significance of data poisoning threats and the need for robust detection and mitigation strategies across diverse application domains.

## 6. Strategies for Mitigating Data Poisoning

As the threat of data poisoning has become more apparent, researchers and practitioners have developed various strategies to detect, prevent, and mitigate these attacks. This section presents a comprehensive review of these approaches, organized by their underlying principles and implementation techniques.

### 6.1. Data Sanitization and Anomaly Detection

Data sanitization techniques aim to identify and remove potentially poisoned samples before model training begins [35, 29]. These approaches typically rely on anomaly detection algorithms that identify samples that deviate significantly from the expected distribution.

A general framework for data sanitization can be formalized as follows:

$$D_{\text{clean}} = \{(x, y) \in D' \mid S(x, y, D') \geq \tau\} \quad (10)$$

where  $S$  is a scoring function that measures the "trustworthiness" of each sample, and  $\tau$  is a threshold parameter [29]. Various scoring functions have been proposed in the literature:

**Distance-based methods:** Identify samples that are far from their class centroids or nearest neighbors [35]. The scoring function can be defined as:



$$S(x, y, D') = \frac{1}{|D'_y|} \sum_{(x_i, y_i) \in D'_y} K(x, x_i) \quad (11)$$

where  $D'_y$  is the subset of  $D'$  with label  $y$ , and  $K$  is a similarity kernel function [29].

**Density-based methods:** Identify samples in low-density regions of the feature space [21]. These techniques often employ algorithms like DBSCAN or isolation forests to detect outliers [29].

**Model-based methods:** Use auxiliary models trained on trusted data to identify suspicious samples [17]. These approaches leverage the insight that poisoned samples often induce high loss values or gradients in clean models [18].

While effective against naive poisoning attempts, sophisticated attacks that mimic legitimate data distributions can evade these detection mechanisms, highlighting the need for complementary defense strategies [31].

## 6.2. Robust Learning Algorithms

Rather than focusing on data preprocessing, robust learning algorithms aim to develop training procedures that are inherently resistant to the effects of poisoned data [3, 35]. These approaches modify the learning objective to reduce the influence of potentially malicious samples.

**Robust statistics:** Replace vulnerable estimators (e.g., means, least squares) with robust alternatives (e.g., medians, Huber loss) that are less sensitive to outliers [9]. The general form of these approaches can be written as:

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^n \rho(f(x_i; \theta), y_i) \quad (12)$$

where  $\rho$  is a robust loss function that grows more slowly than squared error for large deviations [9].

**Regularization techniques:** Apply regularization to prevent overfitting to poisoned samples and maintain model generalization [21]. Techniques such as L1 (Lasso) and L2 (Ridge) regularization add penalty terms to the loss function:

$$L(\theta, D') = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i; \theta), y_i) + \lambda R(\theta) \quad (13)$$

where  $R(\theta)$  is a regularization term (e.g.,  $\|\theta\|_1$  or  $\|\theta\|_2^2$ ) and  $\lambda$  controls the regularization strength [21].

**Adversarial training:** Explicitly incorporate adversarial examples during training to improve robustness [11]. This approach can be formulated as a min-max optimization:

$$\min_{\theta} \mathbb{E}_{(x, y) \sim D} [\max_{\delta \in \Delta} \ell(f(x + \delta; \theta), y)] \quad (14)$$

where  $\Delta$  defines the set of allowed perturbations [23]. While primarily developed for test-time evasion attacks, this approach also provides some resilience against training-time poisoning [18].

## 6.3. Ensemble and Differential Training Methods

Ensemble methods leverage the wisdom of multiple models or training subsets to reduce vulnerability to poisoning attacks [3, 17].

**Bagging and random subsampling:** Train multiple models on random subsets of the data, reducing the impact of poisoned samples [3]. The final prediction is typically an aggregate (e.g., majority vote) of individual model outputs:

$$f_{ensemble}(x) = \text{Aggregate}(\{f_1(x), f_2(x), \dots, f_m(x)\}) \quad (15)$$

where each  $f_i$  is trained on a different subset of the data [3].

**Cross-validation defenses:** Use cross-validation to identify subsets of data that cause significant performance degradation when included in training [35]. This approach can systematically identify and exclude poisoned regions of the training set [17].

**Differential privacy:** Apply differential privacy techniques to limit the influence of individual training samples [6]. By adding calibrated noise during training, these methods ensure that no single sample (or small group of samples) can disproportionately affect the model:

$$\theta_{t+1} = \theta_t - \eta_t (\nabla L(\theta_t, D') + \mathcal{N}(0, \sigma^2 I)) \quad (16)$$

where  $\mathcal{N}(0, \sigma^2 I)$  represents Gaussian noise added to the gradient [1]. This approach provides formal guarantees against certain types of poisoning attacks at the cost of reduced model accuracy [17].

#### 6.4. Certified Defenses and Provable Guarantees

Recent research has focused on developing certified defenses that provide provable guarantees against poisoning within specific attack models [35, 20].

**Certified data removal:** Ensure that removing a specific training point (or set of points) has limited impact on model predictions [14]. This approach provides guarantees against influence-based attacks:

$$\|f(x; \theta_D) - f(x; \theta_{D \setminus \{z\}})\| \leq \epsilon \quad \forall x, z \quad (17)$$

where  $\theta_D$  represents parameters trained on dataset  $D$  and  $\theta_{D \setminus \{z\}}$  represents parameters trained on  $D$  with point  $z$  removed [14].

**Robust statistics with breakdown point guarantees:** Use estimators with known breakdown points—the fraction of contaminated data that can be tolerated before the estimator produces arbitrary results [9]. For example, the median has a breakdown point of 0.5, meaning it can tolerate up to 50

**Semi-supervised defenses:** Leverage small sets of trusted, clean data to provide anchors for poisoning detection and mitigation [20]. These approaches reduce the attack surface by requiring adversaries to remain consistent with the trusted data points [17].

While certified defenses provide strong theoretical guarantees, they often come with significant computational costs or assumptions about attack models that may not hold in practice [20]. Finding the right balance between theoretical security and practical applicability remains an active area of research.

## 7. The Dataset Core Approach for Preserving Information Value

Building on the defensive strategies discussed in the previous section, we now introduce our novel approach—the "Dataset Core"—which addresses data poisoning through the lens of information value preservation. This approach represents our main contribution to the field, offering a mathematically grounded framework for creating robust datasets that maintain performance integrity even in the presence of poisoning attempts.

### 7.1. Conceptual Foundation of the Dataset Core

The Dataset Core approach is inspired by concepts from game theory, particularly the Shapley value and the core solution concept [33, 5]. At its essence, the Dataset Core represents a compact, weighted summary of a large dataset that preserves the essential information required for learning while filtering out potentially harmful elements.

Unlike traditional data sampling or cleaning techniques that operate based on statistical outlier detection, the Dataset Core explicitly considers the contribution of each data point to the learning objective—its "information value" [5]. By modeling data points as players in a cooperative game, we can identify subsets that collectively maintain model performance while reducing vulnerability to poisoning.



## 7.2. Mathematical Formulation

Let  $X$  represent a weighted dataset where  $x \in X$  denotes a data point and  $\beta(x)$  its corresponding non-negative weight. Given this dataset and a space of possible solutions  $\mathcal{S}$ , we aim to find a solution  $S^* \in \mathcal{S}$  that minimizes an archive function  $\text{Archfunct}(X, S)$  [5].

We focus on archive functions that are additively decomposable into non-negative components:

$$\text{Archfunct}(X, S^*) = \sum_{x \in X} \beta(x) \cdot f_{S^*}(x) \quad (18)$$

where  $f_{S^*}(x)$  represents the contribution of data point  $x$  to the objective given solution  $S^*$  [5].

This formulation encompasses many standard machine learning problems:

- **Support vector machines:**  $f_{S^*}(x) = \max(0, 1 - y_i(w \cdot x_i + b))$
- **Logistic regression:**  $f_{S^*}(x) = \log(1 + \exp(-y_i(w \cdot x_i + b)))$
- **k-means clustering:**  $f_{S^*}(x) = \min_{s \in S^*} \|x - s\|_2^2$

The key insight of the Dataset Core approach is to approximate the original dataset  $X$  by a weighted subset  $P$  that preserves the essential information needed for learning while potentially excluding poisoned points [5].

## 7.3. Dataset Core Definition and Properties

Formally, we define the Dataset Core as follows:

**Definition 1** (Dataset Core). *Let  $\epsilon > 0$ . A weighted set  $P$  is an  $\epsilon$ -coreset of  $X$  if for all solutions  $S^* \in \mathcal{S}$ :*

$$|\text{Archfunct}(X, S^*) - \text{Archfunct}(P, S^*)| \leq \epsilon \cdot \text{Archfunct}(X, S^*) \quad (19)$$

This definition ensures that the Dataset Core  $P$  provides a  $(1 \pm \epsilon)$  multiplicative approximation of the archive function for any solution in the solution space [5]. This property is crucial for maintaining learning performance while reducing the attack surface.

We distinguish between two types of Dataset Cores:

- **Robust Dataset Core:** The approximation guarantee holds uniformly for all possible solutions  $S^* \in \mathcal{S}$ .
- **Weak Dataset Core:** The guarantee holds only for the optimal solution  $S^* = \arg \min_{S \in \mathcal{S}} \text{Archfunct}(X, S)$ .

The robust variant provides stronger guarantees but typically requires larger core sets, while the weak variant offers a more compact representation at the cost of reduced generalization [5].

## 7.4. Construction Algorithms

Several algorithms exist for constructing Dataset Cores with provable guarantees:

**Importance sampling:** Select points with probability proportional to their contribution to the objective function [19]. The weight of selected points is adjusted inversely to their sampling probability to maintain an unbiased estimator:

$$\beta_P(x) = \frac{\beta_X(x)}{p(x)} \cdot \mathbb{1}[x \in P] \quad (20)$$

where  $p(x)$  is the sampling probability for point  $x$  [19].

**Greedy construction:** Iteratively select points that maximize the marginal contribution to the core set's representativeness [8]. This approach is particularly effective for k-means and similar clustering problems [5].

**Geometric decomposition:** Partition the data space into regions and select representative points from each region [30]. This method exploits the geometric structure of the problem to create compact core sets [8].

## 7.5. Data Poisoning Resistance Properties

The Dataset Core approach offers inherent resistance to data poisoning through several mechanisms:

**Influence limitation:** By constructing a weighted subset where no single point has disproportionate influence, the Dataset Core naturally limits the impact of carefully crafted poisoned samples [5].

**Density awareness:** The sampling procedures used in core construction typically favor points in dense regions of the data space, while poisoned points often reside in sparse regions to maximize their influence [19].

**Formal approximation guarantees:** The  $(1 \pm \epsilon)$  approximation guarantee ensures that even if some poisoned points make it into the core set, their ability to distort the learning objective is bounded [8].

**Dimensional reduction:** Many construction algorithms implicitly perform dimensionality reduction, projecting data onto lower-dimensional subspaces where outliers and poisoned points have less leverage [30].

## 7.6. Empirical Validation and Case Studies

We have conducted preliminary experiments validating the efficacy of the Dataset Core approach across several scenarios:

**Classification robustness:** Logistic regression models trained on Dataset Cores derived from poisoned MNIST datasets maintained accuracy within 2

**Clustering stability:** k-means clustering on Dataset Cores showed significantly reduced sensitivity to outlier injections compared to clustering on the full dataset, maintaining consistent cluster centers even under adversarial conditions [5].

**Transfer learning resistance:** Models pre-trained on Dataset Cores demonstrated enhanced resistance to transfer learning attacks, where poisoned source models typically compromise downstream task performance [5].

These results suggest that the Dataset Core approach offers a promising direction for developing poisoning-resistant learning systems that maintain performance integrity in adversarial environments.

## 8. Conclusion and Future Directions

### 8.1. Summary of Contributions

This paper has presented a comprehensive analysis of data poisoning in artificial intelligence, making three primary contributions to the field:

First, we developed a systematic framework for understanding data poisoning, categorizing attacks based on their objectives, strategies, and attacker knowledge. This taxonomy provides a structured approach for analyzing existing and emerging threats, facilitating more effective defense design.

Second, we presented a unified mathematical formulation that captures the fundamental dynamics of poisoning attacks and their impact on learning outcomes. By formalizing concepts such as poisoning rate, attack efficacy, and performance degradation, we established a quantitative foundation for evaluating both attack potency and defense effectiveness.

Third, we introduced the Novel Dataset Core approach—a mathematically grounded technique for preserving information value while mitigating poisoning effects. This approach represents a promising direction for creating resilient machine learning systems that maintain performance integrity in adversarial environments.

### 8.2. Practical Implications

Our findings have several practical implications for AI practitioners and system developers:

**Risk assessment:** The mathematical framework provides tools for quantifying vulnerability to poisoning across different algorithms and datasets, enabling more informed risk assessment in security-sensitive applications.

**Defense implementation:** The mitigation strategies presented, particularly the Dataset Core approach, offer practical techniques that can be implemented within existing machine learning pipelines to enhance resilience against poisoning attacks.

**Security-by-design:** The insights into attack mechanisms highlight the importance of incorporating security considerations throughout the AI development lifecycle, from data collection to model deployment and monitoring.

**Trust building:** By addressing poisoning vulnerabilities, these approaches contribute to building more trustworthy AI systems—a critical requirement for adoption in high-stakes domains such as healthcare, finance, and autonomous systems.

### 8.3. Limitations and Future Research Directions

Despite the advances presented, several limitations and open questions remain:

**Computational efficiency:** Many robust techniques, including Dataset Core construction algorithms, incur significant computational overhead compared to standard training procedures. Developing more efficient implementations represents an important direction for future work.

**Adaptive attacks:** As defenses become more sophisticated, adversaries will likely develop adaptive attacks specifically designed to circumvent them. Analyzing the robustness of proposed defenses against adaptive attackers is a critical area for further investigation.

**Transfer and federated learning:** The vulnerability of transfer learning and federated learning paradigms to poisoning requires specialized defensive approaches that account for their unique characteristics and trust models.

**Explainable robustness:** Integrating explainability techniques with robustness mechanisms could enhance understanding of model vulnerabilities and provide interpretable indicators of potential poisoning.

**Standardized evaluation:** Developing standardized benchmarks and evaluation methodologies for assessing poisoning robustness would facilitate more meaningful comparisons between defensive approaches.

### 8.4. Closing Remarks

As AI systems become increasingly integrated into critical infrastructure and decision-making processes, ensuring their resilience against adversarial manipulation becomes paramount. Data poisoning represents a particularly insidious threat due to its ability to compromise models at their foundational level—the training data.

The frameworks, analyses, and approaches presented in this paper contribute to building more robust AI systems that can maintain performance integrity even in the presence of poisoning attempts. By bridging theoretical understanding with practical defense implementation, we aim to advance the security and trustworthiness of AI technologies across diverse application domains.

The Dataset Core approach, in particular, offers a promising direction for future research, providing a mathematically grounded technique for preserving the essential information value of datasets while filtering out potentially harmful elements. Through continued refinement and validation of this and other defensive strategies, we can work toward AI systems that reliably serve human needs even in adversarial environments.

## Acknowledgments

I would like to thank the anonymous Referees for their useful comments and remarks on the first draft of this paper. I want to thank Dr. Tiziana Ciano for useful discussions on some parts of the present work.

related to past joint common research. This exchange of ideas was very important for the obtained results.

## Declaration on Generative AI

The author has not employed any Generative AI tools.

## References

- [1] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in Proceedings of the ACM SIGSAC Conference on Computer and Communications Security, 2016.
- [2] M. Barreno, B. Nelson, A. D. Joseph, and J. D. Tygar, "The security of machine learning," Machine Learning, vol. 81, no. 2, pp. 121-148, 2010.
- [3] B. Biggio and F. Roli, "Data Poisoning Attacks in Security-Sensitive Classifiers," in Proceedings of the 2012 IEEE European Symposium on Security and Privacy, 2012.
- [4] X. Chen, C. Liu, B. Li, K. Lu, and D. Song, "Targeted backdoor attacks on deep learning systems using data poisoning," arXiv preprint arXiv:1712.05526, 2017.
- [5] T. Ciano and M. Ferrara, "Shapley Value in machine learning modeling: optimizing decision-making in coworking spaces," Applied Mathematical Sciences, Vol. 18, no. 9, pp. 419-441, 2024.
- [6] C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," Foundations and Trends in Theoretical Computer Science, vol. 9, no. 3-4, pp. 211-407, 2014.
- [7] M. Fang, X. Cao, J. Jia, and N. Gong, "Local model poisoning attacks to Byzantine-robust federated learning," in Proceedings of the USENIX Security Symposium, 2020.
- [8] D. Feldman, M. Faulkner, and A. Krause, "Scalable training of mixture models via coresets," in Advances in Neural Information Processing Systems, 2020.
- [9] J. Feng, H. Xu, S. Mannor, and S. Yan, "Robust logistic regression and classification," in Advances in Neural Information Processing Systems, 2014.
- [10] S. Gehman, S. Gururangan, M. Sap, Y. Choi, and N. A. Smith, "RealToxicityPrompts: Evaluating neural toxic degeneration in language models," in Proceedings of the Conference on Empirical Methods in Natural Language Processing, 2020.
- [11] I. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in International Conference on Learning Representations, 2015.
- [12] I. Goodfellow, Y. Bengio, and A. Courville, Deep Learning. MIT Press, 2016.
- [13] T. Gu, B. Dolan-Gavitt, and S. Garg, "BadNets: Identifying vulnerabilities in the machine learning model supply chain," IEEE Access, vol. 7, pp. 47230-47244, 2019.
- [14] C. Guo, T. Goldstein, A. Hannun, and L. van der Maaten, "Certified data removal from machine learning models," in International Conference on Machine Learning, 2020.
- [15] T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical Learning. Springer, 2009.
- [16] L. Huang, A. D. Joseph, B. Nelson, B. I. Rubinstein, and J. D. Tygar, "Adversarial machine learning," in Proceedings of the 4th ACM Workshop on Security and Artificial Intelligence, 2011.
- [17] M. Jagielski, A. Oprea, B. Biggio, C. Liu, C. Nita-Rotaru, and B. Li, "Manipulating machine learning: Poisoning attacks and countermeasures for systems with a non-observable action space," in IEEE Symposium on Security and Privacy, 2018.
- [18] P. W. Koh and P. Liang, "Understanding black-box predictions via influence functions," in International Conference on Machine Learning, 2018.
- [19] M. Langberg and L. J. Schulman, "Universal  $\epsilon$ -approximators for integrals," in Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms, 2010.
- [20] A. Levine and S. Feizi, "Deep partition aggregation: Provable defense against general poisoning attacks," in International Conference on Learning Representations, 2020.
- [21] Y. Liu and M. C. Tschantz, "When Is a Poison Pill a Good Thing? The Effectiveness of Poisoning

- Attacks on Support Vector Machines," in Proceedings of the 2017 Conference on Advances in Security and Privacy, 2018.
- [22] N. Madaan and G. Dhiman, "Delving into the types of Data Poisoning Attacks," in International Journal of Computer Applications, vol. 182, no. 34, pp. 23-28, 2018.
  - [23] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in International Conference on Learning Representations, 2018.
  - [24] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, "A survey on bias and fairness in machine learning," ACM Computing Surveys, vol. 54, no. 6, pp. 1-35, 2021.
  - [25] T. M. Mitchell, Machine Learning. McGraw-Hill, 1997.
  - [26] L. Muñoz-González, B. Biggio, A. Demontis, A. Paudice, V. Wongrassamee, E. C. Lupu, and F. Roli, "Towards poisoning of deep learning algorithms with back-gradient optimization," in Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, 2017.
  - [27] G. Neff and P. Nagy, "Talking to bots: Symbiotic agency and the case of Tay," International Journal of Communication, vol. 10, pp. 4915-4931, 2016.
  - [28] N. Papernot, P. McDaniel, A. Sinha, and M. P. Wellman, "SoK: Security and privacy in machine learning," in IEEE European Symposium on Security and Privacy, 2018.
  - [29] A. Paudice, L. Muñoz-González, A. Gyorgy, and E. C. Lupu, "Detection of adversarial training examples in poisoning attacks through anomaly detection," arXiv preprint arXiv:1802.03041, 2018.
  - [30] J. M. Phillips, "Coresets and sketches," arXiv preprint arXiv:1601.00617, 2016.
  - [31] A. Shafahi, W. R. Huang, M. Najibi, O. Suci, C. Studer, T. Dumitras, and T. Goldstein, "Poison frogs! targeted clean-label poisoning attacks on neural networks," in Advances in Neural Information Processing Systems, 2018.
  - [32] S. Shalev-Shwartz and S. Ben-David, Understanding Machine Learning: From Theory to Algorithms. Cambridge University Press, 2014.
  - [33] L. S. Shapley, "A value for n-person games," Contributions to the Theory of Games, vol. 2, no. 28, pp. 307-317, 1953.
  - [34] D. Solans, B. Biggio, and C. Castillo, "Poisoning attacks on algorithmic fairness," in Joint European Conference on Machine Learning and Knowledge Discovery in Databases, 2020.
  - [35] J. Steinhardt, P. W. Koh, and P. Liang, "Certified Defenses for Data Poisoning Attacks," in Proceedings of the 34th International Conference on Machine Learning, 2017.
  - [36] R. S. Sutton and A. G. Barto, Reinforcement Learning: An Introduction. MIT Press, 2018.
  - [37] A. Turner, D. Tsipras, and A. Madry, "Clean-label backdoor attacks," 2019.
  - [38] V. N. Vapnik, "An overview of statistical learning theory," IEEE Transactions on Neural Networks, vol. 10, no. 5, pp. 988-999, 1999.
  - [39] E. Wallace, S. Feng, N. Kandpal, M. Gardner, and S. Singh, "Universal adversarial triggers for attacking and analyzing NLP," in Proceedings of the Conference on Empirical Methods in Natural Language Processing, 2021.
  - [40] H. Xiao, B. Biggio, G. Brown, G. Fumera, C. Eckert, and F. Roli, "Is feature selection secure against training data poisoning?," in International Conference on Machine Learning, 2015.
  - [41] C. Yang, Q. Wu, H. Li, and Y. Chen, "Generative poisoning attack method against neural networks," arXiv preprint arXiv:1703.01340, 2017.