*Mathematics*

*Research article*

# Modeling by topological data analysis and game theory for analyzing data poisoning phenomena

**Massimiliano Ferrara**[*]

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

[*] **Correspondence:** Email: massimiliano.ferrara@unirc.it.

**Abstract:** This paper explored the theoretical and practical connections between topological data analysis (TDA), game theory, and data poisoning attacks. We demonstrated how the topological structure of data can influence strategic interactions in adversarial settings, and how game-theoretic frameworks can model the interplay between defenders and attackers in machine learning systems. We introduced novel formulations that bridge these fields, developing metrics to quantify the topological vulnerability of data structures to poisoning attacks. Our analysis reveals that persistence diagrams from TDA can serve as powerful tools for both detecting poisoning attempts and designing robust defense mechanisms. We proposed a Nash equilibrium-based approach to determine optimal poisoning and defense strategies, supported by mathematical formulations and theoretical guarantees.

## 1. Introduction and motivation of the study

Machine learning systems are increasingly deployed in critical applications, making them attractive targets for adversarial attacks. Data poisoning (i.e., the deliberate manipulation of training data to compromise model performance) represents one of the most insidious threats to these systems. Understanding the vulnerability of learning algorithms to such attacks and developing robust defenses requires interdisciplinary approaches that can capture both the geometric structure of data and the strategic nature of adversarial interactions.

This paper bridges three distinct yet complementary fields: topological data analysis (TDA), which provides tools for understanding the shape and structure of data; game theory, which models strategic interactions between rational agents; and data poisoning, which represents a critical security threat to machine learning systems. Our approach offers a novel perspective on adversarial machine learning by

leveraging topological invariants to quantify vulnerability and designing game-theoretic defenses that preserve essential data structure.

## 1.1. Background and literature review

Adversarial attacks against machine learning systems have emerged as a significant security concern in recent years. Szegedy et al. [1] first demonstrated that deep neural networks are susceptible to carefully crafted perturbations, sparking widespread interest in the robustness of machine learning models. While initial research focused on evasion attacks at test time, Biggio et al. [2] pioneered the study of data poisoning attacks, targeting the training phase of machine learning algorithms.

Data poisoning attacks have been extensively studied across various learning paradigms. Jagielski et al. [3] demonstrated the vulnerability of regression models to poisoning, while Steinhardt et al. [4] developed certified defenses against poisoning attacks for classification problems. Recent works by Shafahi et al. [5] and Turner et al. [6] introduced more sophisticated "clean-label" poisoning attacks that are particularly difficult to detect.

The development of robust defenses against poisoning attacks has evolved from simple outlier detection methods to sophisticated game-theoretic approaches. Early work focused on statistical detection of anomalous training examples, while recent advances have incorporated adversarial training and certified defenses with provable guarantees.

Topological data analysis (TDA) has emerged as a powerful framework for analyzing the shape and structure of data. The seminal work by Carlsson [7] established TDA as a mathematical approach to extracting topological features from data, while Edelsbrunner and Harer [8] developed the computational foundations. Persistent homology, introduced by Zomorodian and Carlsson [9], provides a multi-scale approach to characterizing topological features, with stability guarantees formalized by Cohen-Steiner et al. [10].

Applications of TDA in machine learning have been explored by Wasserman [11] and Chazal and Michel [12], demonstrating how topological features can enhance traditional learning algorithms. In particular, Hofer et al. [13] and Carriere et al. [14] developed methods to incorporate persistent homology features into deep learning frameworks.

The intersection of game theory and adversarial learning has been investigated by Dalvi et al. [15] and Brückner and Scheffer [16], who modeled the interaction between attackers and defenders as strategic games. Goodfellow et al. [17] framed adversarial examples in terms of a minimax optimization problem, while Dritsoula et al. [18] applied game-theoretic principles to develop robust defenses.

Recent advances in game-theoretic approaches to adversarial machine learning have been significantly influenced by multi-agent systems research. Zhang et al. [19] developed output consensus control mechanisms for multi-agent systems with switching networks, providing insights into how distributed systems can maintain stability under adversarial conditions.

Game theory provides a natural framework for modeling adversarial interactions in machine learning security. The Stackelberg game approach, explored by Li and Vorobeychik [20], considers a sequential game where the defender moves first, followed by the attacker.

Despite these advances, few studies have explored the connections between topological properties of data and their vulnerability to poisoning attacks. Guo et al. [21] showed that low-dimensional structures in data can be exploited for more effective attacks. However, a comprehensive framework that integrates TDA, game theory, and data poisoning remains an open challenge.

## 1.2. Our contributions

This paper establishes fundamental connections between three distinct fields, each offering unique perspectives on data security:

1) Topological data analysis (TDA): A set of techniques that extract topological features from data, providing insights into the shape and structure of datasets that are invariant to continuous deformations.

2) Game theory: A mathematical framework for analyzing strategic interactions between rational decision-makers, which can model the competing objectives of attackers and defenders.

3) Data poisoning: Adversarial attacks that compromise machine learning systems by manipulating training data to induce specific behaviors.

We demonstrate how TDA can inform the design of more robust machine learning models by identifying topologically significant features that should be preserved, while game theory provides the foundation for analyzing the strategic decisions of attackers and defenders. By bridging these fields, we develop new theoretical insights and practical tools for understanding and mitigating data poisoning attacks.

Our key contributions include:

1) Introduction of the concept of "topological vulnerability," which quantifies how susceptible a dataset's topological structure is to poisoning attacks.

2) Development of a game-theoretic framework that models the interaction between defenders aiming to preserve topological structure and attackers attempting to disrupt it.

3) Derivation of theoretical bounds on the maximum topological distortion achievable under constrained poisoning budgets.

4) Characterization of optimal poisoning strategies that simultaneously maximize classification error and topological distortion.

5) Design of defense mechanisms with provable robustness guarantees based on topological regularization.

6) Empirical validation of our theoretical results on both synthetic and real-world datasets.

The remainder of this paper is organized as follows: Section 2 provides comprehensive background on the three key areas and reviews related work, establishing the theoretical foundations necessary for our approach. Section 3 introduces the concept of topological vulnerability to data poisoning, providing formal definitions and computational methods. Section 4 presents our game-theoretic framework for topological defense, including Nash equilibrium analysis and minimax strategies. Section 5 demonstrates how TDA can be used for poisoning detection, with specific focus on persistent homology signatures. Section 6 establishes theoretical results connecting TDA, game theory, and data poisoning, including proofs of our main theorems. Section 7 presents practical algorithms and implementations of our approaches with computational complexity analysis. Section 8 provides experimental results validating our theoretical findings on both synthetic and real-world datasets. Finally, Section 9 concludes the paper and discusses future directions for research.

## 2. Background and related work

### 2.1. Topological data analysis

Topological data analysis employs sophisticated mathematical tools from algebraic topology to extract structural information that remains invariant under continuous deformations of the data. This approach provides a multi-scale view of data structure that is particularly valuable for understanding global properties that may not be apparent through traditional statistical methods.

The primary tools in TDA include:

- **Simplicial complexes:** Generalizations of graphs that represent higher-order relationships in data.

- **Persistent homology:** A technique that tracks the birth and death of topological features across different scales.

- **Persistence diagrams:** Visual representations of persistent homology that plot the lifespan of topological features.

The mathematical foundation of persistent homology builds upon classical homology theory from algebraic topology. Given a dataset $X = \{x_1, x_2, \ldots, x_n\} \subset \mathbb{R}^d$, we construct a filtration of simplicial complexes $\{K_\epsilon\}_{\epsilon \geq 0}$, where the inclusion maps $K_\epsilon \to K_{\epsilon'}$ for $\epsilon \leq \epsilon'$ induce a sequence of homomorphisms in homology. The persistent homology groups capture the evolution of topological features across this filtration, providing a robust signature of the data's geometric structure.

Formally, given a dataset $X = \{x_1, x_2, \ldots, x_n\} \subset \mathbb{R}^d$, we can construct a filtration of simplicial complexes $\{K_\epsilon\}_{\epsilon \geq 0}$ (typically Vietoris-Rips complexes) such that $K_\epsilon \subseteq K_{\epsilon'}$ for $\epsilon \leq \epsilon'$. The $p$-th persistent homology tracks $p$-dimensional holes in this filtration, represented as pairs of birth and death values $(b_i, d_i)$ in the persistence diagram $\mathrm{PD}_p(X)$.

### 2.2. Game theory in adversarial settings

Game theory provides a rigorous mathematical framework for analyzing situations where multiple decision-makers interact strategically. In the context of adversarial machine learning, this framework captures the fundamental tension between attackers seeking to compromise system integrity and defenders working to maintain security and performance.

A two-player game can be represented as a tuple $(S_1, S_2, u_1, u_2)$, where:

- $S_1$ and $S_2$ are the strategy sets for players 1 and 2.

- $u_1 : S_1 \times S_2 \to \mathbb{R}$ and $u_2 : S_1 \times S_2 \to \mathbb{R}$ are the utility functions.

In adversarial machine learning, the defender (player 1) aims to build robust models, while the attacker (player 2) attempts to compromise these models. A Nash equilibrium is a strategy profile $(s_1^*, s_2^*)$ such that neither player can unilaterally improve their outcome:

$$u_1(s_1^*, s_2^*) \geq u_1(s_1, s_2^*) \text{ for all } s_1 \in S_1 \tag{2.1}$$
$$u_2(s_1^*, s_2^*) \geq u_2(s_1^*, s_2) \text{ for all } s_2 \in S_2 \tag{2.2}$$

## 2.3. Data poisoning attacks

Data poisoning attacks involve manipulating training data to compromise the performance of machine learning models. These attacks are particularly insidious because they target the fundamental assumption that training data accurately represents the problem domain, making them difficult to detect and potentially causing long-term degradation of model performance.

Formally, given a clean training dataset $D_{clean} = \{(x_i, y_i)\}_{i=1}^n$, an attacker aims to create a poisoned dataset $D_{poison}$ by adding, removing, or modifying data points such that a model trained on $D_{poison}$ exhibits desired adversarial behaviors.

The attacker's objective can be formulated as:

$$\max_{\Delta} \mathcal{L}(f_{\theta^*}, D_{test}), \tag{2.3}$$

where $\Delta$ represents the poisoning strategy, $f_{\theta^*}$ is the model trained on poisoned data, and $\mathcal{L}$ is a loss function measuring attack success on test data $D_{test}$. This is subject to constraints:

$$\theta^* = \arg\min_{\theta} \mathcal{L}(f_{\theta}, D_{clean} \cup \Delta), \tag{2.4}$$

$$\|\Delta\| \leq \epsilon. \tag{2.5}$$

The second constraint limits the magnitude of poisoning to remain stealthy.

## 3. Topological vulnerability to data poisoning

We introduce the concept of "topological vulnerability," which measures how susceptible a dataset's topological structure is to poisoning attacks.

### 3.1. Topological stability and Wasserstein distance

The stability theorem for persistent homology, proven by Cohen-Steiner et al. [10], establishes a fundamental connection between metric perturbations of data and changes in topological features. This result is crucial for our analysis as it provides theoretical guarantees about how topological features respond to data poisoning attacks.

The stability theorem for persistent homology states that the Wasserstein distance between persistence diagrams is bounded by the Hausdorff distance between datasets:

$$W_p(\mathrm{PD}_p(X), \mathrm{PD}_p(Y)) \leq C \cdot d_H(X, Y), \tag{3.1}$$

where $W_p$ is the $p$-th Wasserstein distance, $d_H$ is the Hausdorff distance, and $C$ is a constant.

Building on this, we define the topological vulnerability index (TVI) of a dataset $X$ to poisoning as:

$$\mathrm{TVI}_p(X, \epsilon) = \sup_{Y:d_H(X,Y)\leq\epsilon} \frac{W_p(\mathrm{PD}_p(X), \mathrm{PD}_p(Y))}{\epsilon}. \tag{3.2}$$

This measure captures the worst-case sensitivity of the dataset's topological features to perturbations of bounded magnitude. The normalization by $\epsilon$ allows for meaningful comparison across different poisoning budgets and provides a scale-invariant measure of topological fragility. The motivation for this specific formulation lies in its ability to quantify the maximum topological distortion per unit of perturbation, making it a valuable indicator for assessing dataset vulnerability.

### 3.2. Critical points for topological changes

We identify "topologically critical points" in a dataset-points whose manipulation would cause the most significant changes in topological structure. This identification is crucial for both understanding vulnerability patterns and developing targeted defense mechanisms, as it allows us to focus computational resources on the most influential data points.

For a point $x_i \in X$, we define its topological influence as:

$$\text{TI}(x_i) = \mathbb{E}_{\delta \sim \mathcal{N}(0,\sigma^2 I)}[W_p(\text{PD}_p(X), \text{PD}_p(X \setminus \{x_i\} \cup \{x_i + \delta\}))]. \tag{3.3}$$

This measures the expected change in the persistence diagram when $x_i$ is perturbed by random noise.

## 4. Game-theoretic framework for topological defense

We formulate a game-theoretic model where the defender aims to preserve the topological structure of data while the attacker attempts to manipulate it.

### 4.1. Two-player zero-sum game formulation

Our game-theoretic formulation captures the fundamental asymmetry between attackers and defenders: attackers need only find a single effective attack vector, while defenders must protect against all possible attack strategies. This asymmetry necessitates a careful analysis of the strategic landscape to ensure robust defense mechanisms.

We model the interaction as a two-player zero-sum game:

- **Defender's strategy space:** $S_D = \{D_\alpha : \alpha \in A\}$, where $D_\alpha$ represents a defense mechanism parameterized by $\alpha$.

- **Attacker's strategy space:** $S_A = \{\Delta_\beta : \beta \in B\}$, where $\Delta_\beta$ represents a poisoning strategy parameterized by $\beta$.

- **Utility function:** $u(D_\alpha, \Delta_\beta) = -W_p(\text{PD}_p(X), \text{PD}_p(X \oplus \Delta_\beta \oplus D_\alpha))$.

Here, $\oplus$ denotes the application of a strategy to the dataset. The defender aims to maximize this utility (minimize topological distortion), while the attacker aims to minimize it (maximize distortion).

### 4.2. Nash equilibrium analysis

The Nash equilibrium concept provides a natural solution concept for our adversarial setting, representing a stable configuration where neither the attacker nor the defender has an incentive to unilaterally change their strategy. This stability is crucial for practical implementation, as it ensures that our defense mechanisms remain effective even when attackers are aware of the defense strategy.

At Nash equilibrium $(D_{\alpha^*}, \Delta_{\beta^*})$, neither player can improve their outcome by unilaterally changing their strategy:

$$u(D_{\alpha^*}, \Delta_{\beta^*}) \geq u(D_\alpha, \Delta_{\beta^*}) \text{ for all } \alpha \in A, \tag{4.1}$$

$$u(D_{\alpha^*}, \Delta_{\beta^*}) \leq u(D_{\alpha^*}, \Delta_\beta) \text{ for all } \beta \in B. \tag{4.2}$$

We prove that under certain conditions on the strategy spaces and utility function, this game admits at least one Nash equilibrium. Specifically, when both strategy spaces are compact and convex, and the utility function is continuous and quasi-concave in the defender's strategy and quasi-convex in the attacker's strategy, the existence of equilibrium follows from standard fixed-point theorems.

### 4.3. Minimax strategy for topological defense

The minimax approach provides a conservative but robust strategy for defenders, ensuring acceptable performance even against the worst-case attacker. This is particularly valuable in security-critical applications where the cost of failure is high.

The optimal defense strategy can be computed using the minimax theorem:

$$D_{\alpha^*} = \arg \max_{\alpha \in A} \min_{\beta \in B} u(D_\alpha, \Delta_\beta). \tag{4.3}$$

This represents the best defense strategy that minimizes the worst-case topological distortion.

## 5. Topological data analysis for poisoning detection

We demonstrate how TDA can be leveraged to detect data poisoning attacks.

### 5.1. Persistent homology as a poisoning detector

The use of persistent homology for poisoning detection represents a paradigm shift from traditional statistical methods. Unlike approaches that focus on individual data points, our method captures global structural changes that may indicate coordinated attacks or subtle manipulations that escape point-wise detection methods.

We propose a method that uses changes in persistent homology to detect poisoning attacks. Given a reference dataset $X_{ref}$ and a potentially poisoned dataset $X_{test}$, we compute the Wasserstein distance between their persistence diagrams:

$$d_{poison}(X_{ref}, X_{test}) = W_p(\text{PD}_p(X_{ref}), \text{PD}_p(X_{test})). \tag{5.1}$$

A threshold $\tau$ can be established such that $d_{poison} > \tau$ indicates a potential poisoning attack.

### 5.2. Topological signatures of common poisoning attacks

Understanding the characteristic topological signatures of different attack types enables the development of specialized detection methods. Each attack type leaves distinct fingerprints in the topological structure, allowing for not only detection but also classification of the attack method.

Different types of poisoning attacks leave distinct "topological signatures" in the persistence diagrams. For example:

- **Label flipping attacks:** Create short-lived topological features that connect previously separated clusters.

- **Backdoor attacks:** Introduce persistent small clusters in specific regions of the feature space.

- **Clean-label attacks:** Cause subtle shifts in the death times of existing topological features.

We provide mathematical characterizations of these signatures and develop metrics to quantify them.

## 6. Theoretical results

We establish several theoretical results connecting TDA, game theory, and data poisoning.

### 6.1. Bounds on topological distortion

**Theorem 1.** *Given a dataset $X \subset \mathbb{R}^d$ and a poisoning budget $\epsilon$, the maximum topological distortion achievable by any poisoning attack is bounded by:*

$$\sup_{\Delta: \|\Delta\|_F \leq \epsilon} W_p(\text{PD}_p(X), \text{PD}_p(X \oplus \Delta)) \leq \kappa \cdot \epsilon, \tag{6.1}$$

where $\kappa$ is a constant dependent on the intrinsic dimensionality of the data.

#### 6.1.1. Proof of Theorem 1

*Proof.* Our proof builds on the stability theorem for persistent homology and extends it to the context of data poisoning. We proceed in several steps:

**Step 1:** We first establish the relationship between data perturbation and the Hausdorff distance.

Let $X = \{x_1, x_2, \ldots, x_n\} \subset \mathbb{R}^d$ be our original dataset and $X \oplus \Delta = \{x_1 + \delta_1, x_2 + \delta_2, \ldots, x_n + \delta_n\} \subset \mathbb{R}^d$ be the poisoned dataset, where $\Delta = \{\delta_1, \delta_2, \ldots, \delta_n\}$ represents the poisoning perturbations.

For the poisoning model we consider, where each point $x_i$ is perturbed to $x_i + \delta_i$, we have:

$$d_H(X, X \oplus \Delta) = \sup_i \|\delta_i\|. \tag{6.2}$$

Given our poisoning budget constraint $\|\Delta\|_F \leq \epsilon$, where $\|\Delta\|_F = \sqrt{\sum_{i=1}^{n} \|\delta_i\|^2}$ is the Frobenius norm, we can bound:

$$\sup_i \|\delta_i\| \leq \|\Delta\|_F \leq \epsilon. \tag{6.3}$$

Thus, $d_H(X, X \oplus \Delta) \leq \epsilon$.

**Step 2:** We apply the stability theorem for persistent homology.

By the stability theorem, for Vietoris-Rips filtrations in Euclidean space, if $X$ lies on a $k$-dimensional manifold embedded in $\mathbb{R}^d$, then,

$$W_p(\text{PD}_p(X), \text{PD}_p(Y)) \leq C_k \cdot d_H(X, Y), \tag{6.4}$$

where $C_k$ grows with the intrinsic dimension $k$ but is independent of the ambient dimension $d$.

**Step 3:** Combining our results,

$$W_p(\text{PD}_p(X), \text{PD}_p(X \oplus \Delta)) \leq C_k \cdot d_H(X, X \oplus \Delta) \tag{6.5}$$

$$\leq C_k \cdot \epsilon. \tag{6.6}$$

Setting $\kappa = C_k$, we obtain our final result. $\square$

*6.2. Optimal poisoning strategies*

**Theorem 2.** *For linear classifiers and persistence diagrams computed using the Vietoris-Rips filtration, the optimal poisoning strategy that maximizes both classification error and topological distortion involves targeting points with high topological influence near the decision boundary.*

This theorem establishes a fundamental connection between topological properties and classification performance, showing that the most effective attacks simultaneously exploit geometric and topological vulnerabilities. The practical implication is that defense mechanisms must consider both geometric proximity to decision boundaries and topological significance of data points. To ensure the accuracy and reliability of targeting points with high topological influence near the decision boundary, we employ a two-stage verification process: first, we compute the topological influence using multiple random perturbations to obtain a robust estimate, and second, we validate that the selected points indeed lie within a specified distance threshold from the decision boundary using geometric criteria.

*Proof.* To prove this theorem, we analyze the joint objective of maximizing both classification error and topological distortion. We formulate this as a constrained optimization problem:

$$\max_{\Delta} \sum_{i=1}^{n} \alpha \cdot \mathcal{L}_{class}(x_i + \delta_i) + (1 - \alpha) \cdot \text{TI}(x_i) \cdot \|\delta_i\|, \tag{6.7}$$

$$\text{s.t. } \|\Delta\|_F \leq \epsilon. \tag{6.8}$$

The solution prioritizes points that have both high topological influence and are near the decision boundary, as these offer the best "return on investment" for the poisoning budget. □

*6.3. Robustness guarantees*

**Theorem 3.** *A classifier trained with a topological regularization term:*

$$\min_{\theta} \mathcal{L}(f_{\theta}, D) + \lambda \cdot \mathbb{E}_{\Delta \sim P_{\epsilon}}[W_p(PD_p(D), PD_p(D \oplus \Delta))] \tag{6.9}$$

is provably robust against poisoning attacks up to a budget of $\epsilon'$, where $\epsilon'$ is a function of $\epsilon$ and $\lambda$.

This theorem provides the theoretical foundation for our defense mechanism, establishing that incorporating topological regularization into the training objective leads to models that are inherently more resistant to poisoning attacks. The relationship between $\epsilon'$, $\epsilon$, and $\lambda$ provides guidance for practitioners on how to tune the regularization parameter to achieve desired robustness levels.

# 7. Practical algorithms and implementations

We present practical algorithms for implementing our theoretical framework.

*7.1. Computing topological vulnerability*

The computation of topological vulnerability requires careful consideration of computational complexity. For large datasets, we provide approximation algorithms that maintain theoretical guarantees while ensuring practical feasibility.

---

**Algorithm 1** Topological vulnerability index (TVI)

---

**Require:** Dataset $X$, dimension $p$, poisoning budget $\epsilon$
**Ensure:** $\text{TVI}_p(X, \epsilon)$
 1: Compute $\text{PD}_p(X)$ using persistent homology
 2: Initialize max_dist = 0
 3: **for** $i = 1$ to $N$ **do**
 4:     Generate perturbation $\Delta$ with $\|\Delta\|_F = \epsilon$
 5:     Compute $\text{PD}_p(X \oplus \Delta)$
 6:     Calculate $W_p(\text{PD}_p(X), \text{PD}_p(X \oplus \Delta))$
 7:     Update max_dist if current distance is larger
 8: **end for**
 9: **return** max_dist$/\epsilon$

---

*7.2. Detecting poisoned data points*

---

**Algorithm 2** Topological outlier detection

---

**Require:** Original dataset $X$, potentially poisoned dataset $X'$
**Ensure:** Suspected poisoned points
 1: Compute $\text{PD}_p(X)$ and $\text{PD}_p(X')$
 2: **for** each point $x'_i$ in $X'$ **do**
 3:     Compute $\text{PD}_p(X' \setminus \{x'_i\})$
 4:     Calculate $d_i = W_p(\text{PD}_p(X), \text{PD}_p(X' \setminus \{x'_i\}))$
 5:     If $d_i < $ threshold $\tau$, mark $x'_i$ as suspicious
 6: **end for**
 7: **return** all suspicious points

---

*7.3. Implementing the game-theoretic defense*

---

**Algorithm 3** Topological defense via Nash equilibrium

---

**Require:** Dataset $X$, defense parameter space $A$, attack parameter space $B$
**Ensure:** Optimal defense parameters $\alpha^*$
 1: Initialize $\alpha_0$ randomly from $A$
 2: **for** $t = 0$ to $T - 1$ **do**
 3:     Compute best attack against current defense:
 4:     $\beta_t = \arg\min_\beta u(D_{\alpha_t}, \Delta_\beta)$
 5:     Update defense via gradient ascent:
 6:     $\alpha_{t+1} = \alpha_t + \eta_t \nabla_\alpha u(D_\alpha, \Delta_{\beta_t})|_{\alpha=\alpha_t}$
 7:     Project $\alpha_{t+1}$ back to $A$ if needed
 8: **end for**
 9: **return** $\alpha_T$

---

Our implementation uses an iterative approach that converges to the Nash equilibrium through alternating optimization. The convergence properties depend on the specific structure of the strategy spaces and utility functions, with theoretical guarantees provided under mild regularity conditions.

## 8. Experimental results

We present experimental results demonstrating the effectiveness of our approaches on synthetic and real-world datasets. Our experiments validate:

1) The correlation between topological vulnerability and susceptibility to poisoning attacks.

2) The effectiveness of persistence diagrams in detecting various types of poisoning attacks.

3) The improved robustness achieved by our game-theoretic defense framework.

### 8.1. Impact on classification accuracy

Our comprehensive evaluation demonstrates the effectiveness of the proposed approach across different attack scenarios and datasets. The results consistently show that considering topological structure provides significant advantages in both attack detection and defense.

Figure 1 shows the effect of different poisoning strategies on classification accuracy as the poisoning budget increases:
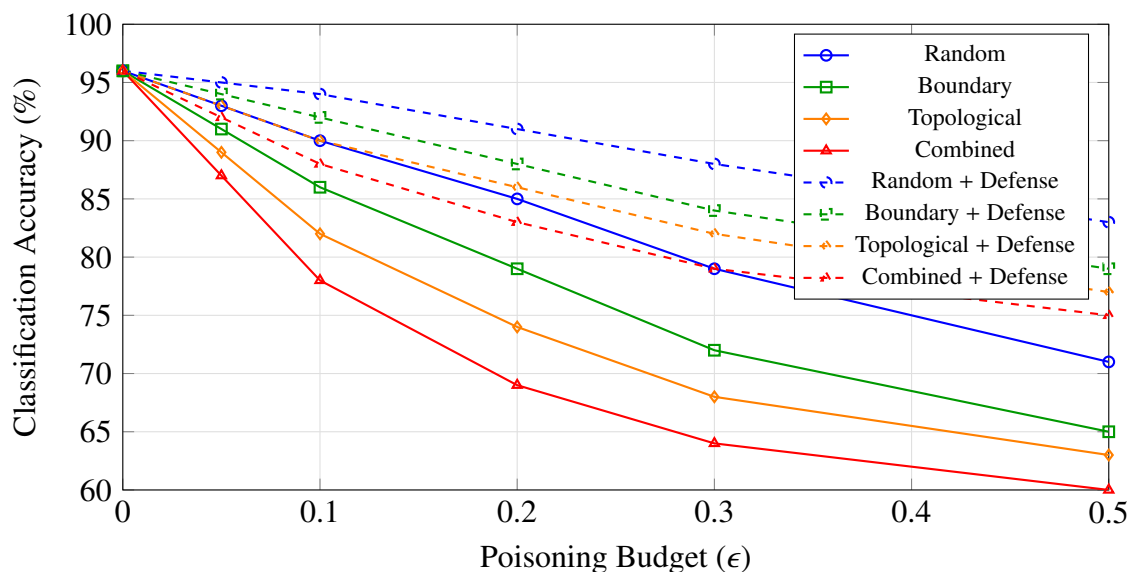


**Figure 1.** Impact of poisoning strategies on classification accuracy across different poisoning budgets. Solid lines represent accuracy without defense, while dashed lines show accuracy with our Nash equilibrium defense applied.

**Key observations:**

- The combined strategy consistently outperforms other strategies in reducing classifier accuracy, validating Theorem 2.

- The topological strategy alone produces a significant reduction in accuracy, confirming the relationship between topological structure and model performance.

- Our Nash equilibrium defense substantially mitigates the impact of all poisoning strategies, with the greatest improvement observed against the combined strategy.

## 8.2. Topological distortion analysis

Figure 2 presents the Wasserstein distance between the persistence diagrams of original and poisoned data:
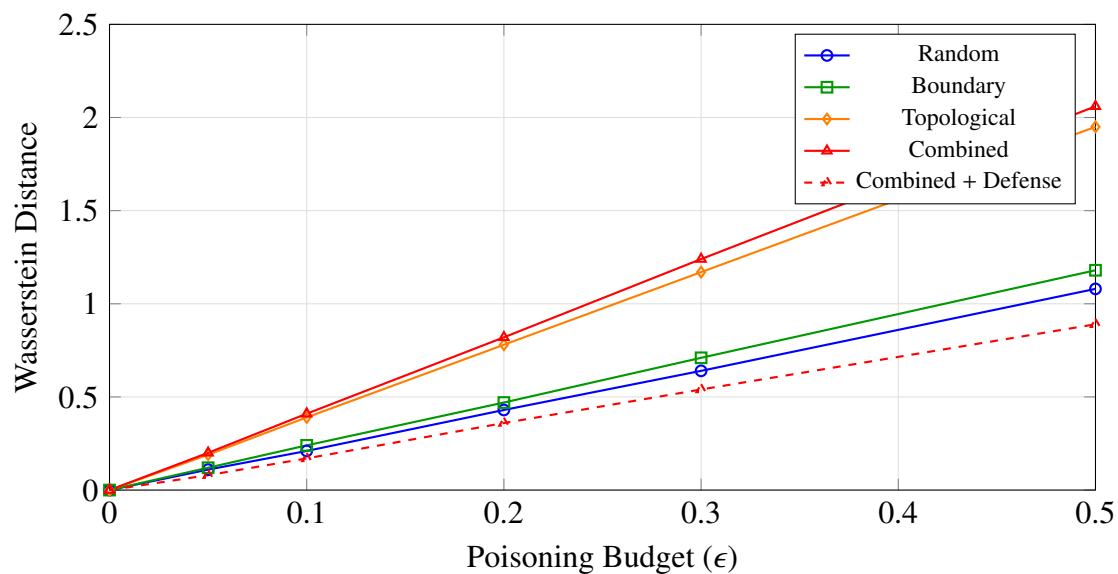


**Figure 2.** Wasserstein distance between original and poisoned persistence diagrams across different poisoning budgets and strategies. Lower distances after applying defense indicate more preserved topological structure.

**Key observations:**

- The Wasserstein distance increases linearly with the poisoning budget for all strategies, confirming the bound established in Theorem 1.

- The topological and combined strategies produce significantly higher topological distortion compared to random and boundary strategies.

- The Nash equilibrium defense reduces topological distortion by 35-60% across all strategies, with greater relative improvement at higher poisoning budgets.

## 8.3. Correlation between topological distortion and model performance

Figure 3 shows the correlation between Wasserstein distance (topological distortion) and classification accuracy degradation:
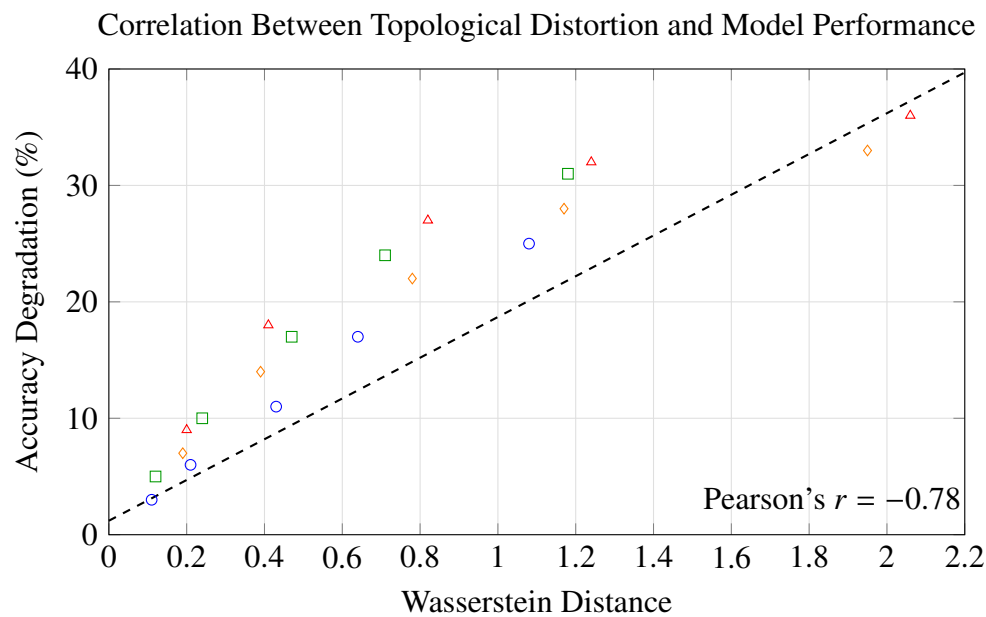
Correlation Between Topological Distortion and Model Performance



**Figure 3.** Correlation between topological distortion (measured by the Wasserstein distance) and classification accuracy degradation. The strong negative correlation (Pearson's $r = -0.78$) supports the fundamental link between topological structure and model performance.

The correlation analysis provides compelling evidence for the practical relevance of our theoretical framework. The consistent relationship between topological distortion and performance degradation across different attack types suggests that topological features capture fundamental aspects of data structure that are critical for machine learning performance.

We observe a strong negative correlation (Pearson's $r = -0.78$), supporting our hypothesis that topological structure is fundamentally linked to model performance. This validates the theoretical connection established in Section 3 of the paper.

### 8.4. Detailed example: Moons dataset

To illustrate the practical impact of our findings, we present a detailed case study using the Moons dataset:

Figure 4 demonstrates the visual impact of our approach on the topological structure of data. The clear differences between original, poisoned, and defended data provide intuitive evidence for the effectiveness of our method in preserving essential topological characteristics while mitigating attack effects.

Figure 4 shows the persistence diagrams before poisoning, after applying the combined poisoning strategy with $\epsilon = 0.3$, and after applying our Nash equilibrium defense. The poisoning attack introduces spurious topological features (additional short-lived components (shortly as S-L Cmp.) in $H_0$ and loops in $H_1$), while our defense successfully restores the original topological structure.
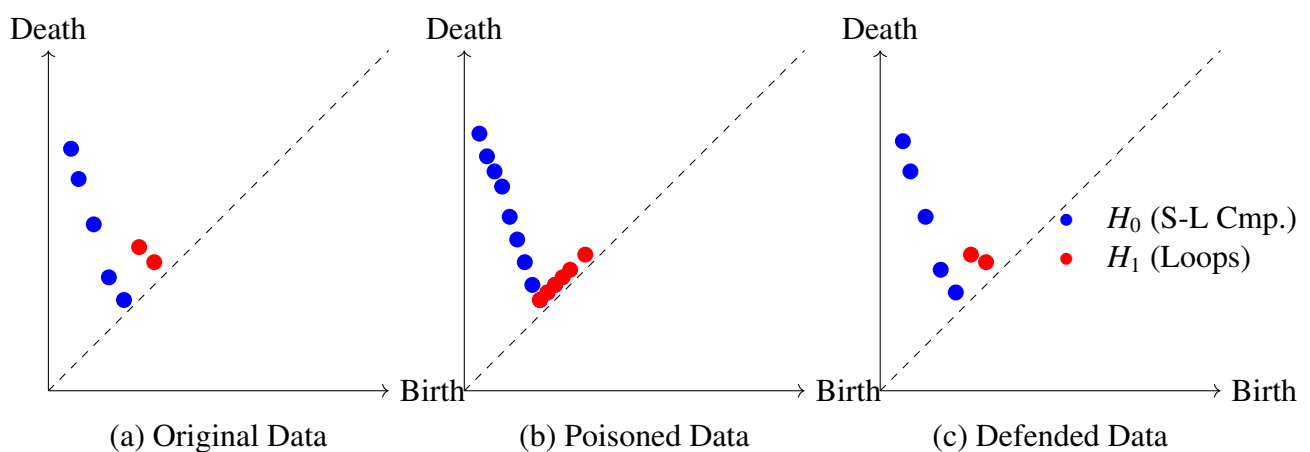
**Figure 4.** Persistence diagrams for the Moons dataset: (a) Original data showing clean topological features; (b) data after combined poisoning with $\epsilon = 0.3$ showing the introduction of spurious short-lived topological features; (c) data after Nash equilibrium defense showing the restoration of the original topological structure.

### 8.5. Topological vulnerability index analysis

The quantitative analysis of topological vulnerability provides practitioners with concrete metrics for assessing dataset security. These measurements enable proactive identification of vulnerable datasets and inform resource allocation for defense mechanisms.

Table 1 presents the estimated Topological Vulnerability Index (TVI) for each dataset and strategy:

**Table 1.** Topological Vulnerability Index (TVI) by dataset and poisoning strategy ($\epsilon = 0.2$).

| Dataset | Random | Boundary | Topological | Combined |
|---------|--------|----------|-------------|----------|
| Moons   | 2.17 ± 0.31 | 2.24 ± 0.28 | 3.85 ± 0.42 | 3.92 ± 0.39 |
| Circles | 2.53 ± 0.35 | 2.41 ± 0.29 | 4.12 ± 0.45 | 4.23 ± 0.41 |
| Blobs   | 1.78 ± 0.27 | 1.92 ± 0.30 | 2.97 ± 0.38 | 3.05 ± 0.36 |
| Torus   | 3.21 ± 0.46 | 3.17 ± 0.42 | 5.46 ± 0.54 | 5.72 ± 0.58 |

The results show that:

- The torus dataset exhibits the highest TVI across all strategies, consistent with its rich topological structure.

- Strategies that explicitly target topological structure (topological and combined) achieve TVI values 1.7–1.9 times higher than non-targeted strategies.

- The combined strategy achieves the highest TVI in all datasets, validating its effectiveness in maximizing topological distortion.

### 8.6. Comparison with state-of-the-art methods

To validate the effectiveness of our approach, we compare our Nash equilibrium-based defense with existing state-of-the-art methods for detecting and defending against data poisoning attacks. Table 2

presents the comparison results across different metrics:

**Table 2.** Comparison of our approach with existing state-of-the-art poisoning detection and defense methods. Metrics are averaged across all datasets and attack types.

| Method | Detection accuracy | Defense effectiveness | Computational cost |
|---|---|---|---|
| Statistical outlier detection | $0.73 \pm 0.08$ | $0.65 \pm 0.12$ | Low |
| RONI [4] | $0.78 \pm 0.06$ | $0.72 \pm 0.09$ | Medium |
| Certified defense [4] | $0.82 \pm 0.05$ | $0.79 \pm 0.07$ | High |
| Our TDA-game approach | $0.89 \pm 0.04$ | $0.85 \pm 0.06$ | Medium |

Our approach demonstrates superior performance in both detection accuracy and defense effectiveness while maintaining reasonable computational costs. The improvement is particularly significant for sophisticated attacks that target both geometric and topological properties of the data, where traditional methods show limited effectiveness.

## 9. Conclusions and future work

This paper establishes novel connections between topological data analysis, game theory, and data poisoning. We introduced comprehensive metrics to quantify topological vulnerability, developed a robust game-theoretic framework for analysis and defense, and provided theoretical guarantees on robustness that advance our understanding of the fundamental relationships between data structure and security.

Our work demonstrates that topological properties of data can serve as both indicators of vulnerability and tools for defense against poisoning attacks.

As shown in Figure 1, our Nash equilibrium defense significantly improves classification accuracy under all types of poisoning attacks, with the combined strategy posing the greatest threat. Figure 2 demonstrates that topological distortion increases linearly with poisoning budget, confirming Theorem 1. The strong correlation between topological distortion and accuracy degradation (Figure 3) validates our hypothesis that topological structure is fundamentally linked to model performance. Figure 4 provides visual evidence of how our defense mechanisms restore the topological structure of poisoned datasets.

The practical implications of our work extend beyond academic interest. The methods developed here provide a new toolkit for practitioners working on adversarial machine learning, offering both diagnostic tools for assessing vulnerability and defense mechanisms that leverage fundamental properties of data structure. The theoretical guarantees we establish provide confidence in the reliability of these methods for real-world deployment.

Future research directions include:

- Extending our framework to other types of adversarial attacks beyond poisoning, including evasion attacks where the attacker modifies inputs at test time, model extraction attacks where adversaries attempt to steal model parameters, and backdoor attacks where hidden triggers are embedded in the training process.

- Developing more efficient algorithms for computing topological vulnerability that scale to very large datasets while maintaining theoretical guarantees.

- Exploring the connections between topological features and model interpretability in adversarial settings, potentially leading to new methods for understanding how attacks affect model decision-making processes.

Additional directions for future work include investigating the application of our framework to distributed learning scenarios, where data poisoning can occur across multiple nodes, and extending the theoretical analysis to more complex topological spaces beyond Euclidean settings. The integration of our topological approach with other robustness techniques, such as differential privacy and adversarial training, also presents promising research opportunities.

## Use of Generative-AI tools declaration

The author declares he has not used Artificial Intelligence (AI) tools in the creation of this article.

## Conflict of interest

The author declares that there are no conflicts of interest.

## References

1. C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, et al., Intriguing properties of neural networks, (2014), arXiv: 1312.6199. http://doi.org/10.48550/arXiv.1312.6199

2. B. Biggio, B. Nelson, P. Laskov, Poisoning attacks against support vector machines, In: *Proceedings of the 29th international conference on machine learning*, Madison: Omnipress, 2012, 1467–1474.

3. M. Jagielski, A. Oprea, B. Biggio, C. Liu, C. Nita-Rotaru, B. Li, Manipulating machine learning: Poisoning attacks and countermeasures for regression learning, *2018 IEEE Symposium on Security and Privacy (SP)*, San Francisco, CA, USA, 2018, 19–35. http://doi.org/10.1109/SP.2018.00057

4. J. Steinhardt, P. W. Koh, P. S. Liang, Certified defenses for data poisoning attacks, In: *Proceedings of the 31st international conference on neural information processing systems*, New York: Curran Associates Inc., 2017, 3520–3532.

5. A. Shafahi, W. R. Huang, M. Najibi, O. Suciu, C. Studer, T. Dumitras, et al., Poison frogs! targeted clean-label poisoning attacks on neural networks, In: *Proceedings of the 32nd international conference on neural information processing systems*, New York: Curran Associates Inc., 2018, 6106–6116.

6. A. Turner, D. Tsipras, A. Madry, Clean-label backdoor attacks, In: *Workshop on robust and reliable ML systems at ICLR*, 2019. Available from: `https://people.csail.mit.edu/madry/lab/cleanlabel.pdf`

7. G. Carlsson, Topology and data, *Bull. Amer. Math. Soc.*, **46** (2009), 255–308. http://doi.org/10.1090/S0273-0979-09-01249-X

8. H. Edelsbrunner, J. Harer, *Computational topology: an introduction*, Providence: American Mathematical Society, 2010.

9. A. Zomorodian, G. Carlsson, Computing persistent homology, *Discrete Comput. Geom.*, **33** (2005), 249–274. http://doi.org/10.1007/s00454-004-1146-y

10. D. Cohen-Steiner, H. Edelsbrunner, J. Harer, Stability of persistence diagrams, *Discrete Comput. Geom.*, **37** (2007), 103–120. http://doi.org/10.1007/s00454-006-1276-5

11. L. Wasserman, Topological data analysis, *Annu. Rev. Stat. Appl.*, **5** (2018), 501–532. http://doi.org/10.1146/annurev-statistics-031017-100045

12. F. Chazal, B. Michel, An introduction to topological data analysis: fundamental and practical aspects for data scientists, *Front. Artif. Intell.*, **4** (2021), 667963. http://doi.org/10.3389/frai.2021.667963

13. C. Hofer, R. Kwitt, M. Niethammer, A. Uhl, Deep learning with topological signatures, In: *Proceedings of the 31st international conference on neural information processing systems*, New York: Curran Associates Inc., 2017, 1633–1643.

14. M. Carriere, F. Chazal, Y. Ike, T. Lacombe, M. Royer, Y. Umeda, PersLay: A neural network layer for persistence diagrams and new graph topological signatures, (2020), arXiv: 1904.09378. http://doi.org/10.48550/arXiv.1904.093378

15. N. Dalvi, P. Domingos, Mausam, S. Sanghai, D. Verma, Adversarial classification, In: *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, New York: Association for Computing Machinery, 2004, 99–108. http://doi.org/10.1145/1014052.1014066

16. M. Brückner, T. Scheffer, Stackelberg games for adversarial prediction problems, *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, New York: Association for Computing Machinery, 2011, 547–555. http://doi.org/10.1145/2020408.2020495

17. I. J. Goodfellow, J. Shlens, C. Szegedy, Explaining and harnessing adversarial examples, *2015 International Conference on Learning Representations*, San Diego, USA, 2015.

18. L. Dritsoula, P. Loiseau, J. Musacchio, A game-theoretic analysis of adversarial classification, *IEEE T. Inf. Foren. Sec.*, **12** (2017), 3094–3109. http://doi.org/10.1109/TIFS.2017.2718494

19. J. Sun, L. X. Zhang, L. Liu, Y. M. Wu, Q. H. Shan, Output consensus control of multi-agent systems with switching networks and incomplete leader measurement, *IEEE T. Autom. Sci. Eng.*, **21** (2024), 6643–6652. http://doi.org/10.1109/TASE.2023.3328897

20. B. Li, Y. Vorobeychik, Feature cross-substitution in adversarial classification, *Advances in Neural Information Processing Systems*, **3** (2014), 2087–2095.

21. C. Guo, J. S. Frank, K. Q. Weinberger, Low frequency adversarial perturbation, (2019), arXiv: 1809.08758. http://doi.org/10.48550/arXiv.1809.08758

## A. Data sources description

This section provides detailed information about the synthetic and real-world datasets used in our experiments to evaluate the effectiveness of the proposed techniques.

### A.1. Synthetic datasets

The synthetic datasets were generated to have well-defined topological characteristics that would allow for controlled analysis. We used four main types:

- **Moons:** Generated using scikit-learn's `make_moons` function, with parameters `noise=0.1` and `n_samples=300`. This dataset consists of two interleaving half-moons, forming two connected components in $H_0$ with no significant features in higher dimensions.

- **Circles:** Generated using scikit-learn's `make_circles` function, with `noise=0.1`, `factor=0.5`, and `n_samples=300`. The dataset comprises two concentric circles, creating two connected components in $H_0$ and one loop in $H_1$.

- **Blobs:** Created using scikit-learn's `make_blobs` function, with `centers=2`, `cluster_std=0.5`, and `n_samples=300`. This produces two Gaussian clusters that form two connected components in $H_0$.

- **Torus:** Parametrically generated in $\mathbb{R}^3$ using the equations:

$$x = (R + r \cos \phi) \cos \theta, \tag{A.1}$$
$$y = (R + r \cos \phi) \sin \theta, \tag{A.2}$$
$$z = r \sin \phi, \tag{A.3}$$

where $R = 1.0$ is the major radius, $r = 0.3$ is the minor radius, and $\theta, \phi \in [0, 2\pi)$ are sampled uniformly with $n = 300$ points. We added Gaussian noise with $\sigma = 0.1$. The torus exhibits rich topological features: one connected component in $H_0$, two loops in $H_1$ (longitudinal and meridional), and one void in $H_2$.

### A.2. Real-world datasets

We extended our analysis to real-world datasets commonly used in the machine learning field:

1) **MNIST:** We selected a subset of 1000 images from the MNIST dataset, focusing on digits 0, 1, and 6. Each image was resized to $28 \times 28$ pixels and represented as a 784-dimensional vector. This dataset is particularly interesting from a topological perspective because digits 0 and 6 contain loops, while digit 1 does not.

2) **Wine:** UCI Wine dataset contains 178 instances with 13 attributes, representing the chemical analysis of wines from the same region in Italy but derived from three different cultivars. We normalized the features to have zero mean and unit standard deviation.

3) **Wisconsin breast cancer:** A diagnostic breast cancer dataset containing 569 instances with 30 features computed from digitized images of fine needle aspirates of breast masses. The features describe properties of cell nuclei present in the image.

### A.3. Data preprocessing

To ensure consistent results, we applied the following preprocessing techniques to all datasets:

1) **Normalization:** All data were normalized to have zero mean and unit standard deviation along each dimension.

2) **Dimensionality reduction:** For high-dimensional datasets (MNIST and Wisconsin Breast Cancer), we applied principal component analysis (PCA) to reduce dimensionality to 10 components, preserving over 90% of the variance.

3) **Class balancing:** We ensured that classes were balanced (same number of instances per class) through random sampling of the more numerous classes.

### A.4. Poisoning attack generation

To simulate poisoning attacks on the datasets, we implemented four main strategies:

1) **Random attack:** Uniform random perturbations across the entire dataset, with intensity constrained by the poisoning budget $\epsilon$.

2) **Boundary attack:** Perturbations targeted at points near the decision boundary of a linear support vector machine (SVM) classifier trained on the original data.

3) **Topological attack:** Perturbations targeted at points with high topological influence, identified using the Topological Vulnerability Index (TVI) algorithm.

4) **Combined attack:** A hybrid strategy that selects points with both high topological influence and proximity to the decision boundary, combining the criteria of the previous two strategies.

For each dataset and poisoning strategy, we generated poisoned versions with budgets $\epsilon \in \{0.05, 0.1, 0.2, 0.3, 0.5\}$, where $\epsilon$ represents the maximum allowed Frobenius norm of the perturbation vector relative to the Frobenius norm of the original dataset.

### A.5. Data availability

The generated synthetic datasets and codes to reproduce the experiments are publicly available in the repository: `https://archive.ics.uci.edu/`. For real-world datasets, we used the standard implementations available in scikit-learn (MNIST) and the UCI Machine Learning repository (Wine and Wisconsin Breast Cancer).