

Influence-Driven Data Poisoning in Graph-Based Semi-Supervised Classifiers

Adriano Franci
SnT, Université de Luxembourg
Luxembourg
adriano.franci@uni.lu

Maxime Cordy
SnT, Université de Luxembourg
Luxembourg
maxime.cordy@uni.lu

Martin Gubri
SnT, Université de Luxembourg
Luxembourg
martin.gubri@uni.lu

Mike Papadakis
SnT, Université de Luxembourg
Luxembourg
michail.papadakis@uni.lu

Yves Le Traon
SnT, Université de Luxembourg
Luxembourg
yves.letraon@uni.lu

ABSTRACT

Graph-based Semi-Supervised Learning (GSSL) is a practical solution to learn from a limited amount of labelled data together with a vast amount of unlabelled data. However, due to their reliance on the known labels to infer the unknown labels, these algorithms are sensitive to data quality. It is therefore essential to study the potential threats related to the labelled data, more specifically, label poisoning. In this paper, we propose a novel data poisoning method which efficiently approximates the result of label inference to identify the inputs which, if poisoned, would produce the highest number of incorrectly inferred labels. We extensively evaluate our approach on three classification problems under 24 different experimental settings each. Compared to the state of the art, our influence-driven attack produces an average increase of error rate 50% higher, while being faster by multiple orders of magnitude. Moreover, our method can inform engineers of inputs that deserve investigation (relabelling them) before training the learning model. We show that relabelling one-third of the poisoned inputs (selected based on their influence) reduces the poisoning effect by 50%.

KEYWORDS

Machine learning, semi-supervised learning, data poisoning

ACM Reference Format:

Adriano Franci, Maxime Cordy, Martin Gubri, Mike Papadakis, and Yves Le Traon. 2022. Influence-Driven Data Poisoning in Graph-Based Semi-Supervised Classifiers. In *1st Conference on AI Engineering - Software Engineering for AI (CAIN'22)*, May 16–24, 2022, Pittsburgh, PA, USA. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3522664.3528606>

1 INTRODUCTION

Recent advances in Machine Learning (ML) resulted in an unprecedented interest towards embedding such technologies in software systems. To be effective though, ML algorithms generally require a massive amount of data together with their expected prediction

outcomes (labels). Such labelling activities are expensive and time-consuming as they are typically performed manually by humans. Thus, acquiring labelled data is seen as a major obstacle to the widespread adoption of ML.

To alleviate this problem, Semi-Supervised Learning (SSL) algorithms [31] were proposed to exploit a limited amount of labelled data together with a vast amount of unlabelled data (whose acquisition is generally inexpensive). Their principle (see Figure 1) is to use the labelled inputs to infer the correct label of unlabelled inputs that are similar (i.e., close in the feature space). Such solutions are flexible: they can be used for both *transductive* learning (i.e., inferring the labels of the unlabelled inputs) and *inductive* learning (i.e., using the initial labelled data and the inferred labels to train a supervised model and make predictions on unseen data). Moreover, SSL was shown to produce significant improvements over “fully” supervised algorithms (which learn from labelled data only) and unsupervised ones (which do not use data labels) [3, 13]. It has been successfully applied in a variety of domains including image classification [8, 10], drug interaction discovery [29] and social media mining [25].

A major thrust in SSL comes from its greater reliance on data quality to perform correct predictions [16]. This makes such algorithms vulnerable to data poisoning [17], i.e. intentional or unintentional alterations of training data that mislead the learning of prediction models (e.g., reducing prediction accuracy). Data poisoning is considered a serious threat to ML-enabled systems [12] that is prominent in cases where the data acquisition process is not entirely trustworthy. This phenomenon has been observed and studied in various contexts such as intrusion detection [7, 21], face recognition [4] and crowdsourcing systems [11, 26]. While much research studied data poisoning in supervised models (see, e.g., [18, 27]), data poisoning in SSL remains largely unexplored: in their NeurIPS’19 paper, Liu et al. proposed the first and only data poisoning framework for SSL [17]. Studying this security aspect is not only important to identify potential threats on safety- and business-critical software systems relying on SSL, but it is also a prerequisite to design appropriate countermeasures.

Existing studies typically rely on dedicated algorithms – named *data poisoning attacks* – that aim to find the data alterations with the maximal impact on SSL. Thanks to these algorithms, researchers and engineers can estimate the impact that the worst-case data alterations will have on the target ML-enabled systems. In turn,

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
CAIN’22, May 16–24, 2022, Pittsburgh, PA, USA
© 2022 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-9275-4/22/05.
<https://doi.org/10.1145/3522664.3528606>

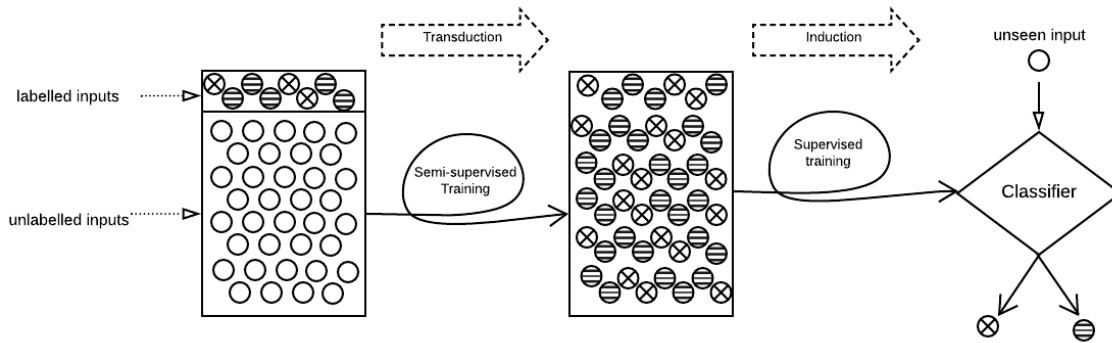


Figure 1: Principles of semi-supervised learning. During transductive learning, the labels of the unlabelled data are inferred from the known labels of the labelled inputs. During inductive learning, the whole data set (with known and inferred labels) are used to train a supervised model, which can be used to classify previously unseen data.

they can design countermeasures to best protect the systems against data poisoning instances.

In this paper, we propose a novel data poisoning attack for SSL which is simple to compute (involves common matrix operations only), effective (in reducing accuracy) and efficient (requires low computation costs). Our attack emulates a common data poisoning scenario where some inputs have received an incorrect label before learning occurs (label poisoning) [19]. Such incorrect labels, in turn, alter the inference capabilities of learning models. The key idea of our method is to exploit mathematical properties of SSL to identify “weak spots” in the labelled data set, i.e. inputs that, if mislabelled, would yield the largest number of wrongly-inferred labels. To do so, we define a metric to approximate the *influence* of labelled inputs onto unlabelled ones. Based on the proposed metric, we can rank labelled inputs wrt. the expected impact of poisoning their label. Inputs ranked higher are the ideal target of attacks since their poisoning yields the highest reduction in prediction accuracy. Hence, we refer to our attack as *influence driven*.

We extensively evaluate our approach on three classification problems from different application domains (image recognition and text classification), in both transductive learning settings and inductive learning settings. As a first step, we empirically demonstrate that our influence metric has medium-to-strong correlations (between 0.31 and 0.99) with the error rate resulting from poisoning the labelled inputs. This indicates that the metric captures well the expected impact of poisoning each of the inputs and, therefore, is convenient to prioritize them. Secondly, we show that our metric remains effective (i) when used to select multiple inputs to poison and (ii) during inductive learning, i.e. when the labelled data resulting from SSL are used to train a supervised model to classify unseen data. Compared to the state of the art, our poisoning attack yields a significantly higher number of new misclassifications than the state of the art (50% higher on average), while being 2 to 3 times faster. Thus, it forms a new, strong baseline for future research in SSL data poisoning.

Finally, we investigate the use case where our influence metric supports engineers in identifying the most critical inputs to relabel (re-investigate the correct label before learning) in order to protect

their system against the effect of label poisoning. We show that ensuring the label correctness of the most influential inputs more than halves the error rate on average, whereas labelling additional inputs brings minor benefits (less than 15% of average error rate reduction).

To summarize, our contribution is four-fold:

- (1) We formalise the concept of input influence in SSL and exploit it to identify the most influential labelled inputs during transductive learning (label inference for unlabelled data).
- (2) We exploit this concept to design a novel data poisoning attack, which consists of altering the label of the most influential inputs.
- (3) We conduct an in-depth experimental study showing the benefits of our approach. Our results indicate that our attack transfers well to different SSL algorithms and to inductive learning settings (using supervised models). Compared to the state of the art [17], our method is significantly more effective (the average increase in model error rate is 50% higher on average) and more efficient (2 to 74 times faster).
- (4) We show that an effective strategy to counter data poisoning is to relabel the most influential inputs, guaranteeing their correctness. Compared to labelling an equivalent number of additional inputs, this is the most effective strategy to alleviate the effect of our attack.

2 BACKGROUND AND RELATED WORK

Semi-supervised learning [31] is a particular form of machine learning that attempts to maximize the benefits of learning from a limited amount of labelled data together with a vast amount of unlabelled data. It contrasts with supervised learning which necessitates all training data to be labelled, and unsupervised learning which does not rely on labels but only exploit statistical relations between the input (e.g., to form clusters).

Among the different families of SSL algorithms (see [31] for an overview), graph-based methods are the most popular because they are effective (in inferring unknown labels from labelled data), efficient (in computation time) and straightforward to implement

(based on common matrix operations). The two established graph-based SSL algorithms are label propagation [32] and label spreading [30]. Both consist of computing the label likelihood of any unlabelled input based on the labels of its neighbouring inputs. This computation follows an iterative process until either the label likelihood values converge or a predefined number of iterations is reached. The key differences between the two algorithms lie in how they compute input similarity and how they propagate the likelihood values from one iteration to the other.

Research on ML security is mainly spearheaded by adversarial machine learning. One distinguishes two types of adversarial attacks [22]: *evasion attacks* and *poisoning attacks*.

Evasion attacks occur at prediction (test) time. They consist of perturbing the input sent to the ML model (yielding an “adversarial input”) to cause misclassification. These attacks have been studied intensively over the recent years [22]. They are also commonly used to improve the thoroughness of ML system test suites by generating failure-inducing inputs [28].

By contrast, poisoning attacks strike at *training time* and aim to degrade the overall performance (e.g., the accuracy) of the ML model [5]. These attacks often consist of either *modifying the training data* (falsifying them) or *injecting new data* (containing intentional mistakes) into the training set. In the former case, the attacker can either modify (selected) features of the poisoned inputs (feature poisoning) or their associated label (label poisoning). Although it has been less studied, the third type of poisoning attack is *algorithm corruption*, during which the attacker alters the logic of the training process, thereby changing the way the ML model learns [9]. The new attack we design in this work belongs to the label poisoning category, which requires fewer assumptions about the attacker capability (e.g., it does not require any knowledge about the features of the inputs or access to the internal structure of the training algorithm).

Data poisoning in SSL has been scarcely studied. To the best of our knowledge, the only approach has been recently designed by Liu et al. [17]. While it supports both regression and binary classification problems, we focus here on classification, which they deem to be the most challenging setting [17]. Like ours, their method leans on the properties of label propagation. It consists of a search algorithm that selects the labelled inputs to poison to maximize the expected increase in error rate, computed as the difference between the ground truth labels and an approximation of the labels that will be inferred. Thus, their method assumes that the attacker has access to the real labels of the unlabelled inputs or can compute them reliably using a surrogate SSL model. We do not make such assumption and select the labelled inputs to be poisoned based only on their relative influence on the unlabelled inputs (without knowing the truth labels). Our evaluation compares their methods with ours on various experimental settings covering different factors like proportion of labelled inputs, poison budget, dataset and used learning algorithms.

3 PROBLEM FORMULATION

3.1 SSL and Performance Metrics

We consider a classification problem where $C = \{c_1, \dots, c_k\}$ is the set of classes. Let $X = (X_L \cup X_U) \subseteq \mathbb{R}^D$ be a set of inputs represented

by D -dimensional feature vectors, such that X is divided into a set X_L of l labelled inputs and a set X_U of u unlabelled inputs. We arbitrarily index those inputs such that the labelled ones are placed first. That is, $X_L = \{x_1, \dots, x_l\}$ and $X_U = \{x_{l+1}, \dots, x_{l+u}\}$. Finally, we denote by $\mathbf{y}_L = (y_1, \dots, y_l) \in C^l$ the label vector associated to x_1, \dots, x_l , respectively. Then, the goal of SSL is to infer the label vector $\mathbf{y}_U = (y_{l+1}, \dots, y_{l+u})$ of X_U given X and \mathbf{y}_L , that is, to learn a function $f(x \in X_U | X, \mathbf{y}_L) \in C$ whose output should be as close as possible to the correct label vector $\mathbf{y}_U^* = (y_{l+1}^*, \dots, y_{l+u}^*)$.

One can measure the performance of SSL according to two different goals. If the goal is transductive learning (i.e., to infer the labels \mathbf{y}_U correctly), then performance can be measured by the *transductive accuracy* of the SSL algorithm, defined as

$$\frac{\sum_{x_i \in X_U} \mathbf{1}_{\{f(x_i) = y_i^*\}}}{|X_U|}$$

where $\mathbf{1}_{\{a=b\}} = 1$ if $a = b$ and 0 otherwise.

The second goal is *inductive learning*. It aims to build an arbitrary supervised model \mathcal{M} that generalizes to some unseen data X_G . In such a case, \mathcal{M} uses as a training set both the labelled inputs with their known labels and the unlabelled inputs with their labels inferred by SSL. Performance can then be measured by the *inductive accuracy*, i.e., the percentage of unseen data that \mathcal{M} classifies correctly. In practice, inductive accuracy is approximated by measuring the accuracy of \mathcal{M} on a carefully collected representative subset $X_{test} \subset X_G$ of the unseen data, named the *test set*.

3.2 Threat Model

We define an attack as a process that aims to reduce the performance of SSL (transductive or inductive accuracy) by poisoning available data. More precisely, we consider label poisoning attacks where the attack can alter the known labels \mathbf{y}_L .

Attack goal. The attack aims to maximize the error rate of the classifier. We target both transductive learning (where the end goal of the system is to correctly infer the label of unlabelled data) and inductive learning (where known and inferred labels are used to train a supervised model). Thus, the achievements of the attack can be measured by the error rate (i.e., 100% minus the accuracy) of the SSL algorithm or the supervised model, respectively, after poisoning. For simplicity, we restrict the scope to binary classification. The attack is indiscriminate as it does not target specific inputs and aims for generic errors (it does not target particular classes).

Attack knowledge. The attack has access to the known dataset (X and \mathbf{y}_L^*). However, our attack is black-box as it has knowledge about neither the graph-based SSL algorithm used during transduction nor the supervised models used for induction.

Attack capability. The attack can alter any label of the labelled inputs and can do so for a limited number k of labels (to minimize its actions). The features of the inputs cannot be altered.

Objective function. In summary, the objective of the attack is to find a perturbation vector $\delta^* \in (C \cup \{0\})^l$ under the constraint $|\{j | \delta_j^* \neq 0\}| = k$, such that altering \mathbf{y}_L into $\mathbf{y}_L \oplus \delta^*$ (where $(\mathbf{y}_L \oplus \delta^*)_j = (\mathbf{y}_L)_j$ if $\delta_j^* = 0$; and $(\mathbf{y}_L \oplus \delta^*)_j = \delta_j^*$ otherwise) results in the maximal (transductive or inductive) error rate.

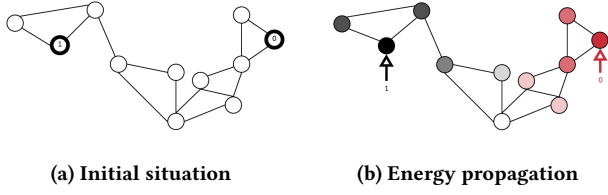


Figure 2: Energy propagation principle used by label propagation to infer the labels of the unlabelled inputs (light circles) from labelled inputs (tick circles with class number inside). Each colour represents a distinct class, while color gradient represents the influence strength of the propagated class label. This influence diminishes as there is more distance between the original labelled input and the unlabelled input.

4 INFLUENCE-DRIVEN DATA POISONING

4.1 Preliminaries

Graph-based SSL is a popular family of SSL algorithms that were shown to reach high accuracy at affordable computation costs [24]. They consist of building a fully connected graph where vertices are (labelled and unlabelled) inputs and where edges are weighted in proportion to the similarity between inputs. Then, SSL uses this representation to infer the label of any unlabelled input based on the weight of its edges.

Label propagation [32, 33] is one of the most popular graph-based SSL algorithms. As illustrated in Figure 2, it utilizes the *energy propagation principle* to achieve transductive learning. Intuitively, any labelled input emits its label towards any unlabelled input to influence its labelling, such that this influence (the “label energy”) fades away as the distance between the two inputs increases. Then, the label of any unlabelled input is inferred from the sum of the influences it receives (both from known labels and other inferred labels). Inferred labels also propagate their own energy. Thus, label propagation is an iterative process that stops after inferred labels have reached a steady state.

Hence, label propagation starts by computing a similarity distance between any pair of inputs. A popular class of functions to measure similarity between feature vectors are the Kernel functions. In particular, label propagation often uses the *Radial Based Function kernel* (RBF kernel) [23]. Accordingly, the similarity w_{ij} between x_i and x_j is given by $w_{ij} = \exp(-\gamma \|x_i - x_j\|^2)$ with $\gamma = (2\sigma^2)^{-1}$. For convenience, we define the $(l+u) \times (l+u)$ weight matrix $W = \begin{bmatrix} W_{LL} & W_{LU} \\ W_{UL} & W_{UU} \end{bmatrix}$ such that $W_{ij} = w_{ij}$. The unknown labels are then predicted via the energy minimization principle:

$$y = \operatorname{argmin}_{\hat{y}} \sum_{i,j} W_{ij} (\hat{y}_i - \hat{y}_j)^2 = \hat{y}^T (\mathbf{D} - W) \hat{y}, \text{ s.t. } \hat{y}_L = y_L$$

where $\mathbf{D} = \operatorname{diag}_{\sum_k W_{ik}}$. This problem admits a closed form solution:

$$y_U = (D_{[(l+1):(l+u), (l+1):(l+u)]} - W_{UU})^{-1} W_{UL} y_L.$$

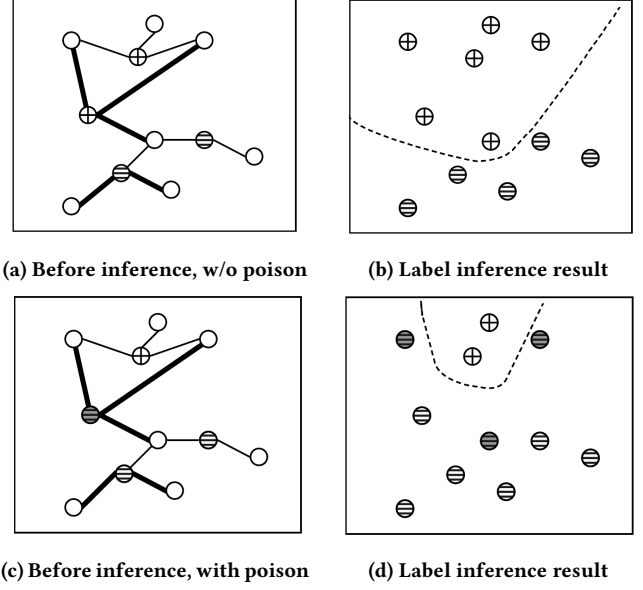


Figure 3: Principles of our influence-based poisoning attack.

Let $T = \begin{bmatrix} T_{LL} & T_{LU} \\ T_{UL} & T_{UU} \end{bmatrix}$ be the $(l+u) \times (l+u)$ label transition matrix defined as

$$T_{ij} = \begin{cases} 1, & i \leq l \wedge i = j; \\ 0, & i \leq l \wedge i \neq j; \\ \frac{w_{ij}}{\sum_{k=1}^{l+u} w_{kj}}, & l < i < l+u. \end{cases}$$

and \bar{T} the matrix obtained by row-normalizing T . Then the above solution can be rewritten as:

$$y_U = (I - \bar{T}_{UU})^{-1} \bar{T}_{UL} y_L$$

which can be computed as the fixed point of an iterative algorithm:

$$\begin{aligned} (1) & Y \leftarrow \bar{T} \cdot Y \\ (2) & Y \leftarrow \begin{bmatrix} Y_L^0 & Y_U \end{bmatrix}^T \end{aligned}$$

where Y_L^0 denotes the initial submatrix of Y corresponding to the known (clamped) labels.

4.2 Influence-driven Poisoning Attack

We show in Figure 3 the principles of our influence-based poisoning attack. The left figure shows the influence graph between inputs (both labelled – filled – and unlabelled – empty – inputs), where a thicker edge indicates a major influence (no edge means infinitesimal influence). Right figures show the result of label inference. The key idea of our approach is to identify the labelled inputs that have a major influence over a maximum number of unlabelled inputs, such that poisoning their label maximizes the number of wrong inferences. Here, the greyed input was the most influential and selected as the target to poison. This resulted in three labels inferred incorrectly.

To build our method, we associate matrix \bar{T} with a probabilistic interpretation: \bar{T}_{ij} is the estimated probability that x_i takes its label from x_j . Remember that \bar{T}_{ij} takes into account both the spread of all inputs’ label energy and their competitive influence onto

determining the inferred labels. In other words, \bar{T}_{ij} measures the relative *influence* of x_j onto x_i (compared to the other inputs). By restricting j to labelled inputs, we obtain the direct influence of any labelled input x_j onto any unlabelled input x_i .

Definition 1. Let $x_i \in X_U$ and $x_j \in X_L$. The *direct influence* of j onto i , noted e_{ij} , is given by $e_{ij} = \bar{T}_{ij}$.

The above influence metric disregards any indirect influence of x_j onto x_i , i.e. the energy received by x_i which transits from x_j to any intermediary input x_k before reaching x_i . This means that e_{ij} is only an approximation of the total influence of x_j onto x_i , which would be obtained by repeatedly multiplying \bar{T} by itself until it converges, that is by computing $\lim_{n \rightarrow \infty} \bar{T}^n$. Thus, the advantage of our direct influence metric is to avoid those costly computations while remaining a good approximation of the total influence (as later revealed by our evaluation results in Section 6).

Leaning on the direct influence metric, we are able to identify which labelled input has the highest direct influence onto a given unlabelled input.

Definition 2. The *most influential* input of $x_i \in X_U$ is given by

$$s_i^* = \operatorname{argmax}_{x_j \in X_L} e_{ij}$$

Thus, s_i^* is the labelled input such that altering its label has the highest likelihood (amongst all labelled inputs) to change the inferred label of x_i . In particular, if s_i^* holds more than half of the influence towards i , then poisoning its labels guarantees that the label of i is inferred incorrectly.¹

Definition 3. Let $x_i \in X_U$ and s_i^* be the most influential input of x_i . Then s_i^* is a *major influencer* of x_i iff

$$\frac{e_{is_i^*}}{\sum_{j \in X_L} e_{ij}} > 0.5$$

By generalizing this concept to the whole set of unlabelled inputs X_U , we can compute the number of unlabelled inputs for which a given labelled input is a major influencer.

Definition 4. Let X_U be the set of unlabelled inputs and $x_j \in X_L$ be a labelled input. The *Major Influence Range* (MIR) of x_j w.r.t. X_U is given by

$$\operatorname{MIR}(x_j) = |\{x_i \in X_U \text{ s.t. } \frac{e_{ij}}{\sum_{x_k \in X_L} e_{ik}} > 0.5\}|$$

Thus, $\operatorname{MIR}(x_j)$ captures the number of unlabelled inputs that would receive an incorrect label if the label of x_j was poisoned. This paves the way to design our *influence-driven data poisoning attack*: given a budget of k labels to poison, our attack alters the label of the k labelled inputs with the highest MIR.

Overall, our approach only requires computing the matrix \bar{T} and extracting the MIR value of each labelled input from it. Thus, its worst-case time complexity is $O((l+u)^2)$.

¹Actually, this may not be the case if indirect influences reduce the relative influence of s_i^* onto x_i . The key idea of our approach is to heuristically ignore the indirect influence between inputs to reduce the computation time of the influence metrics. Hence, it inherently remains an approximation.

5 EXPERIMENTAL PROTOCOL

5.1 Methodology and Research Questions

Our data poisoning method relies on the assumption that the MIR of any input captures well the number of incorrectly inferred labels that would result from poisoning that input. Accordingly, our first step is to validate this hypothesis. This is important to ensure that poisoning inputs ranked higher (according to their MIR) yields a higher number of incorrectly inferred labels. We ask:

RQ1 *Can the MIR metric identify the inputs whose poisoning yields the highest error rate?*

To answer this question, we measure the correlation, over the whole labelled input set, between the MIR of any input and the error rate of transductive learning after altering the label of that input. Thus, given l labelled inputs, we apply label propagation l times (once per poisoned labelled input).

To measure the correlation, we use the Kendall coefficient because, as an ordinal association metric, it focuses on how well MIR ranks the most impactful inputs first (irrespective of the actual error rates). Thus, if one has to select a limited number of inputs to poison, one should select the inputs of higher ranking. To complement our analysis we also use the Pearson's correlation coefficient, which captures linear relationships between two variables (here, MIR and the resulting error rate). Thus, a high Pearson correlation would indicate that MIR can also estimate the relative impacts of poisoning those inputs (taking into account the induced error rate).

As the MIR metric calculation lies on the same principle used by the label propagation algorithm, we expect that it performs well on this particular SSL algorithm. To assess the generalization potential of our approach, we repeat the experiments using a second SSL algorithm, label spreading, which uses a different Laplacian matrix to represent the input graph. Again, we report the resulting correlation coefficients. Having the same level of correlation would indicate that our method to select samples to poison transfers well to SSL algorithm that uses a different learning procedure.

Our next question assesses the practical effects that MIR-based data poisoning can have on inductive learning, that is, when supervised classification models are trained with both (known) labelled inputs and inputs whose labelled was inferred by semi-supervised learning. Thus, we ask:

RQ2 *How effective is MIR-based label poisoning in increasing the error rate of inductive learning?*

We consider the use case where an attacker aims to achieve the maximal effect (increase of error rate) within a limited budget of input labels to poison. Thus, given such a budget k , we study by how much the error rate of the supervised model on an independent test set increases after poisoning the top- k inputs ranked by MIR. Starting from the set of labelled inputs, we poison k labels before running an SSL algorithm (label propagation or label spreading) to infer the labels of the unlabelled data. Then, using the whole training set (labelled data and unlabelled data with inferred labels), we train a supervised model (defined independently of our poisoning attack) and compute its error rate on unseen test data (which have no intersection with the training data).

Following this, we compare our method to the state-of-the-art data poisoning attack for SSL, i.e. the two methods proposed by Liu

Name	$ X $	$ X_{test} $	# features
MNIST (1,7)	13,007	2,163	784
CIFAR-10 (cat, ship)	10,000	2,000	3072
rcv1	20,242	677,399	47,236

Table 1: Characteristics of the datasets we use in our experiments. X is the set of (labelled and unlabelled) inputs used during transductive learning. X_{test} is the independent set of unseen data used for testing the inductive learning.

et al. [17], both in terms of effectiveness (how much it increases error rate) and efficiency (computation time). We ask:

RQ3 *How does our approach compare to the state of the art data poisoning attack?*

We compare the effectiveness of Liu et al.’s method and ours on both transductive learning and inductive learning. In each case, we record the error rates that result from poisoning the inputs suggested by each method. To compare both methods on a fair ground, we use label propagation to perform transductive learning, since both methods lean on the mathematical properties of this SSL algorithm.

Having shown that our method constitutes an effective attack, we aim to determine if we can use the same principle of influence to guide engineers in setting up effective countermeasures. Thus, we ask:

RQ4 *Can MIR drive the design of countermeasures to reduce the effect of label poisoning?*

To answer this question, we simulate a scenario where engineers *relabel* the most influential inputs before transductive learning. Relabelling is indeed a common countermeasure to label poisoning [19]. We compare this with the alternative countermeasure of labelling additional inputs to offset the poisoning effect. Such comparison will reveal the most effective allocation of the engineers’ effort to reduce the effects of label poisoning.

Here it must be noted that the above are general settings common to all RQs we investigate. Specific and detailed settings required to answer each RQ are given at the beginning of the dedicated sections answering them.

5.2 Test Subjects

We run our experiments on three datasets involving two image classification problems (MNIST [14] and CIFAR-10 [2]) and one text classification problem (rcv1 [15]). MNIST and CIFAR-10 are widely used in research and considered as a good baseline to observe key trends, in addition to requiring affordable computation cost. Using both brings variety within the same application domain: MNIST is about black-and-white images of handwritten digits while CIFAR-10 includes colour images of animals and vehicles. Due to this difference, the different class manifolds in MNIST are closer in the feature space than in CIFAR-10. The last dataset, rcv1, allows observing results in another application domain and on a more challenging dataset (computation wise). It involves inputs with 47,236 features, which drastically increases the time to run our experiments.

It is to be noted that these three datasets involve 10 classes (MNIST and CIFAR-10) and 103 classes (rcv1). For simplicity we restrict the scope of our empirical study to binary classification. We use binary variants of the datasets. For MNIST, we keep the images showing digits 1 and 7, following the protocol of [17]. For CIFAR-10, we keep the classes ship and cat. For rcv1, we directly reuse the binary variant of the dataset independently proposed in the LIBSVM collection of datasets [6]. The characteristics of the resulting datasets are summarized in Table 1.

To perform transductive learning, we consider label propagation and label spreading with the default parameters of their scikit-learn implementation. Both label propagation and label spreading use an RBF kernel, a convergence threshold of 10^{-3} and a maximum of 30 iterations. Label spreading also uses a clamping factor of 0.1.

When our experiments involve inductive learning, we use two supervised models: a random forest and a multilayer perceptron. The random forest is composed of 100 trees, uses Gini impurity as the quality split criterion and considers up to $\sqrt{\#features}$ when looking for the best split. The Multilayer Perceptron is composed of an input layer with $\#features$ neurons, a hidden layer with 128 neurons activated by a rectifier function and an output layer with one neuron using a sigmoid activation function.

We chose those two supervised learning methods as their learning process significantly differs from that of graph-based SSL. This allows us to observe how well poisoning transfers to models based on different learning principles.

5.3 Implementation and hardware

We implemented our methods in a prototype tool on top of Python 3.7.0. Our tool is open source and publicly available together with our datasets and results.² As for label propagation and label spreading, we rely on their implementation in the open-source library *scikit-learn*³ [20]. The implementation of the supervised models we use (viz. random forest) is also based on *scikit-learn*.

To compare with the state-of-the-art method [17], we reuse the original implementation provided by the authors.⁴ It includes two variants of the method: one is deterministic and the other is stochastic.

All experiments were run on Google Cloud using a virtual machine with 12 VCPU Intel Xeon Skylake 2.0 GHz and 45GB of RAM. Running once all experiments of all RQs required approximately 15 days of computation.

6 DETAILED SETUP AND RESULTS

6.1 RQ1: Correlation with Transductive Learning Error Rate

6.1.1 Detailed Setup. We measure the Kendall and Pearson coefficients between the MIR of any labelled input x and the error rate of transductive learning (% of unlabelled inputs incorrectly inferred by the SSL algorithm) after poisoning the label of x only. These coefficients take their value between -1 and +1 (negative and positive correlations). Coefficient values greater than 0.5 demonstrate

²https://anonymous.4open.science/r/adversarial_label_propagation-FC83

³<https://github.com/scikit-learn/scikit-learn>

⁴<https://github.com/xuanqing94/AdvSSL>

MNIST	$l = 650$ (5%)		$l = 1950$ (15%)		$l = 3250$ (25%)	
L. propagation	0.96	0.99	0.95	0.97	0.95	0.98
L. spreading	0.62	0.83	0.63	0.80	0.59	0.71
CIFAR-10	$l = 500$ (5%)		$l = 1500$ (15%)		$l = 2500$ (25%)	
L. propagation	0.42	0.52	0.43	0.49	0.41	0.47
L. spreading	0.35	0.35	0.31	0.32	0.30	0.33
RCV1	$l = 1012$ (5%)		$l = 1036$ (15%)		$l = 5060$ (25%)	
L. propagation	0.47	0.73	0.51	0.81	0.56	0.82
L. spreading	0.42	0.53	0.43	0.53	0.43	0.57

Table 2: Correlation, across all labelled inputs, between the MIR metric of any input and the error rate of transductive learning after altering the input label. Left numbers are Kendall’s τ coefficients while right numbers are Pearson’s.

a strong correlation; between 0.3 and 0.49 correspond to medium correlations; less than 0.29 are interpreted as weak correlations.

To perform transductive learning, we consider both label propagation and label spreading. This allows us to observe how well our methods transfers to another SSL algorithm computing input influences differently.

We measure the correlation for all datasets, with different proportions of labelled/unlabelled inputs, i.e., with 5%, 15% and 25% of labelled inputs (see Table 2 for exact numbers for each dataset). We start from 5% because this is the smallest percentage where the SSL algorithms yield acceptable accuracy values (above 80%) for all three datasets. At the opposite end, 25% is aligned with current research on SSL (see, e.g. [1]). Considering different proportions of labelled inputs allows us to observe whether the MIR metric keeps the same effectiveness when an increasingly smaller set of unlabelled inputs are inferred from an increasingly larger set of labelled inputs. For each proportion, we randomly select which inputs are labelled. To account for random variations resulting from this split, we repeat each experiment 3 to 10 times (depending on the computation cost involved) and report the average of the correlation coefficients.

6.1.2 Results. General observations. As observed in Table 2, the Kendall and Pearson coefficients show medium to strong correlations in all experiments (Kendall coefficients are between 0.30 and 0.96, while Pearson coefficients are between 0.32 and 0.99). All reported correlation coefficients are associated with a p-value less than 10^{-7} , meaning that we can reject the null hypothesis of an overall absence of correlation with a type-1 error $< 1\%$. Overall, this means that *our MIR metric captures well the estimated impact of poisoning any labelled input*. Thus, MIR has the potential to prioritize the inputs that an attacker should target to maximize error rate (which we assess in RQ2 and RQ3). Moreover, the high Pearson coefficients demonstrate that MIR also enables a relative comparison of the input impact (i.e., a doubled MIR indicates a doubled increase in error rate).

Dataset. The correlation values vary much across the datasets. This can be explained by the fact that SSL heavily relies on the relative position of the class clusters during label inference. Thus, label poisoning is more efficient as the clusters are closer (the

influence of the negative class is stronger on the positive inputs). For instance, MNIST inputs (white digits on dark background) are closer in the feature space than CIFAR-10 inputs (coloured images with arbitrary background), which explains that the correlations in MNIST are stronger than in CIFAR-10.

Transductive learning algorithm. Stronger correlations are observed when transductive learning is performed by label propagation. This was expected since our approach follows the same energy propagation principle as label propagation does. This reveals that the direct influence of any labelled input on any unlabelled input is a good approximation of its total influence (including indirect propagation through intermediary inputs). When label spreading is used instead, the correlations are slightly lower but remains within the same range of correlation strength. This indicates that *the MIR metric transfers well over different graph-based SSL algorithms*. For engineers, this means that relying on one algorithm over the other increases poisoning resilience only slightly.

Proportion of labelled inputs. The relative proportion of labelled inputs (compared to unlabelled inputs) does not significantly affect the correlations. Interestingly, this means that *MIR remains a strong metric even when there are more labelled inputs competing to influence unlabelled inputs*. In other words, (manually) labelling more inputs does not hinder the capability of MIR to select the most impactful inputs.

6.2 RQ2: Impact on Inductive Error Rate

6.2.1 Detailed setup. We measure the effectiveness of our poisoning method by computing the error rate of inductive learning (using a Random Forest – RF – and a Multi-Layer Perceptron – MLP). We label 25% of inputs.⁵ Then, we poison the k labelled inputs with the highest MIR, for k ranging from 5% to 20% of the labelled inputs (by steps of 5%), and we measure the error rate. Compared to RQ1, these experiments allow us to measure (a) the effectiveness of our method when selecting a set of inputs (rather than a single one) and (b) how much the attack transfers during induction.

To perform the inductive learning, we first apply an SSL algorithm (i.e., each amongst label propagation and label spreading) to infer the labels of the unlabelled input set. Then, we train a supervised model (i.e., an RF or an MLP) using both labelled inputs (with their known labels) and the unlabelled inputs (with their inferred labels). This yields, per dataset, a total of four combinations (2 SSL algorithms \times 2 supervised models).

We show the inductive error rate for the four combinations, on all datasets, with a fixed proportion of labelled/unlabelled inputs (i.e., 25% of labelled inputs in the training set), and with different budget k of those labelled inputs that our method can poison (5%, 10%, 15% and 20% of the labelled inputs). For instance, with 25% of labelled inputs in CIFAR-10, a 5% poison budget means that our method has poisoned the 250 labelled inputs with the highest MIR out of the 5000 labelled inputs.

To account for random variations when splitting the training set into labelled and unlabelled inputs, we repeat each experiment five times and report the average of the obtained inductive error rates.

⁵Because the proportion of labelled inputs does not affect our conclusions and because of lack of space, we do not report here the results for 5% and 15%. For the sake of reproducibility, the results for 5% and 15% of labelled inputs can be computed with our replication package.

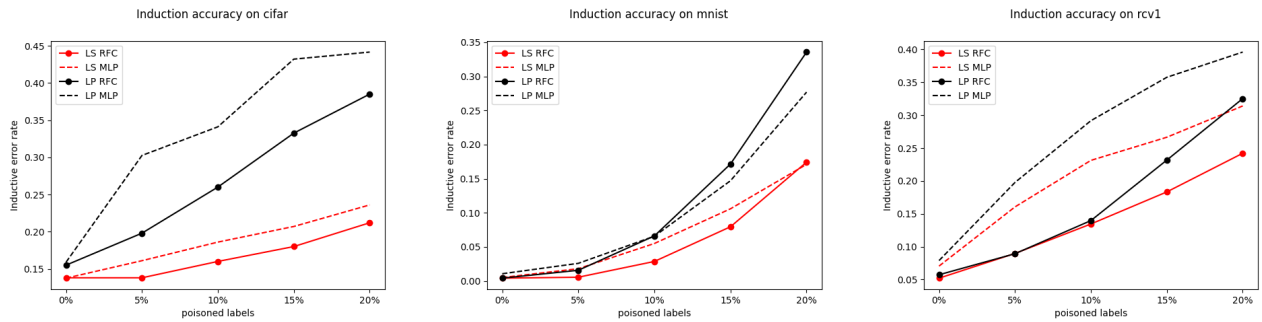


Figure 4: Error rate (Y axis) achieved by our data poisoning attack wrt. an increasing poison budget (X axis, expressed in % of the labelled inputs), measured on random forest classifier (RFC plain curves) and multilayer perceptron (MLP dotted curves) trained on both known labels and labels inferred with label propagation (LP black curves) or label spreading (LS red curves).

6.2.2 Results. General observations. As shown in Figure 4, our poisoning attack significantly increases the error rate of the supervised models, regardless of the SSL algorithm used for transductive learning. For example, when poisoning 20% of the labelled inputs, the total increase in error rate ranges from 16% (MNIST, label spreading, RF) to 45% (CIFAR-10, label propagation, MLP). These results indicate that (1) *the MIR metric remains effective when multiple inputs are poisoned* and (2) *the induced errors transfer to inductive learning*. This second point demonstrates that introducing a different learning algorithm in the induction process does not impede the altered labels to “poison the well”.

SSL algorithm. The choice of the SSL algorithm to infer the labels of the unlabelled inputs has a major effect on the transferability of MIR-based data poisoning. Indeed, using label spreading alleviates not only the total error rate but also the marginal effect of increasing the poison budget. For instance, we observe in left figure of Figure 4 (which concerns CIFAR-10) that the error rate increase for both supervised models is twice steeper when using label propagation than when using label spreading. This trend is also observed on MNIST and RCV1, where the error rate curve is steeper in the case of label propagation. These results shed new light on our findings from RQ1: *while the choice of the SSL algorithm does not change the impact of the first poisoned inputs on transductive learning, it allows increasing the resilience of inductive learning when more inputs are poisoned*.

6.3 RQ3: Comparison with the State of the Art

6.3.1 Detailed setup. We compare the effectiveness our method to the state of the art (i.e., the greedy method and the probabilistic method of Liu et al. [17]), in terms of effectiveness (error rate of transductive learning and inductive learning) and efficiency (computing time required to select the inputs to poison).

To compare efficiency, we measured the time required by each method to select a subset of 20 labelled inputs to poison, on all datasets. For our method, this involves computing the MIR value of all labelled inputs and rank them accordingly. For the greedy method, this means performing 20 iterations of their greedy search algorithm. For the probabilistic method it means estimating the poisoning probability for each labelled samples based on their objective function and returning the most probable 20 samples. In

	MNIST		RCV1		CIFAR-10	
	μ	σ	μ	σ	μ	σ
Greedy	86.47	4.58	711.25	77.69	479.02	11.05
Probabilistic	29.65	0.80	109.51	0.82	17.52	0.71
MIR	7.67	0.19	55.45	0.38	6.44	0.75

Table 3: Efficiency of our method (MIR) compared to the state of the art (greedy and probabilistic). μ and σ are the median and the standard deviation of the runtime (in s.).

case the poisoning budget exceeds 20 inputs, the greedy method takes proportionally more time, whereas the computing time of our method and of the probabilistic one are not affected by the number of poisoned inputs. To account for random variations, we repeat the experiments 50 times and report the median and the standard deviation.

To compare effectiveness, we apply the same experimental protocol as RQ2, with the addition that we also compute the error rate of label propagation during transductive learning.

6.3.2 Results. Efficiency. The median and standard deviation of each method’s runtime on each dataset is given in Table 3. We observe that our method significantly outperforms the other two in every case, being 2 to 3 times faster than the probabilistic method and an order of magnitude faster than the greedy method, while having a smaller standard deviation.

Effectiveness. The error rates achieved by the three methods are shown in Tables 4. When applied to transductive learning, our MIR-based method achieves a higher error rate than the other two for every dataset and poison budget. The difference increases as more inputs are poisoned. This indicates that our method exploits additional poison budget better than the other methods.

When applied to inductive learning, the general trends remain: our method generally outperforms the others and it does so better as the poison budget increases. For MNIST and CIFAR-10, the difference is substantial: our method achieves an error rate 1.22 to 30.8 times higher.

As for RCV1, the three methods perform comparatively well, though the probabilistic method works best when a RF is used for

induction. A Wilcoxon signed rank test reveals that the difference between our method and the probabilistic one (over all RCV1 results) is not statistically significant. Nonetheless, by extrapolating from the efficiency results (Table 3), we estimate that MIR would run 17 times faster on RCV1 than the deterministic method already with a poison budget of 5% (even more as the budget increases).

Taken together, all our results show that our method results in a higher error rate than the state of the art, with statistical significance (the Wilcoxon test over all results rejects the null hypothesis that the two methods perform equally, with a p-value less than 10^{-11}). Thus, *our influence-driven attack outperforms the state of the art in both effectiveness and efficiency and forms a new baseline for SSL poisoning.*

6.4 RQ4: Countermeasures

6.4.1 Detailed setup. We consider again the cases of transductive learning and inductive learning, where 5%, 15%, and 25% of the inputs are labelled. We apply our influence-driven attack with a poison budget totalling 10% of the labelled inputs. Then, we investigate how much the two countermeasures (relabelling inputs with the highest MIR vs. labelling additional (previously unlabelled) inputs). We allocate the same effort to the two countermeasures, i.e., one third of the number of poisoned labels. We repeat the experiments five times and report the average error rates.

6.4.2 Results. Table 5 shows the results. We observe that relabelling the most influential labels systematically achieves a higher reduction in error rate. Interestingly, on average, relabelling one third of the poisoned inputs halves the error rate induced by the poisoning. A Wilcoxon test rejects the null hypothesis that the two methods perform equally (p-value $< 10^{-5}$). Labelling more input actually offers small reductions in error rate ($< 15\%$ on average). Overall, this indicates that *to alleviate the poisoning effects, engineers should focus their effort on relabelling influential inputs.*

6.5 Threats to Validity

Threats to internal validity concern the implementation of the software artefacts used in our study. Some are addressed by the fact that we reuse established implementations of the learning algorithms with typical parameters. The resulting (non-poisoned) models yield a small error rate on state-of-the-art datasets used as is (including their splitting into training and test sets), even when small proportions of labelled inputs are used. This indicates that our setup was appropriate.

The implementation of our approach was tested manually and through various experiments, which provides some confidence regarding its correctness. Moreover, we reused the available implementation of the greedy and the probabilistic method and as provided by its inventors.

The threats to external validity originate from the number of learning algorithms and datasets we used in our experiments. MNIST and CIFAR-10 are established in the ML literature, whereas RCV1 was already used in SSL-related studies [17]. The reduction to binary classification problems has followed the protocol of previous research [6, 17].

The randomness induced by splitting the training set between labelled inputs and unlabelled inputs is another factor affecting

our results. To mitigate its effects, we repeated our experiments multiple times (between 3 and 10, depending on required the computation time) and manually checked the absence of significant variations. Still, it remains possible that additional runs produce outliers (especially for RCV1, the largest dataset of the three).

Nevertheless, in general, it is likely that the effectiveness of our approach varies upon different external factors such as the used datasets (as our experiments already witness). Still, the core principle of measuring input influence is universally applicable and we believe that our key conclusions shall remain valid. Only additional experimentations can alleviate this risk, though. Fortunately, our open-source implementation and the black-box nature of our influence metric facilitate the replication and complementation of our study.

Finally, threats to construct validity come from the factors we measure to draw our conclusions. We studied the correlation between MIR (our metric) and the error rate, which is a natural metric to use. The error rate was also used by previous research [17] to measure the effectiveness of data poisoning.

7 CONCLUSION

We proposed a new label poisoning attack that targets the most influential inputs during semi-supervised learning. Our extensive experiments show that our approach outperforms the state of the art (in effectiveness and efficiency) under various settings and forms a new baseline for future research.

The influence metric it leans on has the advantage of being simple and fast to compute. In addition to allowing the design of efficient and effective attacks, it also provides practical benefits for engineers. Indeed, using the same simple metric, engineers can identify the most critical inputs to investigate in case there is doubt about data integrity (e.g., when data come from untrustworthy sources). This way, they can reduce the risk and effects of poisoning by relabelling influential inputs and, overall, alleviate the data quality threats.

ACKNOWLEDGEMENT

This work is supported by the Luxembourg National Research Funds (FNR) through CORE project C18/IS/12669767/STELLAR/LeTraon.

REFERENCES

- [1] . 2021. . Retrieved August 25, 2020 from <https://paperswithcode.com/sota/semi-supervised-image-classification-on-cifar-2>
- [2] Alex Krizhevsky and Vinod Nair and Geoffrey Hinton. 2009. *The CIFAR-10 dataset*. Retrieved August 25, 2020 from <https://www.cs.toronto.edu/~kriz/cifar.html>
- [3] Shai Ben-David, Tyler Lu, and Dávid Pál. 2008. Does Unlabeled Data Provably Help? Worst-case Analysis of the Sample Complexity of Semi-Supervised Learning. In *COLT*. 33–44.
- [4] Battista Biggio, Luca Didaci, Giorgio Fumera, and Fabio Roli. 2013. Poisoning attacks to compromise face templates. In *International Conference on Biometrics, ICB 2013, 4-7 June, 2013, Madrid, Spain*, Julian Fierrez, Ajay Kumar, Mayank Vatsa, Raymond N. J. Veldhuis, and Javier Ortega-Garcia (Eds.). IEEE, 1–7. <https://doi.org/10.1109/ICB.2013.6613006>
- [5] Battista Biggio, Blaine Nelson, and Pavel Laskov. 2012. Poisoning Attacks against Support Vector Machines. arXiv:1206.6389 [cs.LG]
- [6] Chih-Chung Chang and Chih-Jen Lin. 2011. *LIBSVM*. Retrieved August 25, 2020 from <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>
- [7] Maxime Cordy, Steve Muller, Mike Papadakis, and Yves Le Traon. 2019. Search-based test and improvement of machine-learning-based anomaly detection systems. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2019, Beijing, China, July 15-19, 2019*, Dongmei Zhang and Anders Møller (Eds.). ACM, 158–168. <https://doi.org/10.1145/3293882.3330580>

	clean	5% labels poisoned			10% labels poisoned			15% labels poisoned			20% labels poisoned		
Trans.		MIR	Greed.	Prob.	MIR	Greed.	Prob.	MIR	Greed.	Prob.	MIR	Greed.	Prob.
MNIST													
5% labels	1.25%	16.22%	3.04%	5.41%	27.05%	3.04%	9.69%	35.86%	3.04%	15.13%	43.98%	3.04%	19.46%
15% labels	0.69%	16.09%	1.27%	5.67%	27.48%	1.27%	9.84%	37.07%	1.27%	14.71%	45.71%	1.27%	20.17%
25% labels	0.58%	17.47%	1.16%	5.68%	29.94%	1.16%	9.86%	40.26%	1.16%	14.82%	48.88%	1.16%	19.26%
CIFAR-10													
5% labels	19.12%	33.62%	26.15%	21.90%	53.08%	28.90%	25.20%	61.67%	28.90%	25.00%	66.29%	28.90%	31.90%
15% labels	17.52%	32.03%	26.62%	18.70%	37.98%	26.62%	19.10%	45.00%	26.62%	26.01%	51.66%	26.62%	28.90%
25% labels	17.34%	29.21%	28.11%	19.40%	35.67%	28.99%	20.80%	40.21%	28.99%	19.20%	44.20%	28.99%	21.20%
RCV1													
5% labels	11.66%	22.22%	19.53%	19.66%	29.39%	28.62%	28.72%	35.17%	34.62%	35.02%	42.37%	38.77%	39.85%
15% labels	8.19%	22.00%	18.13%	18.69%	31.26%	27.09%	28.01%	38.89%	34.74%	35.28%	46.26%	39.10%	39.96%
25% labels	6.91%	23.12%	18.41%	19.31%	33.18%	27.98%	28.22%	41.32%	34.62%	35.15%	49.74%	39.57%	40.12%
Ind. R.F.		MIR	Greed.	Prob.	MIR	Greed.	Prob.	MIR	Greed.	Prob.	MIR	Greed.	Prob.
MNIST													
5% labels	1.07%	8.93%	1.25%	0.92%	19.94%	1.25%	1.34%	29.22%	1.25%	2.64%	40.92%	1.25%	4.35%
15% labels	0.56%	2.27%	1.56%	0.55%	10.31%	1.56%	0.55%	19.61%	1.56%	0.83%	34.17%	1.56%	1.66%
25% labels	0.42%	1.62%	0.51%	0.55%	6.85%	1.43%	0.42%	18.50%	1.43%	0.60%	33.71%	1.43%	1.57%
CIFAR-10													
5% labels	18.96%	22.36%	25.35%	17.60%	53.32%	28.81%	19.40%	68.81%	28.81%	20.15%	71.76%	28.81%	21.55%
15% labels	17.25%	22.10%	24.06%	18.55%	29.1%	24.04%	19.25%	43.32%	24.04%	20.65%	53.45%	24.04%	23.75%
25% labels	15.50%	20.55%	18.95%	15.10%	25.70%	20.19%	15.80%	29.60%	20.19%	17.85%	39.90%	20.19%	20.05%
RCV1													
5% labels	9.00%	11.44%	10.49%	13.79%	15.96%	16.45%	36.56%	23.09%	26.36%	48.64%	35.92%	29.36%	52.00%
15% labels	5.85%	9.77%	10.85%	13.28%	15.33%	14.22%	38.88%	26.84%	23.38%	51.16%	44.08%	23.38%	53.40%
25% labels	5.73%	8.92%	11.64%	11.68%	13.96%	20.88%	33.60%	23.17%	30.08%	50.40%	32.60%	30.08%	53.40%
Ind. MLP		MIR	Greed.	Prob.	MIR	Greed.	Prob.	MIR	Greed.	Prob.	MIR	Greed.	Prob.
MNIST													
5% labels	2.17%	10.17%	3.42%	2.31%	20.76%	3.42%	2.40%	33.15%	3.42%	5.41%	45.31%	3.42%	7.12%
15% labels	0.88%	4.53%	1.06%	1.20%	16.18%	1.06%	1.02%	25.98%	1.06%	1.39%	42.07%	1.06%	2.22%
25% labels	0.92%	4.58%	1.06%	0.69%	7.17%	1.06%	0.92%	14.33%	1.06%	1.06%	22.47%	1.06%	1.99%
CIFAR-10													
5% labels	19.00%	21.85%	21.80%	18.80%	64.30%	20.15%	29.70%	66.45%	20.15%	26.90%	69.80%	20.15%	23.90%
15% labels	16.85%	19.05%	31.20%	18.10%	24.25%	32.15%	18.60%	44.15%	32.15%	18.05%	51.30%	32.15%	20.25%
25% labels	16.50%	17.60%	15.95%	16.90%	19.65%	16.85%	17.45%	27.70%	16.85%	17.30%	40.50%	16.85%	19.65%
RCV1													
5% labels	11.68%	17.88%	21.16%	21.52%	26.56%	32.00%	31.72%	32.56%	41.08%	40.68%	39.68%	46.24%	46.92%
15% labels	7.40%	19.16%	18.92%	19.00%	29.08%	32.92%	30.04%	35.88%	40.76%	41.28%	40.48%	45.04%	46.56%
25% labels	6.72%	21.12%	18.96%	20.28%	29.48%	27.64%	27.44%	35.32%	33.68%	36.56%	40.80%	40.08%	42.64%

Table 4: RQ3: Error rate (in %) of transductive learning (with label propagation), inductive learning (with random forest) and inductive learning (with multi-layer perceptron) achieved by our MIR method (MIR), the state of the art deterministic (Greed) and stochastic (Prob) methods. Higher is better. Bold indicates the best method for each experiment.

- [8] Rob Fergus, Yair Weiss, and Antonio Torralba. 2009. Semi-Supervised Learning in Gigantic Image Collections. In *Advances in Neural Information Processing Systems 22*, Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta (Eds.). Curran Associates, Inc., 522–530. <http://papers.nips.cc/paper/3633-semi-supervised-learning-in-gigantic-image-collections.pdf>
- [9] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg. 2019. BadNets: Evaluating Backdoor Attacks on Deep Neural Networks. *IEEE Access* 7 (2019), 47230–47244.
- [10] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. 2019. Label Propagation for Deep Semi-Supervised Learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16–20, 2019*. Computer Vision Foundation / IEEE, 5070–5079. <https://doi.org/10.1109/CVPR.2019.00521>
- [11] Jane Wakefield. 2016. *Microsoft chatbot is taught to swear on Twitter*. Retrieved August 25, 2020 from <https://www.bbc.com/news/technology-35890188>
- [12] Anthony D. Joseph, Pavel Laskov, Fabio Roli, J. Doug Tygar, and Blaine Nelson. 2013. Machine Learning Methods for Computer Security (Dagstuhl Perspectives Workshop 12371). *Dagstuhl Manifestos* 3, 1 (2013), 1–30. <https://doi.org/10.4230/DagMan.3.1.1>
- [13] Matti Kääriäinen. 2005. Generalization error bounds using unlabeled data. In *International Conference on Computational Learning Theory*. Springer, 127–142.
- [14] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [15] David D. Lewis, Yiming Yang, Tony G. Rose, and Fan Li. 2004. RCV1: A New Benchmark Collection for Text Categorization Research. *J. Mach. Learn. Res.* 5 (2004), 361–397. <http://jmlr.org/papers/volume5/lewis04a/lewis04a.pdf>

MNIST	Transduction			Induction		
	MIR	LAB	None	MIR	LAB	None
5% labels	15.80	24.06	27.05	04.38	12.32	19.94
15% labels	16.38	25.64	27.48	01.56	06.02	10.31
25% labels	17.12	27.21	29.94	01.00	03.80	6.85

CIFAR-10	Transduction			Induction		
	MIR	LAB	None	MIR	LAB	None
5% labels	28.76	46.06	53.08	19.46	42.39	54.95
15% labels	23.17	36.35	37.98	19.47	43.44	54.11
25% labels	23.75	34.78	35.67	17.18	41.36	44.06

RCV1	Transduction			Induction		
	MIR	LAB	None	MIR	LAB	None
5% labels	21.64	28.00	29.39	10.56	13.92	15.96
15% labels	19.58	29.62	31.26	09.04	14.04	15.33
25% labels	20.31	31.07	33.18	08.20	12.24	13.96

Table 5: RQ4: Error rate (in %) achieved by double checking one third of the poisoned labels (MIR) compared with labelling the same number of additional inputs (LAB) and no countermeasure (None). Lower is better.

- [16] Yu-Feng Li and De-Ming Liang. 2019. Safe semi-supervised learning: a brief introduction. *Frontiers Comput. Sci.* 13, 4 (2019), 669–676. <https://doi.org/10.1007/s11704-019-8452-2>
- [17] Xuanqing Liu, Si Si, Xiaojin Zhu, Yang Li, and Cho-Jui Hsieh. 2019. A unified framework for data poisoning attack to graph-based semi-supervised learning. *arXiv preprint arXiv:1910.14147* (2019).
- [18] Shike Mei and Xiaojin Zhu. 2015. Using Machine Teaching to Identify Optimal Training-Set Attacks on Machine Learners. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA*, Blai Bonet and Sven Koenig (Eds.). AAAI Press, 2871–2877. <http://www.aaai.org/ocs/index.php/AAAI/AAAI15/paper/view/9472>
- [19] Andrea Paudice, Luis Muñoz-González, and Emil C. Lupu. 2018. Label Sanitization Against Label Flipping Poisoning Attacks. In *ECML PKDD 2018 Workshops - Nemesis 2018, UrbReas 2018, SoGood 2018, IWAISe 2018, and Green Data Mining 2018, Dublin, Ireland, September 10-14, 2018, Proceedings (Lecture Notes in Computer Science, Vol. 11329)*, Carlos Alzate, Anna Monreale, Haytham Assem, Albert Bifet, Teodora Sandra Buda, Bora Caglayan, Brett Drury, Eva Garcia-Martin, Ricard Gavaldà, Stefan Kramer, Niklas Lavesson, Michael Madden, Ian Molloy, Maria-Irina Nicolae, and Mathieu Sinn (Eds.). Springer, 5–15. https://doi.org/10.1007/978-3-030-13453-2_1
- [20] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
- [21] Roberto Perdisci, David Dagon, Wenke Lee, Prahlad Fogla, and Monirul I. Sharif. 2006. MisleadingWorm Signature Generators Using Deliberate Noise Injection. In *2006 IEEE Symposium on Security and Privacy (S&P 2006), 21-24 May 2006, Berkeley, California, USA*. IEEE Computer Society, 17–31. <https://doi.org/10.1109/SP.2006.26>
- [22] Nikolaos Pitropakis, Emmanouil Panaousis, Thanassis Giannetsos, Eleftherios Anastasiadis, and George Loukas. 2019. A taxonomy and survey of attacks against machine learning. *Computer Science Review* 34 (2019), 100199.
- [23] Bernhard Schölkopf, Alex J. Smola, Robert C. Williamson, and Peter L. Bartlett. 2000. New Support Vector Algorithms. *Neural Computation* 12, 5 (2000), 1207–1245. <https://doi.org/10.1162/089976600300015565>
- [24] Matthias Seeger. 2000. Learning with Labeled and Unlabeled Data. (2000). <http://infoscience.epfl.ch/record/161327>
- [25] Michael Speriosu, Nikita Sudan, Sid Upadhyay, and Jason Baldridge. 2011. Twitter Polarity Classification with Label Propagation over Lexical Links and the Follower Graph. In *Proceedings of the First workshop on Unsupervised Learning in NLP*. Association for Computational Linguistics, Edinburgh, Scotland, 53–63. <https://www.aclweb.org/anthology/W11-2207>
- [26] Gang Wang, Tianyi Wang, Haitao Zhang, and Ben Y. Zhao. 2014. Man vs. Machine: Practical Adversarial Detection of Malicious Crowdsourcing Workers. In *Proceedings of the 23rd USENIX Conference on Security Symposium (San Diego, CA) (SEC'14)*. USENIX Association, USA, 239–254.
- [27] Huang Xiao, Battista Biggio, Blaine Nelson, Han Xiao, Claudia Eckert, and Fabio Roli. 2015. Support vector machines under adversarial label contamination. *Neurocomputing* 160 (2015), 53–62. <https://doi.org/10.1016/j.neucom.2014.08.081>
- [28] Jie M. Zhang, Mark Harman, Lei Ma, and Yang Liu. 2019. Machine Learning Testing: Survey, Landscapes and Horizons. *CoRR* abs/1906.10742 (2019). [arXiv:1906.10742](http://arxiv.org/abs/1906.10742) <http://arxiv.org/abs/1906.10742>
- [29] Ping Zhang, Fei Wang, Jianying Hu, and Robert Sorrentino. 2015. Label Propagation Prediction of Drug-Drug Interactions Based on Clinical Side Effects. *Scientific Reports* 5 (2015).
- [30] Dengyong Zhou, Olivier Bousquet, Thomas N Lal, Jason Weston, and Bernhard Schölkopf. 2004. Learning with local and global consistency. In *Advances in neural information processing systems*. 321–328.
- [31] Xiaojin Zhu. 2005. *Semi-Supervised Learning Literature Survey*. Technical Report 1530. Computer Sciences, University of Wisconsin-Madison.
- [32] Xiaojin Zhu and Youbin Ghahramani. 2002. Learning from labeled and unlabeled data with Label propagation. (2002).
- [33] Xiaojin Zhu, Zoubin Ghahramani, and John Lafferty. 2003. Semi-Supervised Learning Using Gaussian Fields and Harmonic Functions. *ICML-03, 20th International Conference on Machine Learning* 3, 912–919.