
Online Strategic Classification with Noise and Partial Feedback

Tianrun Zhao, Xiaojie Mao*, Yong Liang

School of Economics and Management, Tsinghua University, Beijing, China, 100084
ztr23@mails.tsinghua.edu.cn, maobj@sem.tsinghua.edu.cn,
liangyong@sem.tsinghua.edu.cn,

Abstract

In this paper, we study an online strategic classification problem, where a principal aims to learn an accurate binary linear classifier from interactions with sequentially arriving agents. For each agent, the principal announces a classifier. The agent can strategically exercise costly manipulations on his features to be classified as the favorable positive class. The principal is unaware of the true feature-label relationship, but observes all reported features and only labels of positively classified agents. We assume that the true feature-label relationship is given by a halfspace model subject to arbitrary feature-dependent but bounded noise (i.e., Massart noise). This problem faces the combined challenges of agents' strategic feature manipulations, partial feedback observations, and label noise. We tackle these challenges by a novel learning algorithm. We show that the proposed algorithm yields classifiers that converge to the clairvoyant optimal classifier and attains a regret rate of $O(\sqrt{T})$ up to poly-logarithmic and constant factors over T cycles.

1 Introduction

Strategic classification studies the problem of learning robust classifiers in presence of self-interested strategic agents. When subjugated to decision-making aided by classification algorithms, agents may strategically modify their observable features to game the classification algorithms into making decisions that best serve the agents' goals. For example, a bank may use classification to determine whether loan applicants are qualified to grant approvals. The applicants prefer positive classification and loan approvals, so they have the incentive to modify their profiles (e.g., credit score), potentially at certain costs, without actually improving their financial status. It is crucial that classification algorithms used for decision-making be robust to such strategic manipulation.

Besides the strategic feature manipulation, another common challenge in classification-based decision-making is that the decision-maker often only observes partial feedback. In particular, the decision-maker may only observe the true labels of agents who have received the positive decision. For example, the bank can observe the true financial qualification only for applicants who have already been classified as qualified and granted loan approvals, but has no chance to observe the true qualification of rejected applicants. This type of partial feedback is sometimes called one-sided feedback, apple-tasting feedback [e.g., Harris et al., 2023, Helmbold et al., 2000] or selective label feedback [e.g., Lakkaraju et al., 2017, Chen et al., 2025].

In this paper, we study an online strategic classification problem with partial feedback. In this problem, a principal (decision-maker) interacts with sequentially arriving agents. The principal

*Corresponding author.

announces a binary linear classifier to each agent and makes the decision according to the agent’s reported feature that may differ from the truth due to strategic manipulation. Following the existing literature, we assume that the agents manipulate their features to maximize the net utility from the classification decision and the cost of feature manipulation. We assume that the agents’ true feature-label relationship is characterized by a linear halfspace model with arbitrary feature-dependent but bounded noise (i.e., Massart noise). This model is widely adopted in the learning theory literature (see references in Section 1.1). The principal does not know the true feature-label relationship, but needs to learn accurate binary classifiers from observations of agents’ reported features (but not the original true features) and the true labels of only positively classified agents.

Notably, this problem faces the combination of three challenges: agents’ strategic feature manipulations, partial label observations, and label noise. First, because of strategic feature manipulation, the agents’ true features may not be faithfully observed by the principal, which impedes the learning process. This is particularly a challenge in the online setting, as the agents’ strategic behaviors depend on the classifiers announced to them, so their behaviors change over time as the classifiers evolve. Second, the principal can only observe true labels from positively classified agents, without feedback from negatively classified agents. This means that the principal can learn only when a positive classification is made, while the strategic agents are incentivized to manipulate their features to achieve positive classification. Third, the label noise results in noisy feedback, which further complicates the learning process.

Our work contributes to the literature along the following dimensions. First, to the best of our knowledge, our work is the first to study online strategic classification under Massart noise and partial feedback. This advances the literature of learning halfspaces under noise [e.g., Zhang et al., 2020, Diakonikolas et al., 2020] to the strategic setting. Moreover, within the online strategic classification literature, our halfspace model with Massart noise extends the noise-free model of deterministic feature-label relationship in Ahmadi et al. [2021], Shen et al. [2024] and complements the fully adversarial setting [Dong et al., 2018, Chen et al., 2020]. Second, we propose a novel learning algorithm that effectively addresses the aforementioned three key challenges. This algorithm has an initialization-refinement-enhancement pipeline, proceeding in batches and iterations. It features several key components: 1) a localization scheme that iteratively improves the classifiers via online linear optimization, using data within increasingly narrow bands around the classification boundary; 2) a projection-based method to construct proxy features from agents’ reported features; 3) a pairwise contrastive inference technique to infer information of the localization bands by contrasting data from pairs of carefully constructed classifiers. Third, we rigorously prove that the proposed algorithm yields classifiers that converge to the clairvoyant optimal one and attains a regret rate of $O(\sqrt{T})$ up to poly-logarithmic and constant factors over T cycles.

1.1 Related Literature

Strategic Classification Strategic classification, introduced by Hardt et al. [2016], has gained increasing attention. The existing literature has studied strategic classification in both offline settings [e.g., Hardt et al., 2016, Sundaram et al., 2023, Levanon and Rosenfeld, 2021] and online settings. In online strategic classification, a principal sequentially interacts with strategic agents, aiming to learn accurate classifiers in the presence of strategic feature manipulation. Some literature models agents’ strategic behaviors by a manipulation graph that defines agents’ feasible feature manipulations [Ahmadi et al., 2023, 2024, Cohen et al., 2024, Shao et al., 2025]. Meanwhile, other literature considers agents that maximize the utility net the cost of feature manipulation. For example, Dong et al. [2018] derive conditions on the manipulation cost function that enable convex optimization techniques to achieve a sublinear regret rate under different fractions of strategic agents. Chen et al. [2020] consider a distance-based manipulation cost function and a zero-one loss function. Our work considers the same cost function and loss function. However, Chen et al. [2020] studies a fully adversarial setting, while our work studies a stochastic setting where agents’ true features and labels follow some probability distributions. Our work is closely related to Ahmadi et al. [2021] and Shen et al. [2024], as we study similar models for agents’ strategic behaviors and linear classifiers. However, their works focus on the noise-free setting with a deterministic feature-label relationship, while our work tackles label noise.

Notably, nearly all prior studies focus on full feedback settings, whereas our work studies a partial feedback setting. One exception is Harris et al. [2023], where the feedback can be observed also

only under a positive decision. However, they consider continuous feedback following a linear regression model with feature-independent noise and target a different objective. In contrast, our work considers binary classification feedback and studies a halfspace model with potentially feature-dependent bounded noise, which directly extends the models in Ahmadi et al. [2021], Shen et al. [2024].

Learning Halfspaces with Noise Our paper adopts a halfspace model with the label flipped at a potentially feature-dependent bounded probability, i.e., Massart noise [Massart and Nédélec, 2006]. Recent studies find that even in the absence of strategic agent manipulations, learning halfspaces under Massart noise presents significant challenges [Zhang et al., 2020, Diakonikolas et al., 2019, 2020, 2024]. The key challenge stems from the nonconvexity of the 0-1 loss function that characterizes the misclassification error. A standard approach to overcome the non-convexity of 0-1 loss in classification is to use a convex surrogate loss function [Bartlett et al., 2006]. However, Awasthi et al. [2015] show that popular algorithms such as SVM or hinge loss minimization fails to learn a halfspace that achieves arbitrarily small excess error under Massart noise. More generally, Diakonikolas et al. [2019] show that one cannot achieve non-trivial misclassification error for learning halfspaces under Massart Noise by optimizing convex surrogates. Instead, a “localization” scheme has been proposed to learn halfspaces under a variety of noise models [e.g., Shen, 2021a, Awasthi et al., 2017, Zhang and Li, 2021, Shen, 2021b, Awasthi et al., 2017]. In particular, Zhang et al. [2020] and Diakonikolas et al. [2020] apply localization to learn halfspaces with Massart noise. The core idea is to iteratively improve classification via convex optimization, using data within increasingly narrow bands around the classification boundary. This localization scheme focuses more on data near the classification boundary, as data far away from the boundary tend to be less informative since they can be either easily correctly classified or misclassified mainly due to noise. However, naïvely extending this localization scheme to strategic classification poses significant challenges, because data points close to the classification boundary are the most prone to feature manipulation. Our work effectively overcomes these challenges by leveraging carefully constructed proxy data and a novel pairwise contrastive inference approach.

1.2 Notation

We employ the following notation throughout the paper. Boldface letters such as $\mathbf{x}, \mathbf{r}, \mathbf{w}$ denote vectors. The operator $\|\cdot\|_p$ denotes any ℓ_p norm of a vector. The inner product of two vectors is denoted by $\langle \cdot, \cdot \rangle$ and the angle between two vectors is represented by $\theta(\cdot, \cdot)$, i.e., $\theta(\mathbf{v}_1, \mathbf{v}_2) = \arccos\left(\frac{\langle \mathbf{v}_1, \mathbf{v}_2 \rangle}{\|\mathbf{v}_1\|_2 \cdot \|\mathbf{v}_2\|_2}\right)$ for $\forall \mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^d$. Symbols \mathbb{B}^d and \mathbb{S}^d denote the d -dimensional Euclidean unit ball and sphere, respectively. $\mathbb{B}^d(R)$ denotes the ball with radius $R > 0$. For any positive integer N , $[N]$ represents the set $\{1, 2, \dots, N\}$. The indicator function $\mathbb{I}(\cdot)$ gives the value 1 if the event within the parentheses holds and the value 0 otherwise.

2 Problem Setup

We consider a setting where a principal repeatedly interacts with sequentially arriving agents (e.g., applicants). Without loss of generality, time is discretized into T cycles, where one agent arrives in each cycle $t \in [T]$. The agent is characterized by a feature-label pair (\mathbf{x}_t, y_t) , where $\mathbf{x}_t \in \mathbb{R}^d$ denotes a d -dimensional feature vector and $y_t \in \{+1, -1\}$ denotes the agent label (i.e., qualified or not). At the beginning of each cycle $t \in [T]$, the principal announces a classifier $\tilde{h}_t(\cdot)$ as the admission rule for the arriving agent. The agent may strategically manipulate and report feature value $\mathbf{r}_t \neq \mathbf{x}_t$ to the principal at some costs, aiming to get admitted (i.e., classified as the positive class, $\tilde{h}_t(\mathbf{r}_t) = +1$). The principal observes the reported features \mathbf{r}_t , makes the classification decision $\tilde{h}_t(\mathbf{r}_t)$ accordingly, and observes the true label y_t only when this agent is admitted. Importantly, the principal has no chance to observe the true label of rejected agents. Based on the data of reported features and admitted agents’ labels, the principal aims to learn accurate classifiers over the T cycles.

Distributional Assumptions We assume that the feature-label pairs $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$ are independently and identically distributed (i.i.d) draws from a common population denoted by (\mathbf{x}, y) . We need to first impose some distributional assumptions on (\mathbf{x}, y) .

Assumption 1 (Halfspace with Massart noise). *There exists a Boolean-valued function $h^*(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}^*, \mathbf{x} \rangle)$ for a coefficient vector \mathbf{w}^* with $\|\mathbf{w}^*\|_2 = 1$ and a noise level bound $\bar{\eta} \in [0, 1/2)$, such that $y = h^*(\mathbf{x})$ with probability $1 - \eta(\mathbf{x})$ and $y = -h^*(\mathbf{x})$ with probability $\eta(\mathbf{x})$, where $\eta(\mathbf{x})$ characterizes the potentially feature-dependent noise satisfying $0 \leq \eta(\mathbf{x}) \leq \bar{\eta}$ almost surely.*

Assumption 1 is a standard assumption in the literature of learning halfspaces without strategic manipulation [e.g., Zhang et al., 2020, Diakonikolas et al., 2020, Massart and Nédélec, 2006]. It allows for arbitrary feature-dependent label noise with an upper bound $\bar{\eta} \in [0, 1/2)$, and relaxes the assumptions in some existing strategic classification literature that assumes a noiseless halfspace model $y = h^*(\mathbf{x})$ and that the positive and negative classes are strictly separated by a margin [e.g., Ahmadi et al., 2021, Shen et al., 2024]. Our assumption can more aptly model real applications where label noises are common and even feature-dependent. Nonetheless, unlike the existing literature, we need to simultaneously handle both the strategic manipulation and label noise.

Assumption 2 (Regular feature distribution). *Fix constants $R, L_1, L_2, U_1, U_2, \delta, Q > 0$, and let \mathbf{x}_V denote the projection of \mathbf{x} onto any subspace $V \subseteq \mathbb{R}^d$ and ϕ_V denote its probability density function. The distribution of features \mathbf{x} satisfies the following regularity conditions for any 1-dimensional subspace $V_1 \subseteq \mathbb{R}$ and any 2-dimensional subspace $V_2 \subseteq \mathbb{R}^2$:*

1. $\phi_{V_1}(\mathbf{x}_{V_1}) \geq L_1$ and $\phi_{V_2}(\mathbf{x}_{V_2}) \geq L_2$ for any $\mathbf{x}_{V_1} \in V_1 \cap \mathbb{B}^1(R)$, $\mathbf{x}_{V_2} \in V_2 \cap \mathbb{B}^2(R)$.
2. $\phi_{V_1}(\mathbf{x}_{V_1}) \leq U_1$ and $\phi_{V_2}(\mathbf{x}_{V_2}) \leq U_2 e^{-\delta \|\mathbf{x}_{V_2}\|^2}$ for any $\mathbf{x}_{V_1} \in V_1$, $\mathbf{x}_{V_2} \in V_2$.
3. For any $t > 0$ and unit vector $\mathbf{w} \in \mathbb{S}^d$, we have that $\mathbb{P}[\langle \mathbf{w}, \mathbf{x} \rangle \geq t] \leq \exp(1 - Qt)$.

In Assumption 2, condition 1 requires that the densities of any 1-dimensional and 2-dimensional projections of feature \mathbf{x} are lower bounded around the origin. Condition 2 indicates that these densities have proper upper bounds. Condition 3 requires that the inner product of \mathbf{x} with any unit vector \mathbf{w} has a sub-exponential tail bound. These conditions generalize the feature distribution conditions in a large body of literature on learning halfspaces with noise [e.g., Diakonikolas et al., 2020, 2021, Zhang et al., 2020, Dasgupta, 2005, Yan and Zhang, 2017, Shen, 2021a, Awasthi et al., 2017]. This existing literature typically assumes that the feature \mathbf{x} has an isotropic log-concave distribution, such as a uniform distribution over a unit sphere. In Appendix A.1, we show that Assumption 2 accommodates even non-isotropic log-concave distributions, including many common distributions such as uniform, Gaussian, exponential, logistic distributions, etc. Notably, we impose distributional assumptions on the feature-label pairs, which differ from and complement the fully adversarial setting in the literature [e.g., Dong et al., 2018, Chen et al., 2020, Ahmadi et al., 2024].

Agent Feature Manipulation We assume that each agent gains a utility of +1 for admission (classified as +1) and −1 for rejection (classified as −1). An agent with true feature \mathbf{x} may report his feature as \mathbf{r} to sway the classifier’s decision. Following Shen et al. [2024], Ahmadi et al. [2021], we assume that this misreporting or manipulation incurs a cost $\text{Cost}(\mathbf{x}, \mathbf{r}) = 2\|\mathbf{x} - \mathbf{r}\|_2/\gamma$, where $\gamma > 0$ indicates the maximum manipulation distance. Therefore, upon the principal announcing a classifier $\tilde{h}(\cdot)$, the agent’s optimal reported feature that maximizes the net utility would be $\mathbf{r}^*(\mathbf{x}, \tilde{h}) = \arg\max_{\mathbf{r} \in \mathbb{R}^d} \tilde{h}(\mathbf{r}) - 2\|\mathbf{x} - \mathbf{r}\|_2/\gamma$.

Given the linear model in Assumption 1, we restrict the principal’s classifier \tilde{h} to linear classifiers parameterized by $(\mathbf{w}, m) \in \mathbb{S}^d \times \mathbb{R}$, i.e., $\tilde{h}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}, \mathbf{r} \rangle + m)$. In this case, an agent’s optimal reported feature is given in the following lemma [Shen et al., 2024, Ahmadi et al., 2021].

Lemma 1. *Given an announced classifier $\tilde{h}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}, \mathbf{r} \rangle + m)$, the optimal reported feature for an agent with true feature \mathbf{x} is*

$$\mathbf{r}^*(\mathbf{x}, \tilde{h}) = \begin{cases} \mathbf{x} - (\langle \mathbf{w}, \mathbf{x} \rangle + m)\mathbf{w}, & -\gamma \leq \langle \mathbf{w}, \mathbf{x} \rangle + m < 0; \\ \mathbf{x}, & \text{otherwise.} \end{cases}$$

The Clairvoyant Optimal Classifier Under the manipulated feature in Lemma 1, the misclassification rate of a classifier \tilde{h} can be measured by $\text{Err}(\tilde{h}) := \mathbb{P}(\tilde{h}(\mathbf{r}^*(\mathbf{x}, \tilde{h})) \neq y)$. We hope to characterize a clairvoyant optimal classifier achieving the minimal misclassification rate: $\tilde{h}^* \in \arg\min_{\tilde{h}: \mathbb{R}^d \rightarrow \{\pm 1\}} \text{Err}(\tilde{h})$. To this end, we first connect a classifier \tilde{h} under the manipulated feature \mathbf{r} with a hypothetical classifier h under the corresponding true feature \mathbf{x} .

Proposition 1. For any $(\mathbf{w}, m) \in \mathbb{S}^d \times \mathbb{R}$, the output of $\tilde{h}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}, \mathbf{r} \rangle + m - \gamma)$ for $\mathbf{r} = \mathbf{r}^*(\mathbf{x}, \tilde{h})$ is identical to the output of $h(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + m)$ for any $\mathbf{x} \in \mathbb{R}^d$.

According to Assumption 1, the optimal classifier in absence of manipulation is $h^*(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}^*, \mathbf{x} \rangle)$. Following Proposition 1, we can achieve the same classification by a corresponding classifier subject to manipulation, which gives the clairvoyant optimal classifier. This structural knowledge of a clairvoyant optimal classifier will guide our algorithm design in Section 3.

Corollary 1. The classifier $\tilde{h}^*(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}^*, \mathbf{r} \rangle - \gamma)$ minimizes $\text{Err}(\tilde{h}) = \mathbb{P}(\tilde{h}(\mathbf{r}^*(\mathbf{x}, \tilde{h})) \neq y)$.

Notably, the clairvoyant optimal classifier on the manipulated feature \mathbf{r} has a higher threshold to classify an agent into +1 than the corresponding optimal classifier $h^*(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}^*, \mathbf{x} \rangle)$ on the true feature \mathbf{x} . Indeed, the principal would like to raise the bar for positive classification, in order to avoid errors due to unqualified agents (label -1) who game the classifier by manipulating their features.

Principal’s Regret Over the T cycles, the principal learns a sequence of classifiers $\tilde{\mathbf{h}} = (\tilde{h}_1, \dots, \tilde{h}_T)$, where each \tilde{h}_t only depends on the observed data of reported features and admitted agents’ labels prior to cycle t . The goal is to achieve a small cumulative misclassification rate over all cycles. This is equivalent to achieving a small total suboptimality gap, or regret, relative to the clairvoyant optimal classifier. Formally, the regret is defined as:

$$\text{Reg}(\tilde{\mathbf{h}}; T) := \sum_{t=1}^T \text{Err}(\tilde{h}_t) - T \times \text{Err}(\tilde{h}^*). \quad (1)$$

This regret corresponds to the “Stackelberg regret” in the strategic classification literature, where the term “Stackelberg” emphasizes that agents consistently choose their best feature manipulation in response to the principal’s announced classifiers [Dong et al., 2018, Chen et al., 2020, Ahmadi et al., 2024]. In the next section, we will propose a learning algorithm that effectively tackles the combined challenges of agents’ feature manipulations, partial feedback observations, and label noise. We prove that this algorithm achieves a \sqrt{T} -regret rate up to poly-logarithmic and constant factors.

3 The Algorithm

3.1 Overview of our Algorithm

Algorithm 1: Main-Algorithm

Input: Maximum manipulation distance γ , noise level bound $\bar{\eta}$, lengths $\{T_{\text{init}}\} \cup \{T_k\}_{k=0}^K$, bandwidths $\{b_k\}_{k=0}^K$, stepsizes $\{\alpha_k\}_{k=0}^K$, feature dimension d

```

1  $\bar{\mathbf{w}}_0 = \text{Initialization}(T_{\text{init}})$  // See Algorithm 2
2  $\mathbf{w}_1 = \text{Refinement}(\bar{\mathbf{w}}_0, \bar{\eta}, T_0, b_0, \alpha_0, d)$  // See Algorithm 3
3 for  $k \leftarrow 1$  to  $K$  do
4    $\mathbf{w}_{k+1} = \text{Batched-Enhancement}(\gamma, \mathbf{w}_1, \bar{\eta}, k, T_k, b_k, \alpha_k, d)$  // See Algorithm 4
```

Our main Algorithm, outlined in Algorithm 1, comprises three sub-algorithms: an Initialization Algorithm (Algorithm 2), a Refinement Algorithm (Algorithm 3) and a Batched Enhancement Algorithm (Algorithm 4). These algorithms are executed sequentially to generate a sequence of coefficient vectors such that the corresponding classifiers converge to the clairvoyant optimal classifier as specified in Corollary 1. Specifically, we partition the horizon of T cycles (one agent arrives in each cycle) into consecutive batches indexed by $k \in \{\text{init}, 0, 1, 2, \dots, K\}$. Index “init” and “0” denote the batches executing the Initialization and Refinement Algorithms, respectively, while indices “1” to “ K ” represent the K batches that run the Enhancement Algorithm iteratively. Each batch k takes the result of the previous batch $k - 1$ as input. Cycles in each batch k are further grouped into iterations indexed by $i \in \{1, 2, \dots, T_k\}$, where each iteration i performs an update for the coefficient vector \mathbf{w} . At the end of batch k , the algorithms output the (normalized) average vectors of the T_k iterations in the batch. During Refinement, each iteration consists of only one cycle. In contrast, during Initialization and Enhancement, each iteration contains two cycles, denoted by the superscript $j = 1$ or 2 to differentiate between the first and second cycles within the same iteration. Note that the indices (k, i, j) can be mapped to the corresponding cycle t , for convenience, we will use these indices in the remainder of this paper.

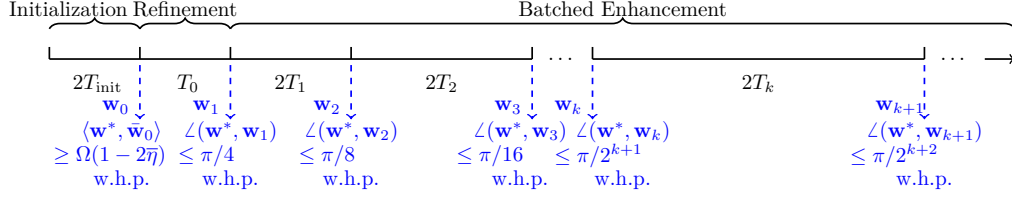


Figure 1: Roles of the three sub-algorithms

The roles of the three sub-algorithms are summarized in Figure 1. First, the *Initialization Algorithm* runs for $T_{\text{init}} = O(\ln T / (1 - 2\bar{\eta})^2)$ iterations (with $2T_{\text{init}}$ cycles) to find a coefficient vector $\bar{\mathbf{w}}_0$ such that $\theta(\mathbf{w}^*, \bar{\mathbf{w}}_0) \leq \frac{\pi}{2}$ with high probability (see Proposition 2). Second, the *Refinement Algorithm* takes $\mathbf{w}_0 = \bar{\mathbf{w}}_0 / \|\bar{\mathbf{w}}_0\|_2$ as the initial vector and runs for $T_0 = O(d \ln d \ln T / (1 - 2\bar{\eta})^8)$ iterations (with T_0 cycles) to obtain a refined vector \mathbf{w}_1 such that $\theta(\mathbf{w}^*, \mathbf{w}_1) \leq \frac{\pi}{4}$ with high probability (see Proposition 3). Third, the *Batched Enhancement Algorithm* runs for $K = O(\log_4(1 - 2\bar{\eta})^4 T / (\gamma d \ln d \ln T))$ batches, where each batch k enhances its initial coefficient vector \mathbf{w}_k through $T_k = O(4^k d \ln d \ln T / (1 - 2\bar{\eta})^4)$ iterations (with $2T_k$ cycles), yielding a vector \mathbf{w}_{k+1} such that $\theta(\mathbf{w}^*, \mathbf{w}_{k+1}) \leq \frac{\pi}{2^{k+2}}$ with high probability (see Proposition 4). The specification of the algorithms involves absolute constants c_0 to c_7 , which are derived from the parameters in Assumption 2. Detailed calculations are available in Appendix A.5.

3.2 Initialization

Algorithm 2: Initialization

Input: Iteration length T_{init}

```

1 for  $i \leftarrow 1$  to  $T_{\text{init}}$  do
2   Uniformly draw  $\mathbf{w}_{\text{init},i} \in \mathbb{S}^d$ 
3   for  $j \leftarrow 1$  to 2 do
4     Declare  $\tilde{h}_{\text{init},i}^{(j)}(\mathbf{r}) = (-1)^{j-1} \text{sgn}(\langle \mathbf{w}_{\text{init},i}, \mathbf{r} \rangle)$ , agent  $(\mathbf{x}_{\text{init},i}^{(j)}, y_{\text{init},i}^{(j)})$  arrives and reports  $\mathbf{r}_{\text{init},i}^{(j)}$ 
5     Make classification decision  $\tilde{h}_{\text{init},i}^{(j)}(\mathbf{r}_{\text{init},i}^{(j)})$  and collect label  $y_{\text{init},i}^{(j)}$  if  $\tilde{h}_{\text{init},i}^{(j)}(\mathbf{r}_{\text{init},i}^{(j)}) = 1$ 
6 return  $\bar{\mathbf{w}}_0 = \frac{1}{T_{\text{init}}} \sum_{i=1}^{T_{\text{init}}} \sum_{j=1}^2 y_{\text{init},i}^{(j)} \mathbf{r}_{\text{init},i}^{(j)} \mathbb{I} \left( (-1)^{(j-1)} \langle \mathbf{w}_{\text{init},i}, \mathbf{r}_{\text{init},i}^{(j)} \rangle > 0 \right)$ 
```

The initialization algorithm runs for $T_{\text{init}} = O(\ln T / (1 - 2\bar{\eta})^2)$ iterations. In each iteration, we randomly explore a coefficient vector $\mathbf{w}_{\text{init},i} \in \mathbb{S}^d$ and offer two opposing classifiers based on $\mathbf{w}_{\text{init},i}$ to two successive agents. Using the reported features $\mathbf{r}_{\text{init},i}^{(1)}, \mathbf{r}_{\text{init},i}^{(2)}$ and true labels $y_{\text{init},i}^{(1)}, y_{\text{init},i}^{(2)}$ of all positively classified agents over the T_{init} iterations, we construct an initial coefficient vector $\bar{\mathbf{w}}_0$.

The design of our initialization algorithm stems from the well-known ‘‘averaging’’ technique for learning halfspaces [Servedio, 2001]. In the non-strategic, noiseless and full feedback setting, $y\langle \mathbf{w}^*, \mathbf{x} \rangle = \langle \mathbf{w}^*, y\mathbf{x} \rangle \geq 0$ for all $(\mathbf{x}, y) \in \mathbb{R}^d \times \{\pm 1\}$, so $y\mathbf{x}$ forms an acute angle with the optimal normal vector \mathbf{w}^* almost surely (since the set $\{\mathbf{x} \mid \langle \mathbf{w}^*, \mathbf{x} \rangle = 0\}$ is zero-measure). Analogously, in the noisy feedback setting, we have $\langle \mathbf{w}^*, \mathbb{E}[y\mathbf{x}] \rangle > 0$, so $\mathbb{E}[y\mathbf{x}]$ forms an acute angle with \mathbf{w}^* . In a non-strategic and full feedback setting, the literature uses the sample average of $y\mathbf{x}$ to approximate $\mathbb{E}[y\mathbf{x}]$ as an initial estimate of \mathbf{w}^* [Zhang et al., 2020]. However, this estimator is unavailable for us because of agents’ feature manipulation and the partial feedback setting. Instead, our algorithm declares pairs of opposite classifiers. We collect and average the $y\mathbf{r}$ of agents whose reported feature \mathbf{r} falls above the hyperplane. Note that these agents report their features truthfully ($\mathbf{r} = \mathbf{x}$), so we are able to form $\bar{\mathbf{w}}_0$ from these agents’ $y\mathbf{r}$ as a proper approximation of $\mathbb{E}[y\mathbf{x}]$. We can show that, for large enough T_{init} , this vector $\bar{\mathbf{w}}_0$ forms an acute angle with the optimal \mathbf{w}^* with high probability.

Proposition 2. *For some constants $c_0, c_1 > 0$, when Algorithm 2 runs for $T_{\text{init}} = c_0 \ln T / (1 - 2\bar{\eta})^2$ iterations, its output $\bar{\mathbf{w}}_0$ satisfies $\langle \mathbf{w}^*, \bar{\mathbf{w}}_0 \rangle > c_1(1 - 2\bar{\eta}) > 0$ and $\theta(\mathbf{w}^*, \bar{\mathbf{w}}_0) \leq \frac{\pi}{2}$ with probability at least $1 - 2/T^2$.*

3.3 Refinement of the Initial Coefficient Vector

Algorithm 3: Refinement

Input: Initial vector $\bar{\mathbf{w}}_0$, noise level $\bar{\eta}$, iteration length T_0 , bandwidth b_0 , step size α_0 , feature dimension d

Initialization: $\mathbf{w}_{0,1} = \bar{\mathbf{w}}_0 / \|\bar{\mathbf{w}}_0\|_2$

```

1 for  $i \leftarrow 1$  to  $T_0$  do
2   Declare classifier  $\tilde{h}_{0,i}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}_{0,i}, \mathbf{r} \rangle)$ , agent  $(\mathbf{x}_{0,i}, y_{0,i})$  arrives and reports  $\mathbf{r}_{0,i}$ 
3   Make classification decision  $\tilde{h}_{0,i}(\mathbf{r}_{0,i})$  and collect label  $y_{0,i}$  if  $\tilde{h}_{0,i}(\mathbf{r}_{0,i}) = 1$ 
4   Compute gradient:  $\tilde{\mathbf{g}}_{0,i} = [-\bar{\eta}\mathbf{r}_{0,i}\mathbb{I}(y_{0,i} = 1) + (1 - \bar{\eta})\mathbf{r}_{0,i}\mathbb{I}(y_{0,i} = -1)]\mathbb{I}(0 < \langle \mathbf{w}_{0,i}, \mathbf{r}_{0,i} \rangle \leq b_0)$ 
5   Set constraint set:  $\mathcal{W}_0 = \{\mathbf{w} \mid \|\mathbf{w}\|_2 \leq 1, \langle \mathbf{w}, \bar{\mathbf{w}}_0 \rangle \geq c_1(1 - 2\bar{\eta})\}$ 
6   Update  $\mathbf{w}$ :  $\hat{\mathbf{w}}_{0,i+1} = \arg \min_{\mathbf{w} \in \mathcal{W}_0} \langle \tilde{\mathbf{g}}_{0,i}, \mathbf{w} \rangle + \frac{1}{\alpha_0} \frac{\|\mathbf{w} - \mathbf{w}_{0,i}\|_p^2}{2(p-1)}$ , where  $p = \frac{\ln(8d)}{\ln(8d)-1}$ 
7   Normalize:  $\mathbf{w}_{0,i+1} = \hat{\mathbf{w}}_{0,i+1} / \|\hat{\mathbf{w}}_{0,i+1}\|_2$ 
8 Compute mean vector:  $\bar{\mathbf{w}}_1 = \frac{1}{T_0} \sum_{i=1}^{T_0} \mathbf{w}_{0,i}$ 
9 return  $\mathbf{w}_1 = \bar{\mathbf{w}}_1 / \|\bar{\mathbf{w}}_1\|_2$ 

```

The refinement algorithm adopts a “localization” scheme to refine the output $\bar{\mathbf{w}}_0$ of Algorithm 2 to better approximate \mathbf{w}^* . In every iteration i , we consider only data within a band $0 < \langle \mathbf{w}_{0,i}, \mathbf{r} \rangle \leq b_0$ adjacent to the boundary of the current classifier $\tilde{h}_{0,i}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}_{0,i}, \mathbf{r} \rangle)$. Agents in this band are positively classified and have no incentives for feature manipulation, allowing us to observe both the true feature $\mathbf{x} = \mathbf{r}$ and the true label y . Moreover, this effectively probes the localized region $D_{0,i} = \{\mathbf{x} : 0 < \langle \mathbf{w}_{0,i}, \mathbf{x} \rangle \leq b_0\}$ in the true feature space. Similar “localization” is widely used in the literature of learning half-spaces with label noises (see references in Section 1.1), since data near the classification boundary is the most informative, while data far from the boundary are either correctly classified with ease or are misclassified mainly due to noises, providing little information.

We formulate an online linear optimization problem with constructed losses $\{\mathbf{w} \mapsto \langle \mathbf{w}, \tilde{\mathbf{g}}_{0,i} \rangle\}_{i=1}^{T_0}$ over a proper constraint set $\mathcal{W}_0 = \{\mathbf{w} \mid \|\mathbf{w}\|_2 \leq 1, \langle \mathbf{w}, \bar{\mathbf{w}}_0 \rangle \geq c_1(1 - 2\bar{\eta})\}$. This constraint set, according to Proposition 2, contains \mathbf{w}^* with high probability. We then solve this problem by online mirror descent with a stepsize α_0 and regularizer $\|\mathbf{w} - \mathbf{w}_{0,i}\|_p^2 / 2(p-1)$ for $p = \ln(8d) / (\ln(8d) - 1)$, perform proper normalization in each iteration, and normalize the average of all iterates to obtain the output \mathbf{w}_1 . By focusing on the band $\{\mathbf{r} : 0 < \langle \mathbf{w}_{0,i}, \mathbf{r} \rangle \leq b_0\}$, we can perfectly observe $\mathbf{r}_{0,i} = \mathbf{x}_{0,i}$ and $y_{0,i}$ and the gradients² $\tilde{\mathbf{g}}_{0,i}$ coincide with the counterparts in Zhang et al. [2020]. As a result, we can follow their analysis to bound the error of the output \mathbf{w}_1 .

Proposition 3. *For the constant c_1 in Proposition 2 and some constants $c_2, c_3, c_4 > 0$, when the initial vector $\bar{\mathbf{w}}_0$ satisfies $\langle \mathbf{w}^*, \bar{\mathbf{w}}_0 \rangle \geq c_1(1 - 2\bar{\eta})$ and Algorithm 3 runs with bandwidth $b_0 = c_2(1 - 2\bar{\eta})^2$ for $T_0 = c_3 d \ln d (\ln T)^2 / (1 - 2\bar{\eta})^8$ iterations with step size $\alpha_0 = c_4 \sqrt{d \ln d} / (\sqrt{T_0} \ln T)$, then its output \mathbf{w}_1 satisfies $\theta(\mathbf{w}^*, \mathbf{w}_1) \leq \pi/4$ with probability at least $1 - 3/T^2$.*

The main idea in proving Proposition 3 is outlined as follows. By the theory of online convex optimization, we can upper bound the cumulative regret for the constructed loss in this stage, i.e., $\sum_{i=1}^{T_0} \langle \mathbf{w}_{0,i}, \tilde{\mathbf{g}}_{0,i} \rangle - \langle \mathbf{w}^*, \tilde{\mathbf{g}}_{0,i} \rangle$. This regret bound, together with a bound on $\sum_{i=1}^{T_0} \langle \mathbf{w}_{0,i}, \tilde{\mathbf{g}}_{0,i} \rangle$ and a concentration bound on $\sum_{i=1}^{T_0} \langle \mathbf{w}^*, -\tilde{\mathbf{g}}_{0,i} \rangle$, leads to a high probability upper bound on $\sum_{i=1}^{T_0} \mathbb{E}[\langle \mathbf{w}^*, -\tilde{\mathbf{g}}_{0,i} \rangle]$. Moreover, it can be shown that $\mathbb{E}[\langle \mathbf{w}^*, -\tilde{\mathbf{g}}_{0,i} \rangle]$ is lower bounded by $\theta(\mathbf{w}^*, \mathbf{w}_{0,i})$ up to some proportional factors. This is why we expect to obtain a high probability upper bound on $\theta(\mathbf{w}^*, \mathbf{w}_1)$. Importantly, the gradients $\tilde{\mathbf{g}}_{0,i}$ are carefully constructed to ensure that $|\langle \mathbf{w}_{0,i}, \tilde{\mathbf{g}}_{0,i} \rangle|$ is small and meanwhile $\mathbb{E}[\langle \mathbf{w}^*, -\tilde{\mathbf{g}}_{0,i} \rangle]$ upper bounds $\theta(\mathbf{w}^*, \mathbf{w}_{0,i})$.

Notably, while Algorithm 3 collects true feature-label data and implements localization by focusing on local bands around the origin-crossing classification hyperplanes $\tilde{h}_{0,i}$ ’s, this approach can be costly. According to Lemma 1, unqualified agents with true features \mathbf{x} satisfying $-\gamma \leq \langle \mathbf{w}_{0,i}, \mathbf{x} \rangle \leq 0$ would manipulate their features to achieve positive classifications, resulting in constant instantaneous regret. Fortunately, Algorithm 3 runs for only $\tilde{O}(\ln T)$ cycles, so this refinement algorithm obtains an improved coefficient \mathbf{w}_1 for the next stage at the cost of at most only $\tilde{O}(\ln T)$ regret.

²It can be verified that $\tilde{\mathbf{g}}_{0,i}$ is the gradient of a Leaky ReLU loss restricted to the band $D_{0,i}$.

3.4 Batched Enhancement: Proxy Features and Pairwise Contrastive Inference

Algorithm 4: Batched Enhancement

Input: Maximum manipulation distance γ , initial vector \mathbf{w}_k , noise level $\bar{\eta}$, batch index k , iteration length T_k , bandwidth b_k , step size α_k , feature dimension d

Initialization: $\mathbf{w}_{k,1} = \mathbf{w}_k$

- 1 **for** $i \leftarrow 1$ **to** T_k **do**
 - 2 Construct classifiers $\tilde{h}_{k,i}^{(1)}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}_{k,i}, \mathbf{r} \rangle - \gamma)$ and $\tilde{h}_{k,i}^{(2)}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}_{k,i}, \mathbf{r} \rangle - \gamma - b_k)$
 - 3 **for** $j \leftarrow 1$ **to** 2 **do**
 - 4 Declare classifier $\tilde{h}_{k,i}^{(j)}$, agent $(\mathbf{x}_{k,i}^{(j)}, y_{k,i}^{(j)})$ arrives and reports $\mathbf{r}_{k,i}^{(j)}$
 - 5 Make classification decision $\tilde{h}_{k,i}^{(j)}(\mathbf{r}_{k,i}^{(j)})$ and collect label $y_{k,i}^{(j)}$ if $\tilde{h}_{k,i}^{(j)}(\mathbf{r}_{k,i}^{(j)}) = 1$
 - 6 Construct proxy data: $\hat{\mathbf{x}}_{k,i}^{(j,+)} = \text{Proj}_{D_{k,i}}^+(\mathbf{r}_{k,i}^{(j)}) \mathbb{I}(y_{k,i}^{(j)} = 1, \mathbf{r}_{k,i}^{(j)} \in \tilde{D}_{k,i}^{(j)}), \hat{\mathbf{x}}_{k,i}^{(j,-)} = \text{Proj}_{D_{k,i}}^-(\mathbf{r}_{k,i}^{(j)}) \mathbb{I}(y_{k,i}^{(j)} = -1, \mathbf{r}_{k,i}^{(j)} \in \tilde{D}_{k,i}^{(j)})$
 - 7 Use the proxy data to compute the gradient: $\hat{\mathbf{g}}_{k,i} = -\bar{\eta}(\hat{\mathbf{x}}_{k,i}^{(1,+)} - \hat{\mathbf{x}}_{k,i}^{(2,+)} + (1 - \bar{\eta})(\hat{\mathbf{x}}_{k,i}^{(1,-)} - \hat{\mathbf{x}}_{k,i}^{(2,-)})$
 - 8 Update: $\hat{\mathbf{w}}_{k,i+1} \leftarrow \arg \min_{\mathbf{w} \in \mathcal{W}_k} \langle \hat{\mathbf{g}}_{k,i}, \mathbf{w} \rangle + \frac{1}{\alpha_k} \frac{\|\mathbf{w} - \mathbf{w}_{k,i}\|_p^2}{2(p-1)}$, where $p = \frac{\ln(8d)}{\ln(8d)-1}$, the constraint set $\mathcal{W}_k = \{\mathbf{w} \mid \|\mathbf{w}\|_2 \leq 1, \langle \mathbf{w}, \mathbf{w}_k \rangle \geq \cos \theta_k\}$, starting angle $\theta_k = \frac{\pi}{2^{k+1}}$
 - 9 Normalize: $\mathbf{w}_{k,i+1} = \hat{\mathbf{w}}_{k,i+1} / \|\hat{\mathbf{w}}_{k,i+1}\|_2$
 - 10 Compute mean vector $\bar{\mathbf{w}}_{k+1} = \frac{1}{T_k} \sum_{i=1}^{T_k} \mathbf{w}_{k,i}$
 - 11 **return** $\mathbf{w}_{k+1} = \bar{\mathbf{w}}_{k+1} / \|\bar{\mathbf{w}}_{k+1}\|_2$
-

In the non-strategic and full feedback setting, after obtaining the refined coefficient \mathbf{w}_1 , Zhang et al. [2020] further improves it by solving a sequence of adaptively constructed online linear optimization problems $\min_{\mathbf{w} \in \mathcal{W}_k} \sum_{i=1}^{T_k} \langle \mathbf{w}, \mathbf{g}_{k,i} \rangle$ with $\mathbf{g}_{k,i} = [-\bar{\eta} \mathbf{x}_{k,i} \mathbb{I}(y_{k,i} = 1) + (1 - \bar{\eta}) \mathbf{x}_{k,i} \mathbb{I}(y_{k,i} = -1)] \mathbb{I}(-b_k < \langle \mathbf{w}_{k,i}, \mathbf{x}_{k,i} \rangle \leq b_k)$ via mirror descent over $k = 1 \dots, K$ batches, using local data within increasingly narrow bands $\{\mathbf{x} \mid -b_k < \langle \mathbf{w}_{k,i}, \mathbf{x} \rangle \leq b_k\}$ around the classification hyperplanes. This process can geometrically reduce the error of the coefficient estimates, outputting a final classifier that approaches the optimal classifier after enough batches. The key ingredient underlying this guarantee is that the gradients $\mathbf{g}_{k,i}$ are well constructed so that $|\langle \mathbf{w}_{k,i}, \mathbf{g}_{k,i} \rangle|$ is small and meanwhile $\mathbb{E}[\langle \mathbf{w}^*, -\mathbf{g}_{k,i} \rangle]$ upper bounds $\theta(\mathbf{w}^*, \mathbf{w}_{k,i})$ (see discussions below Proposition 3). One may consider directly implementing this batched enhancement approach in our strategic classification. In particular, one may again use classifiers $\tilde{h}_{k,i}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}_{k,i}, \mathbf{r} \rangle)$ and focus on the band $D_{k,i} = \{\mathbf{x} \mid 0 < \langle \mathbf{w}_{k,i}, \mathbf{x} \rangle \leq b_k\}$ in each batch k and iteration i , since this enables us to collect the true feature-label data and probe the localized region $D_{k,i}$. However, as we discussed at the end of Section 3.3, this approach may result in constant instantaneous regret in every cycle due to unqualified strategic agents, so that $O(T)$ regret accumulates over the $O(T)$ cycles in this stage.

To avoid excessive errors due to feature manipulation, we can instead employ classifiers $\tilde{h}_{k,i}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}_{k,i}, \mathbf{r} \rangle - \gamma)$, mimicking the form of the clairvoyant optimal strategic classifier \tilde{h}^* and raising the bar for positive classification to tackle strategic behaviors (see Corollary 1). However, this gives rise to new challenges: it is unclear how to construct the gradients $\mathbf{g}_{k,i}$ and probe the localized regions $D_{k,i}$, since both depend on the true features, but all agents in the localized regions $D_{k,i}$ misreport their features. This means that we know neither which agents' true feature values belong to the regions $D_{k,i}$ nor their true feature values. To tackle these challenges, we propose two key ideas: proxy features and pairwise contrastive inference.

Proxy Features Even if we assume, for the sake of argument, that we can identify agents whose true features lie in $D_{k,i}$, their true feature values remain unobservable, since they all misreport their features to secure positive classification (so their reported feature values fall on the hyperplane of the announced classifier). To resolve this, we construct proxy features from the reported features.

Specifically, consider an agent with true feature value $\mathbf{x} \in D_{k,i}$ and reported feature value \mathbf{r} . This agent will manipulate his feature to get positively classified, and thus we can observe his true label. If his true label is $y = +1$, then we construct his proxy feature $\tilde{\mathbf{x}}$ as the projection of \mathbf{r} onto the *upper* boundary of $D_{k,i}$, i.e., $\tilde{\mathbf{x}} = \text{Proj}_{D_{k,i}}^+(\mathbf{r}) := \mathbf{r} + (b_k - \langle \mathbf{w}_{k,i}, \mathbf{r} \rangle) \mathbf{w}_{k,i}$. On the contrary, if his true label is $y = -1$, then we construct his proxy feature $\tilde{\mathbf{x}}$ as the projection of \mathbf{r} onto the *lower*

boundary of $D_{k,i}$, i.e., $\tilde{\mathbf{x}} = \text{Proj}_{\bar{D}_{k,i}}(\mathbf{r}) := \mathbf{r} - \langle \mathbf{w}_{k,i}, \mathbf{r} \rangle \mathbf{w}_{k,i}$. As a result, this agent's proxy feature value, like his true feature value, also belongs to $D_{k,i}$, and the proxy feature value under a positive label (i.e., projection onto the upper boundary of $D_{k,i}$) is more aligned with the direction of positive classification than the proxy feature value under a negative label (i.e., projection onto the lower boundary of $D_{k,i}$). See the illustration in Figure 2(a).

Using the proxy features, we can approximate the ideal gradient $\mathbf{g}_{k,i}$ by a proxy gradient $\tilde{\mathbf{g}}_{k,i} = [-\bar{\eta} \text{Proj}_{D_{k,i}}^+(\mathbf{r}_{k,i}) \mathbb{I}(y_{k,i} = 1) + (1 - \bar{\eta}) \text{Proj}_{\bar{D}_{k,i}}(\mathbf{r}_{k,i}) \mathbb{I}(y_{k,i} = -1)] \mathbb{I}(\mathbf{x}_{k,i} \in D_{k,i})$. Although this may not exactly recover the ideal gradient, it is still effective, in that $|\langle \mathbf{w}_{k,i}, \tilde{\mathbf{g}}_{k,i} \rangle|$ is small and $\mathbb{E}[\langle \mathbf{w}^*, -\tilde{\mathbf{g}}_{k,i} \rangle] \geq \mathbb{E}[\langle \mathbf{w}^*, -\mathbf{g}_{k,i} \rangle]$ also upper bounds $\theta(\mathbf{w}^*, \mathbf{w}_{k,i})$ (see Appendix A.5). Therefore, we can use the proxy gradients $\tilde{\mathbf{g}}_{k,i}$ in the algorithm to achieve similar guarantees. Nevertheless, these proxy gradients require knowing whether an agent's true feature value belongs to the localized region $D_{k,i}$ or not, which is still infeasible in our setting. This motivates our second key idea.

Pairwise Contrastive Inference We propose to offer two classifiers $\tilde{h}_{k,i}^{(1)}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}_{k,i}, \mathbf{r} \rangle - \gamma)$ and $\tilde{h}_{k,i}^{(2)}(\mathbf{r}) = \text{sgn}(\langle \mathbf{w}_{k,i}, \mathbf{r} \rangle - \gamma - b_k)$ successively in each iteration. Under classifier $\tilde{h}_{k,i}^{(1)}(\mathbf{r})$, we consider only agents with reported features in $\tilde{D}_{k,i}^{(1)} = \{\mathbf{r} : \gamma \leq \langle \mathbf{w}_{k,i}, \mathbf{r} \rangle \leq \gamma + b_k\}$, while under classifier $\tilde{h}_{k,i}^{(2)}(\mathbf{r})$, we consider only agents with reported features in $\tilde{D}_{k,i}^{(2)} = \{\mathbf{r} : \langle \mathbf{w}_{k,i}, \mathbf{r} \rangle = \gamma + b_k\}$. These agents are all classified into the positive class, so their true labels are observed. Moreover, according to the feature manipulation rule in Lemma 1, these agents have true feature values in $D_{k,i}^{(1)} = \{\mathbf{x} : 0 \leq \langle \mathbf{w}_{k,i}, \mathbf{x} \rangle \leq \gamma + b_k\}$ and $D_{k,i}^{(2)} = \{\mathbf{x} : b_k \leq \langle \mathbf{w}_{k,i}, \mathbf{x} \rangle \leq \gamma + b_k\}$, respectively. Since $D_{k,i} = D_{k,i}^{(1)} \setminus D_{k,i}^{(2)}$ up to a measure-zero set, we can expect to infer distributional properties of the data within the region $D_{k,i}$ of interest by contrasting the data within $D_{k,i}^{(1)}$ and the data within $D_{k,i}^{(2)}$. We call this a *pairwise contrastive inference* approach, which is illustrated in Figure 2(b).

We can use this approach to infer the two key components in the proxy gradient $\mathbf{g}_{k,i}$. Note

$$\mathbb{E} \left[\text{Proj}_{D_{k,i}}^+(\mathbf{r}_{k,i}) \mathbb{I}(y_{k,i} = 1, \mathbf{x}_{k,i} \in D_{k,i}) \right] = \mathbb{E} \left[\hat{\mathbf{x}}_{k,i}^{(1,+)} - \hat{\mathbf{x}}_{k,i}^{(2,+)} \right],$$

where $\hat{\mathbf{x}}_{k,i}^{(j,+)} = \text{Proj}_{D_{k,i}}^+(\mathbf{r}_{k,i}^{(j)}) \mathbb{I}(y_{k,i}^{(j)} = 1, \mathbf{r}_{k,i}^{(j)} \in \tilde{D}_{k,i}^{(j)})$ for $j = 1, 2$. This means that we can use $\hat{\mathbf{x}}_{k,i}^{(1,+)} - \hat{\mathbf{x}}_{k,i}^{(2,+)}$ to unbiasedly infer one key component of $\tilde{\mathbf{g}}_{k,i}$ in expectation. Similarly, we can construct $\hat{\mathbf{x}}_{k,i}^{(1,-)} - \hat{\mathbf{x}}_{k,i}^{(2,-)}$ to infer the other component. This gives our gradient estimate $\hat{\mathbf{g}}_{k,i}$ in Algorithm 4 Line 7, satisfying that $\mathbb{E}[\langle \mathbf{w}^*, -\hat{\mathbf{g}}_{k,i} \rangle] = \mathbb{E}[\langle \mathbf{w}^*, -\tilde{\mathbf{g}}_{k,i} \rangle]$ also upper bounds $\theta(\mathbf{w}^*, \mathbf{w}_{k,i})$.

After getting the gradient estimator $\hat{\mathbf{g}}_{k,i}$, we again conduct online mirror decent with a regularizer similar to that in Algorithm 3 and constraint set $\mathcal{W}_k = \{\mathbf{w} \mid \|\mathbf{w}\|_2 \leq 1, \langle \mathbf{w}, \mathbf{w}_k \rangle \geq \cos \theta_k\}$, where $\theta_k = \frac{\pi}{2^{k+1}}$. Then we output the normalized average coefficient vector \mathbf{w}_{k+1} for the next batch. Our constraint set ensures that $\theta(\mathbf{w}_{k,i}, \mathbf{w}_k) \leq \pi/2^{k+1}$ for all $i \in [T_k]$. Then, when the input vector \mathbf{w}_k satisfies $\theta(\mathbf{w}^*, \mathbf{w}_k) \leq \pi/2^{k+1}$, we have $\theta(\mathbf{w}^*, \mathbf{w}_{k,i}) \leq \theta(\mathbf{w}^*, \mathbf{w}_k) + \theta(\mathbf{w}_{k,i}, \mathbf{w}_k) \leq \pi/2^k$ by a triangular inequality shown in Appendix A.3, Lemma 12. This statement is critical: First, it controls the expected cumulative error in batch k to be $O(\frac{1}{2^k} \cdot T_k)$. Second, the condition that $\theta(\mathbf{w}^*, \mathbf{w}_{k,i}) \leq \pi/2^k$, together with our localized online mirror descent method, ensures that batch k outputs a vector \mathbf{w}_{k+1} that satisfy $\theta(\mathbf{w}^*, \mathbf{w}_{k+1}) \leq \pi/2^{k+2}$ with high probability (see Proposition 4), which is in turn required by the next batch.

Proposition 4. *For some constants $c_5, c_6, c_7 > 0$, when Algorithm 4 runs with an initial vector \mathbf{w}_k satisfying $\theta(\mathbf{w}^*, \mathbf{w}_k) \leq \theta_k = \pi/2^{k+1}$, bandwidth $b_k = c_5(1 - 2\bar{\eta})2^{-k}$ for $T_k = c_6 4^k(\gamma + 1)d \ln d(\ln T)^2/(1 - 2\bar{\eta})^4$ iterations with step size $\alpha_k = c_7 \sqrt{d \ln d} \theta_k/(\sqrt{T_k} \ln T)$, its output \mathbf{w}_{k+1} satisfies $\theta(\mathbf{w}^*, \mathbf{w}_{k+1}) \leq \theta_{k+1} = \frac{\theta_k}{2}$ with probability at least $1 - 6/T^2$.*

Proposition 4 shows that Algorithm 4 enhances its input by reducing the error by half in every batch, generating a sequence of coefficient estimates $(\mathbf{w}_k)_{k=1}^K$ with geometrically decaying errors. Notably, we achieve the enhancement by classifiers $\tilde{h}_{k,i}^{(1)}, \tilde{h}_{k,i}^{(2)}$ that use at least γ classification thresholds and are hence more resilient to errors due to strategic classification, which results in only a sublinear regret, as we will show in Theorem 1.

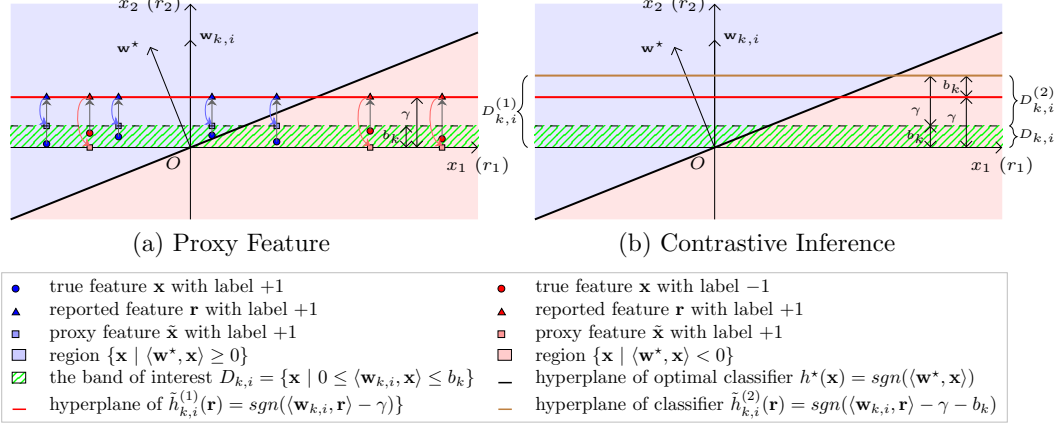


Figure 2: (a) The gray arrows indicate how agents with true feature values (circles) within the band $D_{k,i}$ manipulate their features (triangles). The blue and red arrows indicate how we construct proxy features (squares) from the reported features of agents with labels +1 and -1, respectively. (b) By declaring classifiers $\tilde{h}_{k,i}^{(1)}$ and $\tilde{h}_{k,i}^{(2)}$, we collect data from agents with true values in $D_{k,i}^{(1)}$ and $D_{k,i}^{(2)}$ respectively, through which we infer the information for agents in the region $D_{k,i}$ of interest.

4 Regret Guarantee

We now provide a formal regret guarantee of Algorithm 1, showing that it achieves a sublinear regret dependent on the noise level $\bar{\eta}$ and feature dimension d .

Theorem 1. *For any instance of our online strategic classification problem with noise level $\bar{\eta}$, maximum manipulation distance γ , and feature dimension d , the expected regret of classifiers $\tilde{\mathbf{h}}$ from Algorithm 1 over T cycles satisfies*

$$\mathbb{E}[\text{Reg}(\tilde{\mathbf{h}}; T)] = O\left(d \ln d \times (\ln T)^2 / (1 - 2\bar{\eta})^8 + \sqrt{(\gamma + 1)d \ln d \times T \ln T / (1 - 2\bar{\eta})^2}\right).$$

We prove the theorem by analyzing the regret incurred by each of the three sub-algorithms in Section 3. The full proof is outlined in Appendix A.6. We also conduct numerical experiments to evaluate our proposed algorithm, with results presented in Appendix A.2.

5 Concluding Remarks

In this paper, we study an online strategic classification problem under Massart Noise with partial feedback. The settings are of practical relevance yet theoretically challenging. We introduce a novel algorithm that concurrently learns a linear classifier and manages instantaneous prediction errors. The algorithm leverages localization to mitigate the complexities induced by Massart noise. The strategic manipulation of agents poses a critical challenge by limiting access to reliable training data; thus, the core innovation of our approach lies in using carefully designed classifier pairs to collect some proxy data and contrasting their data for effective learning. This pairwise contrastive inference approach with proxy data effectively addresses the challenges in online strategic classification. This paper has some limitations. First, our algorithm is specifically designed for Massart Noise. Second, this paper assumes that agents' utility functions are homogeneous and known to the principal. Third, we adopt as an objective the traditional classification accuracy metric. Future research directions include extending the algorithm to overcome these limitations.

Acknowledgement

The authors thank the anonymous review team for their insightful comments and suggestions. Yong Liang acknowledges the support from the National Natural Science Foundation of China (72325001). Xiaojie Mao acknowledges the support from the National Natural Science Foundation of China (72322001 and 72201150) and the National Key R&D Program of China (2022ZD0116700).

References

- Saba Ahmadi, Hedyeh Beyhaghi, Avrim Blum, and Keziah Naggita. The strategic perceptron. In *Proceedings of the 22nd ACM Conference on Economics and Computation*, pages 6–25, 2021.
- Saba Ahmadi, Avrim Blum, and Kunhe Yang. Fundamental bounds on online strategic classification. In *Proceedings of the 24th ACM Conference on Economics and Computation*, pages 22–58, 2023.
- Saba Ahmadi, Kunhe Yang, and Hanrui Zhang. Strategic littlestone dimension: Improved bounds on online strategic classification. *Advances in Neural Information Processing Systems*, 37:101696–101724, 2024.
- Pranjal Awasthi, Maria-Florina Balcan, Nika Haghtalab, and Ruth Uerner. Efficient learning of linear separators under bounded noise. In *Conference on Learning Theory*, pages 167–190. PMLR, 2015.
- Pranjal Awasthi, Maria Florina Balcan, and Philip M Long. The power of localization for efficiently learning linear separators with noise. *Journal of the ACM (JACM)*, 63(6):1–27, 2017.
- Maria-Florina Balcan and Phil Long. Active and passive learning of linear separators under log-concave distributions, 2013.
- Peter L Bartlett, Michael I Jordan, and Jon D McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- Jiafeng Chen, Zhen Li, and Xuan Mao. Learning with selectively labeled data from multiple decision-makers. In *Proceedings of the 42nd International Conference on Machine Learning*, 2025.
- Yiling Chen, Yang Liu, and Chara Podimata. Learning strategy-aware linear classifiers. *Advances in Neural Information Processing Systems*, 33, 2020.
- Lee Cohen, Yishay Mansour, Shay Moran, and Han Shao. Learnability gaps of strategic classification. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 1223–1259. PMLR, 2024.
- Sanjoy Dasgupta. Coarse sample complexity bounds for active learning. *Advances in Neural Information Processing Systems*, 18, 2005.
- Ilias Diakonikolas, Themis Gouleakis, and Christos Tzamos. Distribution-independent pac learning of halfspaces with massart noise. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Ilias Diakonikolas, Vasilis Kontonis, Christos Tzamos, and Nikos Zarifis. Learning halfspaces with massart noise under structured distributions. In *Conference on Learning Theory*, pages 1486–1513. PMLR, 2020.
- Ilias Diakonikolas, Daniel M Kane, Vasilis Kontonis, Christos Tzamos, and Nikos Zarifis. Efficiently learning halfspaces with tsybakov noise. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 88–101, 2021.
- Ilias Diakonikolas, Vasilis Kontonis, Christos Tzamos, and Nikos Zarifis. Online learning of halfspaces with massart noise. *arXiv preprint arXiv:2405.12958*, 2024.
- Jinshuo Dong, Aaron Roth, Zachary Schutzman, Bo Waggoner, and Zhiwei Steven Wu. Strategic classification from revealed preferences. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 55–70, 2018.
- Moritz Hardt, Nimrod Megiddo, Christos Papadimitriou, and Mary Wootters. Strategic classification. In *Proceedings of the 2016 ACM conference on innovations in theoretical computer science*, pages 111–122, 2016.
- Keegan Harris, Chara Podimata, and Steven Z Wu. Strategic apple tasting. *Advances in Neural Information Processing Systems*, 36:79918–79945, 2023.

- David P Helmbold, Nicholas Littlestone, and Philip M Long. Apple tasting. *Information and Computation*, 161(2):85–139, 2000.
- Adam R Klivans, Philip M Long, and Alex K Tang. Baums algorithm learns intersections of half-spaces with respect to log-concave distributions. In *International Workshop on Approximation Algorithms for Combinatorial Optimization*, pages 588–600. Springer, 2009.
- Himabindu Lakkaraju, Jon Kleinberg, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan. The selective labels problem: Evaluating algorithmic predictions in the presence of unobservables. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 275–284, 2017.
- Sagi Levanon and Nir Rosenfeld. Strategic classification made practical. In *International Conference on Machine Learning*, pages 6243–6253. PMLR, 2021.
- László Lovász and Santosh Vempala. The geometry of logconcave functions and sampling algorithms. *Random Structures & Algorithms*, 30(3):307–358, 2007.
- Pascal Massart and Élodie Nédélec. Risk bounds for statistical learning. *The Annals of Statistics*, 34(5):2326–2366, 2006.
- Francesco Orabona. A modern introduction to online learning, 2023.
- Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
- Rocco Anthony Servedio. *Efficient algorithms in computational learning theory*. Harvard University, 2001.
- Shai Shalev-Shwartz. *Online learning: Theory, algorithms, and applications*. Hebrew University, 2007.
- Han Shao, Shuo Xie, and Kunhe Yang. Should decision-makers reveal classifiers in online strategic classification? In *Forty-second International Conference on Machine Learning*, 2025.
- Jie Shen. On the power of localized perceptron for label-optimal learning of halfspaces with adversarial noise. In *International Conference on Machine Learning*, pages 9503–9514. PMLR, 2021a.
- Jie Shen. Sample-optimal pac learning of halfspaces with malicious noise. In *International Conference on Machine Learning*, pages 9515–9524. PMLR, 2021b.
- Jie Shen. Pac learning of halfspaces with malicious noise in nearly linear time. In *International Conference on Artificial Intelligence and Statistics*, pages 30–46. PMLR, 2023.
- Lingqing Shen, Nam Ho-Nguyen, Khanh-Hung Giang-Tran, and Fatma Klnç-Karzan. Mistake, manipulation and margin guarantees in online strategic classification, 2024.
- Ravi Sundaram, Anil Vullikanti, Haifeng Xu, and Fan Yao. Pac-learning for strategic classification. *Journal of Machine Learning Research*, 24(192):1–38, 2023.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- Songbai Yan and Chicheng Zhang. Revisiting perceptron: Efficient and label-optimal learning of halfspaces. *Advances in Neural Information Processing Systems*, 30, 2017.
- Chicheng Zhang and Yinan Li. Improved algorithms for efficient active learning halfspaces with massart and tsybakov noise. In *Conference on Learning Theory*, pages 4526–4527. PMLR, 2021.
- Chicheng Zhang, Jie Shen, and Pranjal Awasthi. Efficient active learning of sparse halfspaces with arbitrary bounded noise. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 7184–7197. Curran Associates, Inc., 2020.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: The abstract accurately reflects the paper's contributions and scope by stating that the paper addresses the online strategic classification problem and proposes a novel learning algorithm that converges to the optimal classifier and achieves a regret rate of $O(\sqrt{T})$ (up to poly-logarithmic and constant factors). It also clearly outlines the combined challenges of agents' strategic feature manipulations, partial label observations, and label noises.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: In Section 5, we acknowledge our limitations and point out some future directions.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: All assumptions are clearly stated in Assumption 1 and Assumption 2. We provide a complete proof in our appendix and provide proof sketches in Section 3 and Section 4.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We use Python 3.9 to conduct our numerical experiments. All settings and results are listed in Appendix A.2. We guarantee that our results are genuine and credible.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The paper provides the code and data required for the experiments in the supplementary material. The code is well-organized and well-documented, which facilitates the reproduction process.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so No is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: All details of our numerical experiment is listed in Appendix A.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We present the mean regret rate of all algorithms, calculated across 10 independent experimental replications in Appendix A.2. The experimental findings are consistent with the theoretical assurances we provide in Section 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We state in Appendix A.2 that all experiments can be conducted locally using a standard CPU without requiring specialized hardware, making reproduction accessible and straightforward.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Yes, the research performed in the paper conforms with the NeurIPS Code of Ethics. The paper does not involve human subjects or sensitive personal information, so issues like privacy and consent are not applicable. The research focuses on developing an algorithm for strategic classification under specific noise conditions, and it does not present any foreseeable risks of harm, discrimination, or other unethical consequences as outlined in the Code of Ethics. Our numerical experiment only uses simulated data, so it does not involve any deprecated datasets or copyright violations.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.

- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [\[Yes\]](#)

Justification: We discuss in Section 1 that agents’ strategic feature manipulation can hurt a certain classification rule, and we design an algorithm to prevent this, which is a potential positive societal impact. This algorithm aims to enhance the fairness and accuracy of automated decision-making systems, potentially reducing financial losses from misclassifications. While the research focuses on foundational aspects and does not directly address all potential negative societal impacts, we acknowledge the importance of considering such implications as the technology evolves.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [\[NA\]](#)

Justification: The research in this paper focuses on developing an algorithm for online strategic classification and does not involve the release of data or models with high misuse risks, such as pretrained language models, image generators, or scraped datasets. Therefore, no specific safeguards for such releases are applicable.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.

- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: The research in this paper focuses on developing a novel algorithm for online strategic classification. We do not use any existing assets such as code, data, or models from external sources. All methods and experiments are designed and implemented by us, so there are no original owners or licenses to credit.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The research in this paper does not release any new assets such as datasets, code, or models. Therefore, no new assets are introduced that would require documentation.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The research in this paper does not involve crowdsourcing experiments or research with human subjects. Therefore, there are no instructions, screenshots, or compensation details to provide.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This research does not involve crowdsourcing or research with human subjects. Therefore, there are no risks to participants, disclosures, or IRB approvals to report.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core methodology of this research does not involve the use of LLMs as any important, original, or non-standard components. LLMs were not utilized in developing the algorithms or conducting the experiments.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.