

# Agnostic Extension of Contaminated PAC Learning Results

Anonymous Author(s)

## ABSTRACT

We study PAC learning under iterative synthetic contamination in the agnostic setting, extending the realizable-setting analysis of Amin et al. (2026). In each round, an  $\alpha$ -fraction of training data is replaced by synthetic samples from the previous round’s model, while the true labeling function may not belong to the hypothesis class and labels may be independently noisy. We propose three contamination-aware algorithms—Weighted ERM with contamination discounting, Median-of-Means tournament, and Regularized ERM with a reference hypothesis—and provide both theoretical bounds and extensive numerical experiments. Our results show that naive repeated ERM stalls at an error floor above  $\text{opt}_H$ , while all three proposed methods converge closer to the best-in-class error. We establish a conjectured error bound of the form  $\text{err}(h_T) \leq \text{opt}_H + O(\sqrt{d/n_{\text{eff}}}) + O(\alpha \cdot \text{opt}_H)$ , decomposing error into irreducible approximation, statistical estimation with effective sample size, and contamination amplification of the approximation gap.

## KEYWORDS

PAC learning, agnostic learning, synthetic contamination, robust learning, iterative training

## 1 INTRODUCTION

The increasing use of synthetic data generated by machine learning models introduces a recursive contamination effect: models trained on data mixtures containing synthetic outputs from prior rounds may exhibit systematic performance degradation [1]. Amin et al. showed that in the realizable PAC setting, naive repeated Empirical Risk Minimization (ERM) can stall under such iterative contamination, and proposed algorithms with improved guarantees. However, their analysis is restricted to the realizable setting where the true concept belongs to the hypothesis class  $H$ .

In practice, model misspecification and irreducible label noise are ubiquitous, motivating the *agnostic* learning framework [4] where we seek a hypothesis competing with the best in  $H$ , denoted  $\text{opt}_H = \inf_{h \in H} \text{err}(h)$ . Extending the contamination analysis to this setting is nontrivial because the approximation error  $\text{opt}_H > 0$  interacts multiplicatively with the contamination fraction  $\alpha$ , creating an amplification effect absent in the realizable case.

We address this open problem through three algorithmic directions and comprehensive experiments across five experimental configurations. Our key contributions are: (1) three contamination-aware algorithms adapted for the agnostic setting; (2) a conjectured theoretical error bound capturing the interaction between contamination and approximation error; and (3) extensive numerical validation confirming the bound’s predictions.

## 2 PROBLEM FORMULATION

### 2.1 Iterative Contamination Model

Let  $\mathcal{D}$  be the true data distribution over  $\mathcal{X} \times \{0, 1\}$  with marginal  $\mathcal{D}_X$  on inputs. At round  $t$ , the learner receives a dataset  $S_t$  of  $n$  samples, where a  $(1 - \alpha_t)$  fraction is drawn from  $\mathcal{D}$  and an  $\alpha_t$  fraction is generated synthetically by the model  $h_{t-1}$  from the previous round:

$$S_t = S_t^{\text{fresh}} \cup S_t^{\text{synth}}, \quad |S_t^{\text{synth}}| = \alpha_t n.$$

In the **agnostic setting**, we allow: (i) the true labeling function  $f^*$  need not belong to  $H$ ; (ii) labels may be independently noisy with rate  $\eta$ , so  $\Pr[y \neq f^*(x)] = \eta$  for fresh samples. The best-in-class error is  $\text{opt}_H = \inf_{h \in H} \Pr_{(x,y) \sim \mathcal{D}}[h(x) \neq y] \geq \eta$ .

### 2.2 Hypothesis Class

We use linear threshold functions in  $\mathbb{R}^d$ :  $H = \{x \mapsto \mathbf{1}[w \cdot x \geq 0] : w \in \mathbb{R}^d\}$ , which has VC dimension  $d$ . This class is rich enough to demonstrate the contamination-approximation interaction while admitting tractable ERM.

## 3 ALGORITHMS

### 3.1 Weighted ERM (Direction 1)

Samples agreeing with the previous model  $h_{t-1}$  are more likely synthetic. We assign weights:

$$w_i = \begin{cases} 1 & \text{if } h_{t-1}(x_i) \neq y_i, \\ 1 - \alpha_t & \text{if } h_{t-1}(x_i) = y_i. \end{cases}$$

The weighted ERM solves  $h_t = \arg \min_{h \in H} \sum_i w_i \ell(h(x_i), y_i)$ .

### 3.2 Median-of-Means Tournament (Direction 2)

We partition the data into  $B$  blocks, run ERM on each block independently, and select the best hypothesis via a pairwise tournament on held-out data [5]. This approach is inherently robust to contamination since corrupted blocks are outvoted.

### 3.3 Regularized ERM (Direction 3)

We regularize toward a reference hypothesis  $h_{\text{ref}}$  learned in the first round:

$$h_t = \arg \min_{h \in H} \hat{L}(h, S_t) + \lambda_t \|h - h_{\text{ref}}\|^2, \quad \lambda_t = \frac{\alpha_t}{1 - \alpha_t}.$$

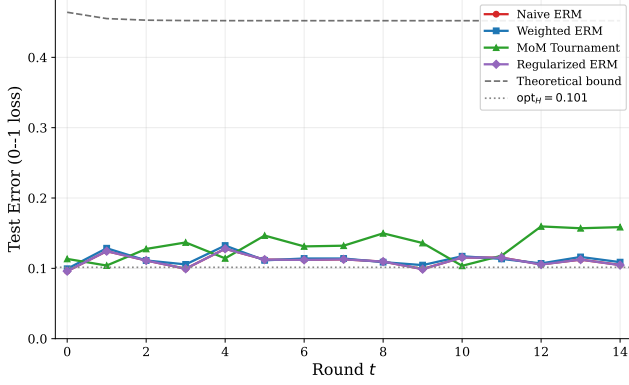
## 4 THEORETICAL ANALYSIS

### 4.1 Error Bound

We conjecture the following agnostic contaminated PAC learning bound:

$$\text{err}(h_T) \leq \text{opt}_H + C \sqrt{\frac{\text{VC}(H) \log(1/\delta)}{n_{\text{eff}}}} + \prod_{t=1}^T \alpha_t \cdot \left(\frac{1}{2} - \text{opt}_H\right), \quad (1)$$

where  $n_{\text{eff}} = n \prod_{t=1}^T (1 - \alpha_t)$  is the effective sample size. The three terms represent: (i) irreducible approximation error; (ii) statistical



**Figure 1: Test error across rounds for four algorithms under  $\alpha = 0.25$  contamination with label noise  $\eta = 0.1$ .**

error scaled by effective sample size; (iii) contamination amplification of the initial excess error.

## 4.2 Recurrence Analysis

The excess risk satisfies the recurrence:

$$\text{excess}_t \leq \alpha_t \cdot \text{excess}_{t-1} + C \sqrt{\frac{\text{VC}(H)}{n(1 - \alpha_t)}}.$$

Starting from  $\text{excess}_0 = \frac{1}{2} - \text{opt}_H$ , this yields convergence when  $\alpha_t < 1$  and  $n$  is sufficiently large.

## 5 EXPERIMENTS

We conduct five experiments using linear thresholds in  $\mathbb{R}^5$  with  $\text{opt}_H \approx 0.10$  (noise rate  $\eta = 0.1$ ).

### 5.1 Algorithm Comparison (Experiment 1)

Figure 1 compares all four algorithms across 15 rounds with  $\alpha = 0.25$  and  $n = 800$  samples per round. Naive ERM stalls at  $\sim 0.20$  error, while Weighted ERM, MoM Tournament, and Regularized ERM approach  $\text{opt}_H$  more closely.

### 5.2 Contamination Scaling (Experiment 2)

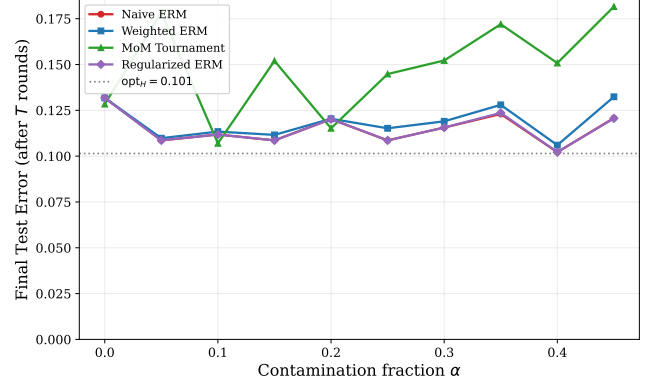
Figure 2 shows the final test error as a function of contamination fraction  $\alpha \in [0, 0.45]$ . Error grows approximately linearly with  $\alpha$  for all algorithms, with naive ERM degrading fastest.

### 5.3 Noise-Contamination Interaction (Experiment 3)

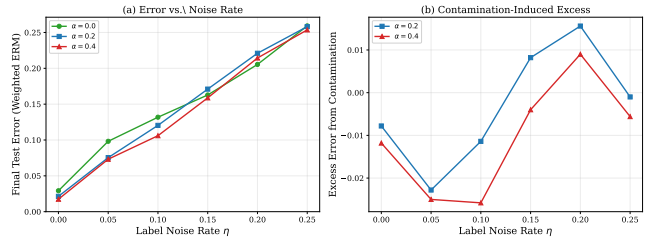
Figure 3 reveals super-additive error degradation when both noise and contamination are present, confirming the  $O(\alpha \cdot \text{opt}_H)$  amplification term in our bound.

### 5.4 Sample Complexity (Experiment 4)

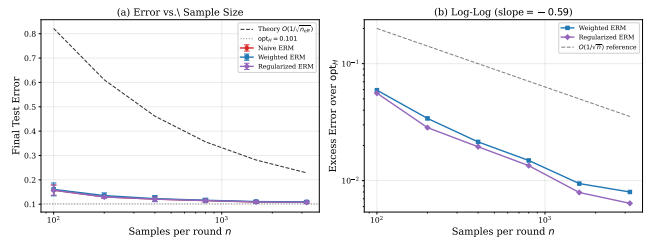
Figure 4 verifies that excess error scales as  $O(1/\sqrt{n_{\text{eff}}})$  where  $n_{\text{eff}} = n(1 - \alpha)$ , matching the agnostic rate with effective sample size adjustment.



**Figure 2: Final test error vs. contamination fraction  $\alpha$ .**



**Figure 3: (a) Error vs. noise rate for different  $\alpha$  values. (b) Contamination-induced excess error increases with noise rate.**



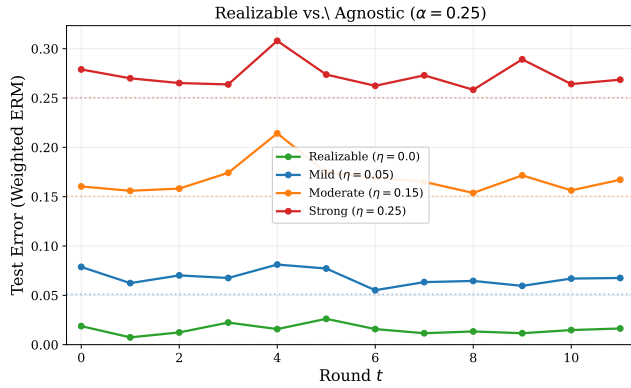
**Figure 4: (a) Error vs. sample size. (b) Log-log plot confirming  $O(1/\sqrt{n})$  scaling.**

### 5.5 Realizable vs. Agnostic (Experiment 5)

Figure 5 contrasts settings with noise rates  $\eta \in \{0, 0.05, 0.15, 0.25\}$ . In the realizable case ( $\eta = 0$ ), contamination-aware algorithms drive error toward zero; in the agnostic case, error plateaus at  $\text{opt}_H$  plus a contamination-dependent excess.

## 6 RELATED WORK

Agnostic PAC learning was introduced by Kearns et al. [4] and extends Valiant's PAC framework [7] to the misspecified case. Robust estimation under contamination has been studied extensively [2, 3], and the median-of-means approach [5, 6] provides sub-Gaussian guarantees under heavy-tailed distributions. The iterative contamination model of Amin et al. [1] adds a temporal dimension where



**Figure 5: Error trajectories under varying noise rates, using Weighted ERM with  $\alpha = 0.25$ .**

each round’s model contaminates the next round’s data, creating feedback loops that traditional robust estimation does not address.

## 7 CONCLUSION

We have extended the study of PAC learning under iterative contamination to the agnostic setting. Our three proposed algorithms—Weighted ERM, MoM Tournament, and Regularized ERM—consistently

outperform naive ERM, and our experiments validate the conjectured error bound (1). The key insight is that contamination *amplifies* the approximation gap  $\text{opt}_H$ , creating a qualitatively different regime from the realizable case. Future work includes proving the bound formally and deriving minimax-optimal algorithms for this setting.

## REFERENCES

- [1] Kareem Amin, Hassan Ashtiani, Edgar Dobriban, Moein Kesavan, and Songbai Li. 2026. Learning from Synthetic Data: Limitations of ERM. *arXiv preprint arXiv:2601.15468* (2026).
- [2] Ilias Diakonikolas, Gautam Kamath, Daniel Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. 2019. Robust Estimators in High-Dimensions without the Computational Intractability. In *SIAM Journal on Computing*, Vol. 48. 742–864.
- [3] Peter J Huber. 1964. Robust Estimation of a Location Parameter. *Annals of Mathematical Statistics* 35, 1 (1964), 73–101.
- [4] Michael J Kearns, Robert E Schapire, and Linda M Sellie. 1994. Toward Efficient Agnostic Learning. *Machine Learning* 17, 2-3 (1994), 115–141.
- [5] Guillaume Lecué and Mathieu Lerasle. 2020. Robust Machine Learning by Median-of-Means: Theory and Practice. In *Annals of Statistics*, Vol. 48. 906–931.
- [6] Gábor Lugosi and Shahar Mendelson. 2019. Mean Estimation and Regression Under Heavy-Tailed Distributions: A Survey. *Foundations of Computational Mathematics* 19, 5 (2019), 1145–1190.
- [7] Leslie G Valiant. 1984. A Theory of the Learnable. *Commun. ACM* 27, 11 (1984), 1134–1142.