

---

# Conformal Margin Risk Minimization: An Envelope Framework for Robust Learning Under Label Noise

---

Yuanjie Shi\*

Peihong Li\*

Zijian Zhang

Jana Doppa

Yan Yan

School of EECS, Washington State University, Pullman, WA, USA

## Abstract

Most methods for learning with noisy labels require privileged knowledge such as noise transition matrices, clean subsets or pretrained feature extractors, resources typically unavailable when robustness is most needed. We propose **Conformal Margin Risk Minimization (CMRM)**, a plug-and-play envelope framework that improves *any* classification loss under label noise by adding a single quantile-calibrated regularization term, with no privileged knowledge or training pipeline modification. CMRM measures the confidence margin between the observed label and competing labels, and thresholds it with a conformal quantile estimated per batch to focus training on high-margin samples while suppressing likely mislabeled ones. We derive a learning bound for CMRM under arbitrary label noise requiring only mild regularity of the margin distribution. Across five base methods and six benchmarks with synthetic and real-world noise, CMRM consistently improves accuracy (up to +3.39%), reduces conformal prediction set size (up to -20.44%) and does not hurt under 0% noise, showing that CMRM captures a method-agnostic uncertainty signal that existing mechanisms did not exploit.

## 1 INTRODUCTION

Deep neural networks have achieved remarkable success across domains such as vision (Zhao et al., 2024; Zhang et al., 2024), language (Naveed et al., 2025; Teubner

et al., 2023), and healthcare (Thirunavukarasu et al., 2023; Celard et al., 2023), but their performance typically depends on clean labels (Arpit et al., 2017). However, in real-world scenarios, labels are often corrupted by human annotation errors (e.g., crowd-sourcing) or automated data collection pipelines (Song et al., 2022; Gupta and Gupta, 2019). Learning with label noise (LNL) formalizes this setting, where corrupted labels distort empirical risk and lead to overfitting (Natarajan et al., 2013; Zhang et al., 2021; Arpit et al., 2017). This issue is particularly severe in high-stakes domains such as medical imaging (Shi et al., 2024a) and autonomous driving (Li et al., 2021a), making robustness to label noise an important problem in modern machine learning (Song et al., 2022; Frénay and Verleysen, 2013).

Learning under label noise has been extensively studied (Song et al., 2022), with existing approaches broadly falling into the following categories. Noise modeling approaches assume a specific corruption process (e.g., symmetric, class-conditional, or instance-dependent) and often require estimating or knowing the noise transition matrix (Natarajan et al., 2013; Patrini et al., 2017; Li et al., 2021b; Cheng et al., 2020), which is typically unobservable and difficult to model under complex noise (Yao et al., 2020). Loss correction methods adjust training objectives to counteract label corruption but typically depend on auxiliary information, such as small trusted clean subsets or accurate noise rate estimates (Hendrycks et al., 2018; Xia et al., 2019), both of which are often unavailable or unreliable in practice. Auxiliary methods leverage additional sources, such as peer networks, semi-supervised learning (SSL) pipelines, or external expert models such as large language models (LLMs), which introduce strong data and supervision requirements (Han et al., 2018; Li et al., 2020; Wang et al., 2024). Table 1 summarizes representative methods and their key assumptions (see Section 3 for details).

Despite this progress, these approaches share a common limitation: they all require some form of privileged knowledge about the noise process or access to aux-

---

\*Equal contribution.

Methods	Assumptions
LNL (Natarajan et al., 2013)	Symmetric noise
Forward (Patrini et al., 2017)	Known label noise transition matrix
GLC (Hendrycks et al., 2018)	Small clean data subsets
Co-teaching (Han et al., 2018)	Two peer nets
CORES (Cheng et al., 2020)	Instance-dependent noise
VolMinNet (Li et al., 2021b)	Class-conditional noise
FINE (Kim et al., 2021)	Eigen-structure noise
ELR+ (Liu et al., 2020)	Clean samples learned first
NCFW (Zhang and Agarwal, 2024)	Known class posteriors
CSGN (Lin et al., 2024)	Latent causal graph transition
NI-ERM (Zhu et al., 2024)	Strong pretrained feature extractor
NoiseGPT (Wang et al., 2024)	External LLM
<b>CMRM (ours)</b>	Smooth CDF + positive density

Table 1: **Classical and recent methods for learning from noisy labels and their assumptions.** CMRM only assumes mild regularity conditions, a smooth CDF with positive density at the target quantile. Details in Section 3. CMRM’s assumptions are automatically satisfied in standard training (see Figure 2(d)).

iliary resources, such as a noise transition matrix, a clean data subset, a peer network or a pretrained feature extractor, which is typically unavailable in the settings where robustness to label noise is most needed, especially under severe or heterogeneous noise (Arazo et al., 2019; Song et al., 2022). This raises the main research question of this paper: *Can we design a flexible envelope framework that enhances the robustness of existing methods under arbitrary label noise, requiring only standard mathematical regularity conditions, instead of any knowledge of the noise process or access to auxiliary supervision?*

To answer this question, we propose **Conformal Margin Risk Minimization (CMRM)**, an envelope framework to improve the robustness and accuracy of prior methods. CMRM does not rely on noise models, clean subsets or auxiliary supervision. CMRM achieves the goal through an uncertainty-aware training objective that integrates confidence margins and conformal quantiles. Confidence margin, defined as the gap between the confidence of the observed label and competing labels, provides a principled signal to distinguish reliable from uncertain training samples (Cui et al., 2019; Pleiss et al., 2020). Conformal quantiles offer distribution-free, statistically valid thresholds for these margins (Angelopoulos and Bates, 2021; Lei et al., 2018).

This algorithm design is motivated by two prior observations: (1) label noise creates a mismatch between corrupted labels and the evolving confidence structure of the model (Pleiss et al., 2020), and (2) ignoring this uncertainty often leads to overfitting on noisy samples and prevents algorithms from leveraging informative signals (Zhang et al., 2021; Arpit et al., 2017). Building on this motivation, CMRM formulates a conformal margin risk that directly incorporates uncertainty into the training objective. Specifically, CMRM computes per-example confidence margins, estimates a conformal

quantile to set an adaptive threshold, and minimizes a conformal risk defined as the average negative margin above that threshold. Importantly, CMRM requires no architectural modification, no additional networks and no clean validation data, as it adds only a single quantile-calibrated regularization term to any existing training objective. We further establish a learning bound for CMRM under arbitrary noise. Our extensive experiments on both binary and multi-class classification benchmarks demonstrate that, without requiring prior knowledge of noise, CMRM consistently improves methods with fundamentally different LNL design principles, indicating that conformal margin calibration captures a method-agnostic uncertainty signal that existing mechanisms fail to fully exploit.

**Contributions.** Our key contributions include:

- We propose CMRM, an envelope framework whose robustness requires only mild regularity of the margin distribution that is automatically satisfied in standard training, instead of assumptions on noise models, clean data or auxiliary resources (see Table 1).
- We derive a learning bound under arbitrary label noise that decomposes the gap between the noisy surrogate and clean conditional margin risk into function-class complexity, quantile estimation error, and distribution shift.
- We demonstrate that CMRM, requiring only a single regularization term, improves accuracy across all base method and dataset combinations tested and reduces prediction set size in nearly all cases, across five base methods (CE, Focal, LDAM, GCE, NI-ERM) and six benchmarks with synthetic and real-world noise, and incurs no accuracy penalty when labels are clean. The CMRM code is available at <https://github.com/YuanjieSh/CMRM>.

## 2 RELATED WORK

**Learning with Label Noise (LNL)** studies how to train models when labels are corrupted by annotation errors, heuristic labeling, or unreliable sources (Frénay and Verleysen, 2013; Song et al., 2022; Johnson and Khoshgoftaar, 2022). Many approaches have been proposed, often under specific assumptions on the noise mechanism or requiring auxiliary supervision (see Table 1 and Section 3). These include loss or posterior correction with known noise models (Natarajan et al., 2013; Patrini et al., 2017; Li et al., 2021b), methods using clean side information or peer networks (Hendrycks et al., 2018; Han et al., 2018; Li et al., 2020; Liu et al., 2020), approaches for instance-dependent or structured

noise (Cheng et al., 2020; Kim et al., 2021; Lin et al., 2024), and methods leveraging external knowledge such as pretrained features or LLMs (Zhang and Agarwal, 2024; Zhu et al., 2024; Wang et al., 2024). In contrast, our envelope approach requires only mild regularity conditions and is compatible with most prior methods.

**Conformal Prediction (CP)** is a distribution-free uncertainty quantification framework that constructs prediction sets with guaranteed finite-sample coverage (Vovk et al., 2005; Lei et al., 2018; Angelopoulos and Bates, 2021; Fontana et al., 2023; Ghosh et al., 2023a,b; Shi et al., 2024b; Shahrokhi et al., 2025). Recent conformal training methods incorporate these principles into the training process, focusing on minimizing prediction set size (Stutz et al., 2021; Shi et al., 2025) or refining coverage calibration (Einbinder et al., 2022; Kiyani et al., 2024), but assume clean training data. In contrast, CMRM leverages conformal principles to enhance discriminative ability under noisy supervision. Importantly, CMRM is not a conformal prediction method in the standard sense; it does not construct prediction sets or provide coverage guarantees at inference time.

### 3 PROBLEM SETUP AND MOTIVATION

**Notations.** Suppose  $X \in \mathcal{X}$  is an input from  $\mathcal{X}$ , and  $Y \in \mathcal{Y} = \{0, 1, \dots, K-1\}$  is the ground-truth label, where  $K$  is the number of candidate classes. Let the underlying data distribution be  $\mathcal{P}(X, Y)$ , which characterizes the relationship between inputs and class labels. We denote by  $\mathcal{P}(X)$  the marginal distribution of inputs, and by  $\mathcal{P}(Y|X)$  the conditional label distribution given inputs. Let  $f : \mathcal{X} \rightarrow \mathbb{R}^K$  denote the logit vector produced by soft classifier  $f \in \mathcal{F}$ , where  $\mathcal{F}$  denotes a hypothesis class. We also define  $P_f(X) = \sigma \circ f(X) : \mathcal{X} \rightarrow \Delta_+^K$  as the corresponding confidence score, where  $\Delta_+^K$  is the  $(K-1)$ -dimensional probability simplex, and  $\sigma$  is the Sigmoid function in binary and Softmax function in multi-class setting, respectively. Let  $P_f(X)_y$  denote the confidence score of class  $y$ . Define  $\mathbb{1}[\cdot]$  as an indicator function. Denote  $\mathcal{D}_{\text{tr}}$  and  $\mathcal{D}_{\text{test}}$  as the training and test sets. Let  $\mathcal{B}$  be a randomly sampled batch of training data of size  $s$ . We denote by  $W_1(\mathcal{P}, \mathcal{Q}) = \inf_{\pi \in \Pi(\mathcal{P}, \mathcal{Q})} \int_{\mathcal{X} \times \mathcal{X}} d(x, x') d\pi(x, x')$  as the Wasserstein-1 distance between two distributions  $\mathcal{P}$  and  $\mathcal{Q}$  on metric space  $(\mathcal{X}, d)$ , where  $\Pi(\mathcal{P}, \mathcal{Q})$  is the set of all joint distributions of  $\mathcal{P}$  and  $\mathcal{Q}$ .

**Learning with Noisy Labels (LNL).** In many real-world scenarios, the training data are corrupted by noisy labels. Instead of the clean label  $Y$ , we only observe a potentially corrupted label  $\tilde{Y} \in \mathcal{Y}$  generated from  $Y$  through an unknown label noise transition matrix  $T(\tilde{Y} | X, Y)$  (Zhu et al., 2024). For example,

symmetric noise corresponds to  $T(\tilde{Y} | X, Y)$  being uniform over all incorrect labels, while class-conditional noise assumes dependence only on  $Y$ . We highlight that *our proposed CMRM framework does not rely on such assumptions and permits  $T$  to be arbitrary.*

Such noisy labels setting arises typically due to annotator mistakes, inconsistent labeling criteria, or spurious labels introduced by large-scale data collection pipelines (Natarajan et al., 2013; Patrini et al., 2017). Formally, we define the noisy training data distribution as:

$$\mathcal{P}_{\text{noisy}}(X, \tilde{Y}) = \sum_{y \in \mathcal{Y}} \mathcal{P}(X, y) T(\tilde{Y} | X, y).$$

Accordingly, the training set  $\mathcal{D}_{\text{tr}} = \{(X_i, \tilde{Y}_i)\}_{i=1}^n$  of size  $n$  is drawn from  $\mathcal{P}_{\text{noisy}}$ , while the clean test set  $\mathcal{D}_{\text{test}}$  is drawn from  $\mathcal{P}$ .

**Limitations of Existing Approaches.** Despite substantial progress on learning from noisy labels, most existing methods remain fundamentally constrained by strong assumptions. As summarized in Table 1, these assumptions fall into three major categories:

First, many approaches impose explicit noise-model assumptions, from symmetric noise or known transition matrices (Natarajan et al., 2013; Patrini et al., 2017) to class-conditional (Li et al., 2021b), instance-dependent (Cheng et al., 2020), or eigen-structure noise (Kim et al., 2021), but these still rarely align with real-world corruption.

Second, a range of methods rely on auxiliary information, such as clean subsets (Hendrycks et al., 2018), noise transition matrices (Patrini et al., 2017), strong pretrained feature extractors (Zhu et al., 2024), or external LLMs (Wang et al., 2024), which are costly and unavailable at scale.

Third, some approaches exploit architectural assumptions or model-specific heuristics, such as peer networks (Han et al., 2018; Li et al., 2020), early-learning assumptions (Liu et al., 2020), or latent causal structures (Lin et al., 2024).

Overall, these assumptions limit existing methods, especially under severe or heterogeneous noise (Arazo et al., 2019; Song et al., 2022). Therefore, an envelope framework that provides robustness guarantees under only mild regularity conditions is needed.

## 4 CMRM FRAMEWORK

In this section, we first describe the general *Conformal Margin Risk Minimization (CMRM)* envelope framework and its variant for binary classification. Next, we develop a practical optimization algorithm for CMRM and theoretically analyze its learning bound.

#### 4.1 General Framework

**Confidence Margin.** The confidence margin quantifies the separation between the confidence assigned to the observed label and the highest confidence among other candidate labels:

$$M_f(X_i, \tilde{Y}_i) = P_f(X_i)_{\tilde{Y}_i} - \max_{y \in \mathcal{Y} \setminus \{\tilde{Y}_i\}} P_f(X_i)_y. \quad (1)$$

This notion of confidence margin has been widely used in the machine learning literature (Cui et al., 2019; Bartlett and Mendelson, 2002). Large margins indicate a strong preference for the observed label, whereas small or negative margins reflect uncertainty (Cui et al., 2019). Margins are widely used to characterize decision boundary distance (Bartlett and Mendelson, 2002; Neyshabur et al., 2017) and to measure predictive uncertainty (Elsayed et al., 2018; Jiang et al., 2018; Lakshminarayanan et al., 2017).

**Conformal Quantile.** Conformal prediction (CP) (Vovk et al., 2005; Angelopoulos and Bates, 2021) provides a distribution-free mechanism for defining quantile thresholds. Let  $V$  be a real-valued random variable, and let  $\{V_i\}_{i=1}^m$  denote an i.i.d. sample from  $V$  with size  $m$ . Given a target level  $\alpha \in (0, 1)$ , the empirical conformal quantile is defined as  $\hat{\tau}_\alpha = Q(\alpha, \{V_i\}_{i=1}^m)$ , where  $Q$  selects the  $\lceil \alpha(m+1) \rceil$ -th largest value in  $\{V_i\}_{i=1}^m$ . Hence, at most an  $\alpha$ -fraction of the sample values lie below  $\hat{\tau}_\alpha$ , i.e., the empirical  $\alpha$ -quantile. Through CP, this quantile enjoys distribution-free validity guarantee, making it a principled tool for UQ.

**Conformal Margin Risk.** While conformal quantiles provide principled thresholds, their effectiveness depends on using a score that reflects label reliability under noise. Confidence margin serves this role: clean samples typically exhibit large margins, whereas noisy labels tend to have small margins (Zhang and Sabuncu, 2018; Liu et al., 2020; Zhang et al., 2021). This is intuitive: correct labels align with dominant class evidence, while corrupted labels conflict with the input features and are outscored by alternative labels. Figure 1 illustrates this separation on CIFAR-100 under different types of noise (class-conditional and human annotation noise). However, such observations have mainly been used for diagnostic analysis or heuristic filtering, without formal risk formulations.

CMRM closes this gap by calibrating confidence margins with conformal quantiles to yield an assumption-light and uncertainty-aware risk. Specifically, we instantiate  $V$  with confidence margin  $M_f(X, \tilde{Y})$ , and define the conformal quantile threshold on  $\{M_f(X_i, \tilde{Y}_i)\}_{i=1}^n$ :

$$\hat{\tau}_\alpha(f) = Q\left(\alpha, \{M_f(X_i, \tilde{Y}_i)\}_{i=1}^n\right), \quad (2)$$

so that at most an  $\alpha$ -fraction of samples fall below

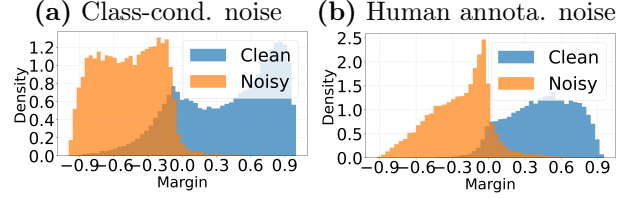


Figure 1: **Confidence margin distributions for clean (blue) and noisy (orange) samples on CIFAR-100.** (a) Class-conditional noise at 20% and (b) human annotation noise at 40%. In both cases, clean samples concentrate on positive margins while noisy samples shift to negative, showing that confidence margins can distinguish clean from noisy labels without assumptions on the noise process.

this threshold, a property that always holds without assumptions on the noise process.  $\alpha$  is a hyper-parameter tuned on validation set. Larger  $\alpha$  filters more aggressively, improving noise robustness but risking loss of clean samples, while smaller  $\alpha$  retains more signal but with corruption.  $\alpha$  is a hyper-parameter selected through validation data (see sensitivity study in Section 5.1).

Next, we define the conformal margin risk  $\hat{\mathcal{L}}_{\text{cr}}(f)$  as:

$$\hat{\mathcal{L}}_{\text{cr}}(f) = \frac{1}{n} \sum_{i=1}^n \hat{\ell}_{\text{cr}}(f, X_i, \tilde{Y}_i), \quad \text{such that} \quad (3)$$

$$\hat{\ell}_{\text{cr}}(f, X_i, \tilde{Y}_i) = -M_f(X_i, \tilde{Y}_i) \cdot \tilde{\mathbb{1}}[M_f(X_i, \tilde{Y}_i) \geq \hat{\tau}_\alpha(f)],$$

where  $\hat{\ell}_{\text{cr}}(f, X_i, \tilde{Y}_i)$  is the per-sample conformal margin loss of  $(X_i, \tilde{Y}_i)$ , and  $\tilde{\mathbb{1}}[x \geq y] = 1/(1 + \exp(-(x - y)/\text{temp}))$  is the smoothed indicator by Sigmoid function with temperature parameter  $\text{temp}$ . This formulation assigns a soft weight in  $(0, 1)$  to each sample rather than hard 0/1 filtering, smoothly downweighting low-margin samples during training, similar to soft reweighting strategies in Tjandra and Wiens (2023).

By construction,  $\hat{\mathcal{L}}_{\text{cr}}(f)$  is the average negative margin over the selected samples, i.e., those with margins above the threshold. Minimizing  $\hat{\mathcal{L}}_{\text{cr}}(f)$  increases the separation on selected high-confidence samples while discarding low-confidence samples. Thus, this principle aligns training with predictive uncertainty and filters out potentially corrupted labels (empirically verified in Figure 2(b)), without assumptions on the noise model or auxiliary supervision.

Finally, we incorporate CMRM as a plug-and-play regularizer for any classification loss  $\hat{\mathcal{L}}_{\text{cl}}(f)$  (e.g., cross-entropy), including prior LNL loss functions, as demonstrated in our empirical evaluation:

$$\hat{\mathcal{L}}(f) = \hat{\mathcal{L}}_{\text{cl}}(f) + \lambda \cdot \hat{\mathcal{L}}_{\text{cr}}(f), \quad (4)$$

where  $\lambda \geq 0$  controls the strength of the CMRM term relative to  $\hat{\mathcal{L}}_{\text{cl}}(f)$ . The above formulation is general

for multi-class classification problems with standard accuracy metrics. However, binary classification tasks require optimizing more specific metrics. Hence, we provide an instantiation for binary classification below.

**CMRM for Binary Classification.** Binary classification often relies on class-conditional performance measures such as false positive and false negative rates (FPR and FNR), which allow users to impose class-conditional tolerances (Zhou and Liu, 2006). Since CMRM’s quantile-based formulation naturally extends to class-conditional settings, we introduce a binary variant that adapts its margin thresholding to each class separately, enabling asymmetric error control while preserving its core structure.

First, we define class-conditional quantile thresholds for tolerances  $\alpha^+$  and  $\alpha^-$ :

$$\hat{\tau}^-(f) = \min \left\{ t : \frac{1}{n_0} \sum_{i: \tilde{Y}_i=0} \mathbb{1}[P_f(X_i)_1 \geq t] \leq \frac{\lceil \alpha^-(n_0+1) \rceil}{n_0} \right\},$$

$$\hat{\tau}^+(f) = \max \left\{ t : \frac{1}{n_1} \sum_{i: \tilde{Y}_i=1} \mathbb{1}[P_f(X_i)_1 \leq t] \leq \frac{\lceil \alpha^+(n_1+1) \rceil}{n_1} \right\},$$

where  $n_0$  and  $n_1$  are the numbers of observed negative and positive samples, respectively. These thresholds control the upper tail of the negative class (potential false positives) and the lower tail of the positive class (potential false negatives).

Next, we define a two-sided hinge formulation relative to these thresholds:

$$\hat{\ell}_{\text{cr}}^{\text{bin}}(f) = \frac{1}{n} \sum_{i=1}^n \hat{\ell}_{\text{cr}}^{\text{bin}}(f, X_i, \tilde{Y}_i), \text{ such that} \quad (5)$$

$$\hat{\ell}_{\text{cr}}^{\text{bin}}(f, X_i, \tilde{Y}_i) = -\lambda^- \tilde{\mathbb{1}}[\tilde{Y}_i = 0] \cdot (P_f(X_i)_1 - \hat{\tau}^-(f))^+ - \lambda^+ \tilde{\mathbb{1}}[\tilde{Y}_i = 1] \cdot (\hat{\tau}^+(f) - P_f(X_i)_1)^+,$$

where  $(z)^+ = \max\{0, z\}$  and  $\lambda^+, \lambda^- > 0$  control the relative strength of the two penalties.

This binary classification variant preserves the two core ingredients of CMRM: conformal quantile-based filtering and margin maximization on retained samples. By separately controlling the class-conditional distribution, it aligns naturally with standard binary performance metrics such as FPR and FNR.

## 4.2 Optimization and Theoretical Analysis

**Optimization Algorithm.** A key challenge in CMRM is that its objective depends on set-level quantiles of confidence margins over the full training data distribution, requiring  $O(n \log n)$  sorting per iteration and thus prohibitive for large-scale training. We address this by replacing set-level quantiles with batch-level quantiles  $\hat{\tau}_\alpha^s(f)$ : each iteration estimates the

---

### Algorithm 1 Conformal Margin Risk Minimization

---

- 1: **Input:** training dataset  $\mathcal{D}_{\text{tr}}$ , regularization parameter  $\lambda$ , batch size  $s$ , learning-rate  $\eta > 0$ , exclusion rate  $\alpha$
  - 2: Randomly initialize the deep neural network  $f_0$
  - 3: **for**  $t \leftarrow 0 : T - 1$  **do**
  - 4: Randomly sample batch  $\mathcal{B}_t \subset \mathcal{D}_{\text{tr}}$
  - 5: Compute  $M_{f_t}(X_i, \tilde{Y}_i)$  on  $\mathcal{B}_t$
  - 6: Compute batch-wise quantile  $\hat{\tau}_\alpha^s(f_t)$  on  $\mathcal{B}_t$
  - 7: Compute classification loss  $\hat{\mathcal{L}}_{\text{cl}}(f_t)$  on  $\mathcal{B}_t$
  - 8: Compute conformal margin risk  $\hat{\mathcal{L}}_{\text{cr}}(f_t)$  on  $\mathcal{B}_t$
  - 9:  $f_{t+1} \leftarrow f_t - \eta \nabla_f (\hat{\mathcal{L}}_{\text{cl}}(f_t) + \lambda \hat{\mathcal{L}}_{\text{cr}}(f_t))$
  - 10: **end for**
  - 11: **Output:** the trained model  $f_T$
- 

quantile from a batch of size  $s$ , reducing the cost to  $O(s \log s)$ . This yields a tractable surrogate objective  $\hat{\mathcal{L}}_{\text{cr}}^s(f)$  by using  $\hat{\tau}_\alpha^s(f)$ :

$$\hat{\mathcal{L}}_{\text{cr}}^s(f) = \mathbb{E}_{\hat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} \left[ \sum_{i=1}^n -M_f(X_i, \tilde{Y}_i) \cdot \tilde{\mathbb{1}}[M_f(X_i, \tilde{Y}_i) \geq \hat{\tau}_\alpha^s(f)] \right], \quad (6)$$

where  $\mathcal{T}_\alpha(f)$  is the underlying distribution of  $\hat{\tau}_\alpha^s(f)$ . Although this introduces an approximation gap relative to population quantiles, our theoretical analysis shows that the gap is bounded.

Algorithm 1 summarizes the CMRM training procedure. The algorithm follows a standard stochastic optimization loop. For each iteration  $t$ , we first compute the confidence margins on batch  $\mathcal{B}_t$  (Line 5) and then estimate the batch-wise conformal quantile threshold (Line 6). Next, we compute both the classification loss (Line 7) and the conformal margin risk (Line 8). The final update (Line 9) jointly minimizes both terms. Importantly, CMRM requires no change to model architecture or optimization, introducing only a quantile-calibrated regularizer that is compatible with arbitrary classifiers and loss functions.

**Theoretical Analysis.** Building on the optimization procedure, we analyze the effect of using batch-level estimates and the corresponding learning bound. Proposition 1 shows that the batch-level quantile in  $\hat{\mathcal{L}}_{\text{cr}}^s(f)$  concentrates around the population quantile at rate  $\tilde{O}(1/\sqrt{s})$  where  $s$  is the batch size. Theorem 1 further bounds the gap between empirical margin risk on noisy data and population risk on clean data. Together, these results characterize the robustness of CMRM under arbitrary label noise with minimal assumptions.

We start with the definition of population quantile

$\tau_\alpha(f)$  on the noisy data distribution:

$$\tau_\alpha(f) = \min \left\{ t : \mathbb{P}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \leq t] \geq \alpha \right\}.$$

Next, we analyze the gap between the batch-level quantile and population one.

**Proposition 1** (Gap between  $\tau_\alpha$  and  $\hat{\tau}_\alpha^s$ ). *Denote by  $G(t)$  the cumulative distribution function (CDF) of  $M_f(X, \tilde{Y})$  under the noisy distribution  $\mathcal{P}_{\text{noisy}}$ . Assume that  $G(t)$  is continuously differentiable in a neighborhood of  $\tau_\alpha(f)$  with density  $g(t)$  and  $g(\tau_\alpha(f)) > 0$ . Then, for any  $\delta \in (0, 1)$ , we have:*

$$\mathbb{P} \left( |\tau_\alpha(f) - \hat{\tau}_\alpha^s(f)| \leq \tilde{O} \left( \frac{1}{\sqrt{s}} \right) \right) \geq 1 - \delta,$$

where  $\tilde{O}$  hides the logarithmic factors.

**Remark 1.** Proposition 1 shows that the batch-level quantile  $\hat{\tau}_\alpha^s(f)$  closely approximates the population quantile  $\tau_\alpha(f)$ , with an error of order  $\tilde{O}(1/\sqrt{s})$ , which quantifies the statistical accuracy of the threshold estimation step in CMRM. This relies on a mild regularity assumption on the CDF (smoothness and positive density), standard in quantile estimation theory (e.g., Serfling (2009); Van der Vaart (2000)). We empirically verify these assumptions in Figure 2(c).

Next, we analyze the learning bound. Instead of evaluating the unconditional classification risk, CMRM focuses on the retained high-margin region where the model is confident. This reflects the mechanism of the method: low-margin samples—more likely to be corrupted under label noise—are suppressed during training, while high-margin samples dominate the learning signal. Accordingly, the appropriate clean-distribution objective is the conditional margin risk restricted to this retained region, analogous in spirit to Conditional Value at Risk (CVaR)-style tail-risk analyses in the robustness literature.

To this end, we first define the conformal margin risk on the clean distribution as:

$$\mathcal{L}_{\text{cr}}(f) = -\mathbb{E}_{(X, Y) \sim \mathcal{P}} [M_f(X, Y) \mid M_f(X, Y) \geq \tau_\alpha(f)].$$

We define its empirical Rademacher complexity as:

$$\hat{\mathfrak{R}}_n(\mathcal{F}) := \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i, \tilde{Y}_i) \right],$$

where  $\{\sigma_i\}_{i=1}^n$  are Rademacher variables. The following theorem analyzes the learning bound of CMRM.

**Theorem 1** (Learning bound). *Suppose the assumptions in Proposition 1 hold. Define  $\delta_w$  as the average Wasserstein-1 distance between the noisy and*

*clean label distributions conditional on  $X$ , where  $\delta_w = \mathbb{E}_{X \sim \mathcal{P}(X)} [W_1(\mathcal{P}_{\text{noisy}}(\cdot \mid X), \mathcal{P}(\cdot \mid X))]$ . Then the learning bound of CMRM is:*

$$\mathcal{L}_{\text{cr}}(f) - \hat{\mathcal{L}}_{\text{cr}}^s(f) \leq \tilde{O} \left( \hat{\mathfrak{R}}_n(\mathcal{F}) + \frac{1}{\sqrt{s}} + \delta_w + \alpha + \text{temp} \right).$$

**Remark 2.** Theorem 1 characterizes how the surrogate objective optimized by CMRM on noisy data relates to the clean retained-region margin risk, and clearly decomposes the generalization gap of CMRM into three terms reflecting its core components: quantile estimation error  $\tilde{O}(1/\sqrt{s})$ , function class complexity  $\hat{\mathfrak{R}}_n(\mathcal{F})$ , and distribution shift  $\delta_w$ . These correspond to the use of batchwise quantiles, empirical-to-population approximation on noisy data, and the mismatch between noisy and clean label distributions, respectively. Together, they provide an assumption-light characterization of CMRM’s robustness under arbitrary label noise.

## 5 EXPERIMENTS AND RESULTS

### 5.1 Multi-class Classification Experiments

**Datasets.** We evaluate CMRM under both synthetic and real-world noisy supervision. For *synthetic noise*, we conduct experiments on CIFAR-100 (Krizhevsky et al., 2009), mini-ImageNet (Vinyals et al., 2016), and Food-101 (Bossard et al., 2014). CIFAR-100 is corrupted with class-conditional label noise, where each label is randomly replaced with another label within the same coarse superclass, following the protocol of (Yao et al., 2023). For mini-ImageNet and Food-101, we apply symmetric label flips following (Jiang et al., 2020). We vary the noise rate from  $\{0\%, 5\%, 10\%, 20\%, 30\%, 40\%\}$ , where 0% corresponds to the clean-label setting; additional implementation details are provided in Appendix B.1. For *real-world noise*, we evaluate on CIFAR-10N and its four variants, as well as CIFAR-100N (Wei et al., 2022), whose labels were collected through human annotation. The noise rates of these datasets range from 9.03% to 40.20%, covering a broad spectrum of labeling quality.

**Baselines and ML Models.** For synthetic noise experiments, we compare CMRM against four representative families of methods: **(i)** cross-entropy (CE), the standard objective; **(ii)** Focal loss (Lin et al., 2017), a standard robustness-oriented variant; **(iii)** LDAM (Cao et al., 2019), a margin-based objective; and **(iv)** GCE (Zhang and Sabuncu, 2018), a general-purpose noise-robust loss that remains a common baseline in the recent learning from noisy labels literature (Engleson and Azizpour, 2024; Nguyen et al.,

Detailed statistics and descriptions of CIFAR-10N and CIFAR-100N are summarized at <http://noisylab.com>.

Dataset	Metric	CE		Focal		LDAM		GCE	
		Base	+CMRM	Base	+CMRM	Base	+CMRM	Base	+CMRM
CIFAR-100	ACC (%) $\uparrow$	65.16	<b>66.32</b> (+1.16)	64.42	<b>65.39</b> (+0.97)	59.63	<b>61.12</b> (+1.49)	62.17	<b>63.65</b> (+1.48)
	M.APSS $\downarrow$	6.67	<b>6.52</b> (-2.25%)	6.89	<b>6.61</b> (-4.06%)	17.85	<b>17.67</b> (-0.78%)	7.70	<b>6.28</b> (-18.44%)
mini-ImageNet	ACC (%) $\uparrow$	57.42	<b>59.41</b> (+1.99)	55.54	<b>57.93</b> (+2.39)	56.60	<b>56.62</b> (+0.02)	55.16	<b>55.51</b> (+0.35)
	M.APSS $\downarrow$	7.40	<b>7.04</b> (-4.86%)	7.67	<b>7.28</b> (-5.08%)	13.07	<b>12.88</b> (-1.45%)	9.90	<b>9.59</b> (-3.13%)
Food-101	ACC (%) $\uparrow$	56.21	<b>58.48</b> (+2.27)	56.35	<b>58.92</b> (+2.57)	55.49	<b>56.42</b> (+0.93)	55.66	<b>59.05</b> (+3.39)
	M.APSS $\downarrow$	7.93	<b>6.96</b> (-12.23%)	7.76	<b>6.89</b> (-11.21%)	11.94	<b>11.53</b> (-3.43%)	8.41	<b>6.93</b> (-20.44%)

Table 2: **Top-1 accuracy (%) and marginal average prediction set size (M.APSS  $\downarrow$ ) on multi-class datasets corrupted by synthetic noise with noise rate 20%.** Each Base objective is paired with its +CMRM counterpart; the better value within each pair is in **bold**. Numbers in parentheses indicate the relative change (%): + denotes accuracy improvement, and - denotes M.APSS reduction compared to the corresponding Base objective. On average across all datasets and objectives, CMRM improves accuracy by 1.58 and reduces M.APSS by 7.28%.

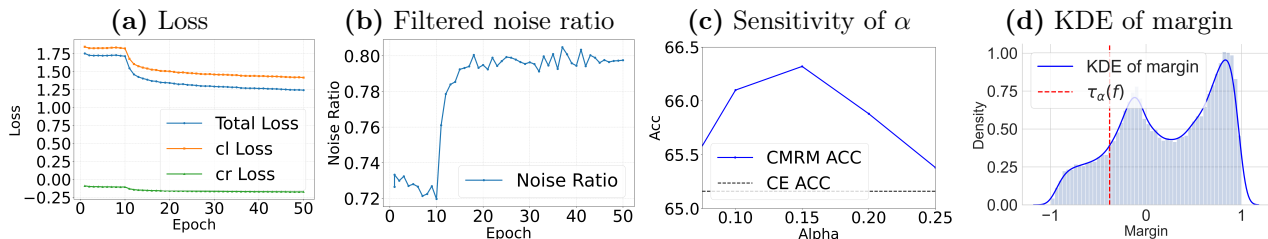


Figure 2: **Justification experiments for multi-class classification** on CIFAR-100 with 20% synthetic label noise. Subfigure (a) shows the training dynamics of total loss (Total), classification loss (cl), and CMRM loss (cr) over epochs. CMRM exhibits stable and monotonic convergence alongside standard loss components. Subfigure (b) reports the ratio of noisy samples among those filtered (with soft weight  $< 0.5$ ) out by CMRM ( $\alpha = 0.15$ ) at each epoch, demonstrating that CMRM consistently suppresses noisy examples by excluding low-margin samples during training. Subfigure (c) examines the sensitivity of  $\alpha$ , showing that CMRM maintains higher accuracy than CE across a range of  $\alpha$  values, indicating robustness to hyperparameter  $\alpha$ . Subfigure (d) depicts the kernel density estimate (KDE) of the margin distribution, with the vertical dashed line indicating the estimated  $\tau_\alpha(f)$  with  $\alpha = 0.15$ . The density curve is smooth and strictly positive around  $\tau_\alpha(f)$ , supporting the differentiability and positive-density assumption in Proposition 1.

Method	CIFAR-10N (worst)		CIFAR-100N	
	ACC(%)	M.APSS	ACC(%)	M.APSS
NI-ERM	95.71	0.93	83.17	1.49
NI-ERM+CMRM	<b>97.19</b>	<b>0.91</b>	<b>83.95</b>	<b>1.29</b>
	(+1.48)	(-2.15%)	(+0.78)	(-13.42%)

Table 3: **Top-1 accuracy (%) and marginal average prediction set size (M.APSS  $\downarrow$ ) on CIFAR-10N and CIFAR-100N corrupted by human annotation noise.** Numbers in parentheses indicate the relative change: + denotes accuracy improvement and - denotes M.APSS reduction. CMRM consistently improves accuracy and reduces uncertainty across CIFAR-N variants, with the largest gains observed on CIFAR-10N Worst and CIFAR-100N. Full results on all CIFAR-10N variants (Aggre, Rand1-3, Worst) and CIFAR-100N are provided in Appendix B.

2024). All methods use a ResNet-20 (He et al., 2016) backbone with standard data augmentation. For real-world noise experiments, we benchmark CMRM against NI-ERM (Zhu et al., 2024), a recent state-of-the-art approach that achieves strong performance on CIFAR-10N and CIFAR-100N by training a linear classifier on frozen DINOv2 (Oquab et al., 2023) features. All methods share the same protocol. For

CMRM, We set `temp` = 1.0 in all experiments and select  $\lambda \in \{0.05, \dots, 0.25\}$  and  $\alpha \in \{0.05, \dots, 0.25\}$  via grid search; See Appendix B.1 for full details.

**Evaluation Metrics.** We evaluate models using marginal Top-1 accuracy and marginal average prediction set size (M.APSS). To assess predictive uncertainty, we adopt CP and define prediction sets and average prediction set size (APSS) (Romano et al., 2020; Angelopoulos et al., 2022). For each input, CP outputs a prediction set containing the true label with high probability, controlled by a target coverage rate of 0.9. For multi-class settings, M.APSS denotes the marginal variant of APSS, i.e., averaged uniformly over all test examples. It also complements top-1 accuracy by capturing how sharply the model distinguishes plausible labels under noisy supervision. We also report class-conditional variants (PC APSS, NC APSS) for binary classification.

**CMRM improves accuracy and reduces uncertainty under different types of noise.** Table 2 reports results on CIFAR-100, mini-ImageNet, and Food-101 corrupted by synthetic label noise. For each objective (CE, Focal, LDAM, GCE), we compare the

Method	Evaluation Metric							
	AUROC ( $\uparrow$ )	AUPRC ( $\uparrow$ )	FNR ( $\downarrow$ )	FPR ( $\downarrow$ )	ACC ( $\uparrow$ )	M.APSS ( $\downarrow$ )	PC APSS ( $\downarrow$ )	NC APSS ( $\downarrow$ )
LR	0.784	0.885	<b>0.073</b>	0.571	0.802	1.223	1.154	1.432
LR + CMRM	<b>0.852</b> (+0.068)	<b>0.925</b> (+0.04)	0.082 (+0.009)	<b>0.422</b> (-0.149)	<b>0.833</b> (0.031)	<b>1.209</b> (-1.15%)	<b>1.109</b> (-3.9%)	<b>1.308</b> (-8.66%)
Focal	0.809	0.890	0.136	0.388	0.801	1.257	1.224	1.356
Focal + CMRM	<b>0.872</b> (+0.063)	<b>0.942</b> (+0.052)	<b>0.128</b> (-0.008)	<b>0.324</b> (-0.064)	<b>0.823</b> (0.022)	<b>1.221</b> -2.87%	<b>1.148</b> -6.21%	<b>1.295</b> -4.5%
SVM	0.808	0.925	<b>0.029</b>	0.807	0.776	1.276	1.370	1.512
SVM + CMRM	<b>0.847</b> (+3.9%)	<b>0.937</b> (+1.2%)	0.048 (+1.9%)	<b>0.585</b> (-22.2%)	<b>0.817</b> (+4.1%)	<b>1.199</b> (-6.03%)	<b>1.322</b> (-3.5%)	<b>1.343</b> (-11.17%)
GCE	0.819	0.904	<b>0.119</b>	0.424	<b>0.804</b>	1.286	<b>1.176</b>	1.396
GCE + CMRM	<b>0.846</b> (+0.027)	<b>0.928</b> (+0.024)	0.172 (+0.053)	<b>0.286</b> (-0.138)	0.800 (-0.004)	<b>1.273</b> (-1.01%)	1.207 (+2.63%)	<b>1.340</b> (-4.01%)

Table 4: **Results for binary classification on Adult** with 20% label noise. We report ranking (AUROC, AUPRC), error rates (FNR, FPR), accuracy (ACC), and uncertainty (M.APSS, PC APSS, NC APSS).  $\uparrow$  indicates higher is better;  $\downarrow$  lower is better. Best results for each base method (LR, Focal, SVM, GCE) are in **bold**. Numbers in parentheses indicate absolute and relative changes (%) of CMRM, where + denotes increase and - denotes decrease. On average, CMRM improves AUROC by 0.049, AUPRC by 0.032, and ACC by 0.023, while reducing FPR by 0.143 and slightly increasing FNR by 0.018. For uncertainty-aware metrics, CMRM reduces M.APSS, PC APSS, and NC APSS by 2.77%, 2.75%, and 7.09% on average, respectively. Results on Email and Credit show similar trends and are provided in Appendix C.2.

Base model with its +CMRM counterpart. On average across all datasets and objectives, CMRM improves accuracy by 1.58 and reduces M.APSS by 7.28%. The improvements are most pronounced for CE, Focal, and GCE, with accuracy gains of up to 3.39 and substantial uncertainty reduction. Even when combined with LDAM, which already encourages margin separation, CMRM consistently yields additional accuracy improvements without increasing uncertainty. Table 3 summarizes results on CIFAR-10N and CIFAR-100N with human-annotated label noise. CMRM consistently outperforms NI-ERM across all CIFAR-N variants, achieving the largest gains on the most challenging settings (CIFAR-10N Worst and CIFAR-100N). Complete results on all CIFAR-10N variants (Aggre, Rand1-3, Worst) are provided in Appendix B. These findings clearly demonstrate that CMRM improves accuracy and reduces uncertainty under both synthetic and real-world noisy supervision.

**CMRM loss convergence results.** Figure 2(a) shows the training dynamics of the classification loss and the CMRM regularization loss. Both components decrease steadily and stabilize as training progresses, indicating smooth joint optimization. The CMRM term integrates with standard objectives and does not introduce instability or slow down of convergence, demonstrating that CMRM can be efficiently optimized.

**CMRM filters out noisy samples during training.** Figure 2(b) shows the fraction of noisy samples among those excluded by CMRM (with soft weight  $< 0.5$ ) at each epoch. This proportion rapidly increases during the early training phase and stabilizes above 78%, indicating that CMRM consistently identifies and filters out mislabeled examples via its margin-based thresholding mechanism.

**CMRM is robust to the choice of hyperparameter  $\alpha$ .** Figure 2(c) examines the sensitivity of CMRM to the hyperparameter  $\alpha$ . Across a range of  $\alpha$  val-

ues, CMRM consistently achieves higher accuracy than CE, indicating that its performance is robust to the choice of  $\alpha$  and does not rely on careful hyperparameter tuning. Notably, CMRM performs well even when  $\alpha$  does not match the true noise rate, suggesting that exact noise rate knowledge is unnecessary. In practice,  $\alpha \in [0.1, 0.2]$  consistently yields strong results across all settings tested.

**Assumptions in Proposition 1 are empirically valid.** Figure 2(d) presents the kernel density estimate (KDE) of the margin distribution, with the vertical dashed line indicating the estimated  $\tau_\alpha(f)$ . The density curve is smooth and strictly positive in the neighborhood of  $\tau_\alpha(f)$ , supporting the differentiability and positive-density assumption in Proposition 1.

**CMRM incurs no penalty on clean labels.** Notably, CMRM also improves or matches accuracy at 0% noise across all objectives (Table 6 in Appendix B.2), confirming that the regularizer incurs no penalty even when labels are clean.

## 5.2 Binary Classification Experiments

**Experiment Setup.** We also evaluate the binary variant of CMRM on three datasets: Email (Hopkins et al., 1999), Credit (Quinlan, 1987), and Adult (Becker and Kohavi, 1996). To simulate label noise, we randomly flip 20% of the training labels while keeping the test labels clean. Additional implementation details are provided in Appendix C.1. Baselines include logistic regression (LR), focal loss, support vector machines (SVM) with hinge loss as a margin-based method, and GCE. LR, Focal, and GCE models use a two-layer MLP, while SVM employs a linear kernel with default regularization. We report AUC-ROC and PR-AUC to assess ranking quality, FPR and FNR to capture class-specific error tendencies, and Accuracy as a general indicator. We also measure the predictive uncertainty using

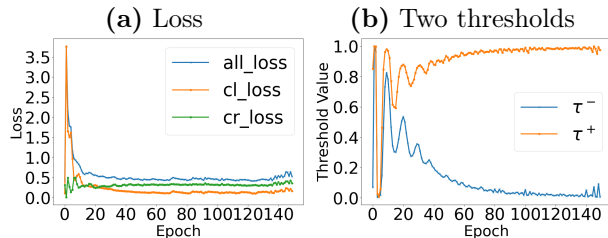


Figure 3: **Training dynamics of LR+CMRM for binary classification** on the Email dataset. Subfigure (a) Training dynamics of total loss (all loss), classification loss (cl loss), and CMRM loss (cr loss) over epochs. CMRM exhibits stable and monotonic convergence alongside standard loss components. Subfigure (b)  $\tau^-$  (negative class threshold) and  $\tau^+$  (positive class threshold) of LR+CMRM during training. The separation between the thresholds increases, indicating that CMRM actively maximizes the margin between positive and negative classes.

APSS, including marginal (M.APSS), positive-class (PC APSS), and negative-class (NC APSS) variants. Larger values of AUC-ROC, PR-AUC, and Accuracy indicate better performance ( $\uparrow$ ), whereas smaller values of FPR, FNR, and APSS metrics (M.APSS, PC APSS, NC APSS) are preferred ( $\downarrow$ ).

**CMRM improves robustness for binary classification.** Table 4 reports results on the Adult dataset with 20% label noise. CMRM consistently improves ranking performance, with AUROC and AUPRC increasing by 0.049 and 0.032 on average, respectively, indicating that it enables models to better separate classes under noisy supervision. These gains are accompanied by improvements in classification metrics, with Accuracy increasing by 0.023 on average and FPR decreasing by 0.143, at the cost of a modest increase in FNR (+0.018). Finally, CMRM reduces predictive uncertainty, as reflected by lower M.APSS, PC APSS, and NC APSS (average reduction of 2.77%, 2.75%, and 7.1%, respectively), indicating sharper and more discriminative predictions. Overall, these improvements demonstrate that CMRM enhances robustness in binary classification under label noise. Similar trends are observed on the Email and Credit datasets (see Appendix C.2 for complete results).

**CMRM optimization dynamics in binary classification.** Figure 3 examines the training behavior of LR + CMRM on the Email dataset with 20% label noise. Subfigure (a) shows that both the classification and CMRM regularization losses decrease steadily and stabilize, indicating that the joint objective can be optimized smoothly on binary classification data. Subfigure (b) tracks the evolution of the class-conditional thresholds  $\tau^+$  and  $\tau^-$  during training.  $\tau^-$  steadily increases while  $\tau^+$  decreases, leading to a widening gap between them. This growing separation reflects that CMRM

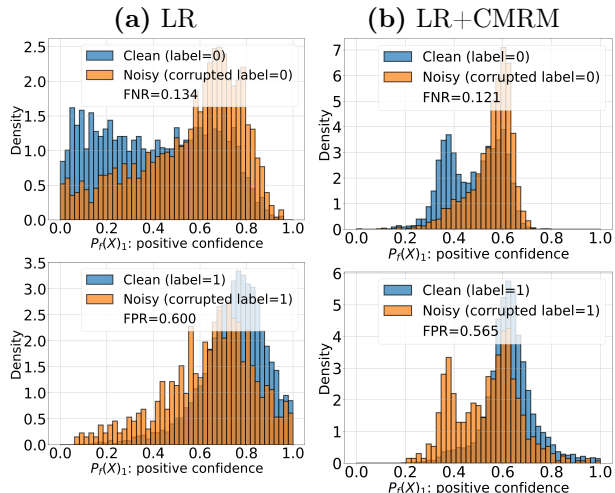


Figure 4: **Histograms of positive confidence distributions for clean (blue) and noisy (orange) samples** on the Credit dataset with 20% label noise. The top and bottom rows correspond to samples with observed labels  $\tilde{Y} = 0$  (negative) and  $\tilde{Y} = 1$  (positive), respectively. Distributions are obtained using LR (left) and LR+CMRM (right). CMRM induces a clearer separation between clean and noisy confidence distributions for both classes.

could actively enlarge the margin between positive and negative classes.

**CMRM separates clean and noisy supervision.** Figure 4 shows the effect of CMRM on positive-class confidence distributions under label noise. With LR (Figure 4a), the distributions of clean and noisy samples overlap substantially for both observed positive and negative label groups, indicating limited ability to distinguish reliable labels from corrupted labels. In contrast, LR+CMRM (Figure 4b) yields a clearer separation: clean samples concentrate in high-confidence regions for their true class, while noisy samples are pushed toward lower-confidence regions, improving their distinguishability in the confidence space.

## 6 CONCLUSION

We propose CMRM, a simple envelope framework for robust learning under label noise, requiring only mild regularity of the margin distribution rather than noise models, clean data or auxiliary resources. CMRM uses a batch-wise conformal quantile on confidence margins to focus training on reliable samples while suppressing likely corrupted ones. Theoretically, we establish a learning bound under arbitrary label noise. Empirically, CMRM integrates smoothly into standard pipelines and consistently improves accuracy and robustness across binary and multi-class benchmarks. As a single regularization term with no architectural cost, CMRM is a natural default when training with noisy labels.

## ACKNOWLEDGEMENTS

The authors gratefully acknowledge the in part support by the USDA-NIFA funded AgAID Institute award 2021-67021-35344, and NSF grants CNS-2312125, IIS-2443828, DUE-2519063. The views expressed are those of the authors and do not reflect the official policy or position of the USDA-NIFA and NSF.

## References

- Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- Anastasios N Angelopoulos, Amit Pal Kohli, Stephen Bates, Michael Jordan, Jitendra Malik, Thayer Alshaabi, Srigokul Upadhyayula, and Yaniv Romano. Image-to-image regression with distribution-free uncertainty quantification and applications in imaging. In *International Conference on Machine Learning*, pages 717–730. PMLR, 2022.
- Eric Arazo, Diego Ortego, Paul Albert, Noel O’Connor, and Kevin McGuinness. Unsupervised label noise modeling and loss correction. In *International conference on machine learning*, pages 312–321. PMLR, 2019.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- Devansh Arpit, Stanislaw Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In *International Conference on Machine Learning*, pages 233–242. PMLR, 2017.
- Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI: <https://doi.org/10.24432/C5XW20>.
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VI 13*, pages 446–461. Springer, 2014.
- Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *Advances in neural information processing systems*, 32, 2019.
- Pedro Celard, Eva Lorenzo Iglesias, José Manuel Sorribes-Fdez, Rubén Romero, A Seara Vieira, and Lourdes Borrajo. A survey on deep learning applied to medical images: from simple artificial neural networks to generative models. *Neural Computing and Applications*, 35(3):2291–2323, 2023.
- Hao Cheng, Zhaowei Zhu, Xingyu Li, Yifei Gong, Xing Sun, and Yang Liu. Learning with instance-dependent label noise: A sample sieve approach. *arXiv preprint arXiv:2010.02347*, 2020.
- Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9268–9277, 2019.
- Richard M Dudley. *Real analysis and probability*. Chapman and Hall/CRC, 2018.
- Bat-Sheva Einbinder, Yaniv Romano, Matteo Sesia, and Yanfei Zhou. Training uncertainty-aware classifiers with conformalized deep learning. *Advances in neural information processing systems*, 35:22380–22395, 2022.
- Gamaleldin Fathy Elsayed, Dilip Krishnan, Hossein Mobahi, Kevin Regan, and Samy Bengio. Large margin deep networks for classification. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.
- Erik Engleson and Hossein Azizpour. Robust classification via regression for learning with noisy labels. In *ICLR 2024—The Twelfth International Conference on Learning Representations, Messe Wien Exhibition and Congress Center, Vienna, Austria, May 7–11t, 2024*, 2024.
- Matteo Fontana, Gianluca Zeni, and Simone Vantini. Conformal prediction: a unified review of theory and new challenges. *Bernoulli*, 29(1):1–23, 2023.
- Benoît Frénay and Michel Verleysen. Classification in the presence of label noise: a survey. *IEEE transactions on neural networks and learning systems*, 25(5):845–869, 2013.
- Subhankar Ghosh, Taha Belkhouja, Yan Yan, and Jannardhan Rao Doppa. Improving Uncertainty Quantification of Deep Classifiers via Neighborhood Conformal Prediction: Novel Algorithm and Theoretical Analysis. In *Proc. of AAAI Conf.*, pages 7722–7730, 2023a.
- Subhankar Ghosh, Yuanjie Shi, Taha Belkhouja, Yan Yan, Jana Doppa, and Brian Jones. Probabilistically Robust Conformal Prediction. In *UAI Conf.*, volume 216 of *Proc. of Machine Learning Research*, pages 681–690. PMLR, 2023b.

- Shivani Gupta and Atul Gupta. Dealing with noise problem in machine learning data-sets: A systematic review. *Procedia Computer Science*, 161:466–474, 2019.
- Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. *Advances in neural information processing systems*, 31, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. Using trusted data to train deep networks on labels corrupted by severe noise. *Advances in neural information processing systems*, 31, 2018.
- Mark Hopkins, Erik Reeber, George Forman, and Jaap Suermondt. Spambase. UCI Machine Learning Repository, 1999. DOI: <https://doi.org/10.24432/C53G6X>.
- Heinrich Jiang, Been Kim, Melody Y. Guan, and Maya R. Gupta. To trust or not to trust a classifier. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 31, 2018.
- Lu Jiang, Di Huang, Mason Liu, and Weilong Yang. Beyond synthetic noise: Deep learning on controlled noisy labels. In *International conference on machine learning*, pages 4804–4815. PMLR, 2020.
- Justin M Johnson and Taghi M Khoshgoftaar. A survey on classifying big data with label noise. *ACM Journal of Data and Information Quality*, 14(4):1–43, 2022.
- Taehyeon Kim, Jongwoo Ko, JinHwan Choi, Se-Young Yun, et al. Fine samples for learning with noisy labels. *Advances in Neural Information Processing Systems*, 34:24137–24149, 2021.
- Shayan Kiyani, George J Pappas, and Hamed Hassani. Length optimization in conformal prediction. *Advances in Neural Information Processing Systems*, 37:99519–99563, 2024.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Toronto, ON, Canada*, 2009.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, 2017.
- Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.
- Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. *arXiv preprint arXiv:2002.07394*, 2020.
- Sijia Li, Yong Fan, Yue Ma, and Ya Pan. Evaluation of dataset distribution and label quality for autonomous driving system. In *2021 IEEE 21st International Conference on Software Quality, Reliability and Security Companion (QRS-C)*, pages 196–200. IEEE, 2021a.
- Xuefeng Li, Tongliang Liu, Bo Han, Gang Niu, and Masashi Sugiyama. Provably end-to-end label-noise learning without anchor points. In *International conference on machine learning*, pages 6403–6413. PMLR, 2021b.
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- Yexiong Lin, Yu Yao, and Tongliang Liu. Learning the latent causal structure for modeling label noise. *Advances in Neural Information Processing Systems*, 37:120549–120577, 2024.
- Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-learning regularization prevents memorization of noisy labels. *Advances in neural information processing systems*, 33:20331–20342, 2020.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning with noisy labels. *Advances in neural information processing systems*, 26, 2013.
- Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language models. *ACM Transactions on Intelligent Systems and Technology*, 16(5):1–72, 2025.
- Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. *Advances in neural information processing systems*, 30, 2017.
- Tri Nguyen, Shahana Ibrahim, and Xiao Fu. Noisy label learning with instance-dependent outliers: Identifiability via crowd wisdom. *Advances in Neural Information Processing Systems*, 37:97261–97298, 2024.

- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1944–1952, 2017.
- Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q Weinberger. Identifying mislabeled data using the area under the margin ranking. *Advances in Neural Information Processing Systems*, 33:17044–17056, 2020.
- J. R. Quinlan. Credit Approval. UCI Machine Learning Repository, 1987. DOI: <https://doi.org/10.24432/C5FS30>.
- Yaniv Romano, Matteo Sesia, and Emmanuel Candes. Classification with valid and adaptive coverage. *Advances in Neural Information Processing Systems*, 33:3581–3591, 2020.
- Robert J Serfling. *Approximation theorems of mathematical statistics*. John Wiley & Sons, 2009.
- Hooman Shahrokhi, Devjeet Raj Roy, Yan Yan, Venera Arnaudova, and Janaradhan Rao Doppa. Conformal Prediction Sets for Deep Generative Models via Reduction to Conformal Regression, 2025. URL <https://arxiv.org/abs/2503.10512>.
- Jialin Shi, Kailai Zhang, Chenyi Guo, Youquan Yang, Yali Xu, and Ji Wu. A survey of label-noise deep learning for medical image analysis. *Medical image analysis*, 95:103166, 2024a.
- Yuanjie Shi, Subhankar Ghosh, Taha Belkhouja, Jana Doppa, and Yan Yan. Conformal Prediction for Class-wise Coverage via Augmented Label Rank Calibration. In *Advances in Neural Information Processing Sys. (NeurIPS)*, 2024b.
- Yuanjie Shi, Hooman Shahrokhi, Xuesong Jia, Xiongzhi Chen, Janardhan Rao Doppa, and Yan Yan. Direct prediction set minimization via bilevel conformal classifier training. *arXiv preprint arXiv:2506.06599*, 2025.
- Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *IEEE transactions on neural networks and learning systems*, 34(11):8135–8153, 2022.
- David Stutz, Ali Taylan Cemgil, Arnaud Doucet, et al. Learning optimal conformal classifiers. *arXiv preprint arXiv:2110.09192*, 2021.
- Timm Teubner, Christoph M Flath, Christof Weinhart, Wil Van Der Aalst, and Oliver Hinz. Welcome to the era of chatgpt et al. the prospects of large language models. *Business & Information Systems Engineering*, 65(2):95–101, 2023.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine*, 29(8):1930–1940, 2023.
- Donna Tjandra and Jenna Wiens. Leveraging an alignment set in tackling instance-dependent label noise. In *Conference on Health, Inference, and Learning*, pages 477–497. PMLR, 2023.
- Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29, 2016.
- Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random world*. Springer Science & Business Media, 2005.
- Haoyu Wang, Zhuo Huang, Zhiwei Lin, and Tongliang Liu. Noisept: Label noise detection and rectification through probability curvature. *Advances in Neural Information Processing Systems*, 37:120159–120183, 2024.
- Jiaheng Wei, Zhaowei Zhu, Hao Cheng, Tongliang Liu, Gang Niu, and Yang Liu. Learning with noisy labels revisited: A study using real-world human annotations. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=TBWA6PLJZQm>.
- Xiaobo Xia, Tongliang Liu, Nannan Wang, Bo Han, Chen Gong, Gang Niu, and Masashi Sugiyama. Are anchor points really indispensable in label-noise learning? *Advances in neural information processing systems*, 32, 2019.
- Jiangchao Yao, Bo Han, Zhihan Zhou, Ya Zhang, and Ivor W Tsang. Latent class-conditional noise model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):9964–9980, 2023.
- Yu Yao, Tongliang Liu, Bo Han, Mingming Gong, Jiankang Deng, Gang Niu, and Masashi Sugiyama. Dual t: Reducing estimation error for transition matrix in label-noise learning. *Advances in neural information processing systems*, 33:7260–7271, 2020.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *Proceedings of the National Academy of Sciences*, 118(3), 2021.

Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 46(8):5625–5644, 2024.

Mingyuan Zhang and Shivani Agarwal. Multiclass learning from noisy labels for non-decomposable performance measures. In *International Conference on Artificial Intelligence and Statistics*, pages 2170–2178. PMLR, 2024.

Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. *Advances in neural information processing systems*, 31, 2018.

Xia Zhao, Limin Wang, Yufei Zhang, Xuming Han, Muhammet Deveci, and Milan Parmar. A review of convolutional neural networks in computer vision. *Artificial Intelligence Review*, 57(4):99, 2024.

Zhi-Hua Zhou and Xu-Ying Liu. Training cost-sensitive neural networks with methods addressing the class imbalance problem. *IEEE Transactions on knowledge and data engineering*, 18(1):63–77, 2006.

Yilun Zhu, Jianxin Zhang, Aditya Gangrade, and Clay Scott. Label noise: Ignorance is bliss. *Advances in Neural Information Processing Systems*, 37:116575–116616, 2024.

## CHECKLIST

1. For all models and algorithms presented, check if you include:
  - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes] We provide it in Section 4.2.
  - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Yes] We provide it in Section 4.2.
  - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes/No/Not Applicable]
2. For any theoretical claim, check if you include:
  - (a) Statements of the full set of assumptions of all theoretical results. [Yes] We provide it in Section 4.2.
  - (b) Complete proofs of all theoretical results. [Yes] We provide it in Appendix A.
  - (c) Clear explanations of any assumptions. [Yes] We provide it in Section 4.2.
3. For all figures and tables that present empirical results, check if you include:
  - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes] we provide it in the supplemental material.
  - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes] we provide it in Appendix B and C.
  - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes] we provide it in Appendix B and C.
  - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Not Applicable]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
  - (a) Citations of the creator If your work uses existing assets. [Yes] we cite it in 5.
  - (b) The license information of the assets, if applicable. [Not Applicable]
  - (c) New assets either in the supplemental material or as a URL, if applicable. [Yes] we provide it in the supplemental material.
  - (d) Information about consent from data providers/curators. [Not Applicable]
  - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Not Applicable]
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
  - (a) The full text of instructions given to participants and screenshots. [Yes/No/Not Applicable]
  - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Yes/No/Not Applicable]
  - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Yes/No/Not Applicable]

## A TECHNICAL PROOFS

### A.1 Technical Proofs for Proposition 1

**Proposition (Re) 1** (Gap between  $\tau_\alpha$  and  $\hat{\tau}_\alpha^s$ ). *Denote by  $G(t)$  the cumulative distribution function (CDF) of  $M_f(X, \tilde{Y})$  under the noisy distribution  $\mathcal{P}_{\text{noisy}}$ . Assume that  $G(t)$  is continuously differentiable in a neighborhood of  $\tau_\alpha(f)$  with density  $g(t)$  and  $g(\tau_\alpha(f)) > 0$ . Then, for any  $\delta \in (0, 1)$ , we have:*

$$\mathbb{P}\left(|\tau_\alpha(f) - \hat{\tau}_\alpha^s(f)| \leq \tilde{O}\left(\frac{1}{\sqrt{s}}\right)\right) \geq 1 - \delta,$$

where  $\tilde{O}$  hides the logarithmic factors.

*Proof.* of Proposition 1.

Before proving Proposition 1, we first define  $\hat{G}^s(t) = \frac{1}{s} \sum \mathbb{1}[M_f(X_i, \tilde{Y}_i) \leq t]$  as the empirical CDF of  $M_f(X, \tilde{Y})$  with  $s$  samples.

By the definition,  $\tau_\alpha(f) = G^{-1}(\alpha)$ , where  $\hat{\tau}_\alpha^s(f) := \inf\{t : \hat{G}^s(t) \geq \alpha\}$ .

**Step 1.** By the definition of  $\hat{\tau}_\alpha^s$ , we have  $\hat{G}^s(\hat{\tau}_\alpha^s) \geq \alpha$ . Moreover, since  $\hat{G}^s$  is a step function with jumps of size  $1/s$ , the value at the first crossing satisfies

$$\alpha \leq \hat{G}^s(\hat{\tau}_\alpha^s) \leq \alpha + \frac{1}{s}. \quad (7)$$

Equivalently,  $\hat{\tau}_\alpha^s$  is an order statistic and the ECDF at an order statistic lies in an interval of width  $1/s$ .

**Step 2.** By the DKW inequality, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\sup_{t \in \mathbb{R}} |\hat{G}^s(t) - G(t)| \leq \varepsilon_s(\delta) := \sqrt{\frac{\log(2/\delta)}{2s}}. \quad (8)$$

On the event (8), evaluating at the *same* point  $t = \hat{\tau}_\alpha^s$  gives

$$|G(\hat{\tau}_\alpha^s) - \hat{G}^s(\hat{\tau}_\alpha^s)| \leq \varepsilon_s(\delta).$$

Combining with (7), we obtain the sandwich

$$\alpha - \varepsilon_s(\delta) \leq G(\hat{\tau}_\alpha^s) \leq \alpha + \frac{1}{s} + \varepsilon_s(\delta). \quad (9)$$

**Step 3.** Let  $U = (\tau_\alpha - \eta, \tau_\alpha + \eta)$  and assume  $\inf_{t \in U} g(t) \geq g_0 > 0$ . Then  $G$  is strictly increasing on  $U$ , hence invertible on  $G(U)$ , and its inverse is  $(1/g_0)$ -Lipschitz on  $G(U)$ : for any  $u_1, u_2 \in G(U)$ ,

$$|G^{-1}(u_1) - G^{-1}(u_2)| \leq \frac{1}{g_0} |u_1 - u_2|. \quad (10)$$

Now we show that on the event (8) and for sufficiently large  $s$  (or more explicitly, whenever  $\varepsilon_s(\delta) + 1/s \leq g_0\eta$ ), the probability interval in (9) lies inside  $G(U)$ , which implies  $\hat{\tau}_\alpha^s \in U$  and legitimizes applying (10). Indeed, since  $G(\tau_\alpha) = \alpha$  and  $g \geq g_0$  on  $U$ , we have

$$G(\tau_\alpha + \eta) - \alpha \geq g_0\eta, \quad \alpha - G(\tau_\alpha - \eta) \geq g_0\eta.$$

Thus if  $\varepsilon_s(\delta) + 1/s \leq g_0\eta$ , then

$$\alpha - \varepsilon_s(\delta) \geq G(\tau_\alpha - \eta), \quad \alpha + \frac{1}{s} + \varepsilon_s(\delta) \leq G(\tau_\alpha + \eta),$$

so (9) implies  $G(\hat{\tau}_\alpha^s) \in G(U)$  and hence  $\hat{\tau}_\alpha^s \in U$ .

**Step 4.** On the event (8) and under  $\varepsilon_s(\delta) + 1/s \leq g_0\eta$ , we can apply (10) with  $u_1 = G(\widehat{\tau}_\alpha^s)$  and  $u_2 = \alpha = G(\tau_\alpha)$ :

$$|\widehat{\tau}_\alpha^s - \tau_\alpha| = |G^{-1}(G(\widehat{\tau}_\alpha^s)) - G^{-1}(\alpha)| \leq \frac{1}{g_0} |G(\widehat{\tau}_\alpha^s) - \alpha|.$$

Using (9), we obtain

$$|G(\widehat{\tau}_\alpha^s) - \alpha| \leq \varepsilon_s(\delta) + \frac{1}{s},$$

hence

$$|\widehat{\tau}_\alpha^s - \tau_\alpha| \leq \frac{1}{g_0} \left( \varepsilon_s(\delta) + \frac{1}{s} \right) = \frac{1}{g_0} \left( \sqrt{\frac{\log(2/\delta)}{2s}} + \frac{1}{s} \right).$$

Finally, since (8) holds with probability at least  $1 - \delta$ , the desired high-probability bound follows.  $\square$

## A.2 Technical Proofs for Theorem 1

**Theorem (Re) 1** (Learning bound). *(Theorem 1 restated.) Suppose the assumptions in Proposition 1 hold. Define  $\delta_w$  as the average Wasserstein-1 distance between the noisy and clean label distributions conditional on  $X$ , where  $\delta_w = \mathbb{E}_{X \sim \mathcal{P}(X)} \left[ W_1(\mathcal{P}_{\text{noisy}}(\cdot | X), \mathcal{P}(\cdot | X)) \right]$ . Then the learning bound of CMRM is:*

$$\mathcal{L}_{\text{cr}}(f) - \widehat{\mathcal{L}}_{\text{cr}}^s(f) \leq \tilde{O} \left( \widehat{\mathfrak{R}}_n(\mathcal{F}) + \frac{1}{\sqrt{s}} + \delta_w + \alpha + \text{temp} \right).$$

*Proof.* of Theorem 1.

Before proving Theorem 1, we first define conformal margin risk on the noisy distribution as:

$$\mathcal{L}_{\text{cr}}^{\text{noisy}}(f) = -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} \left[ M_f(X, \tilde{Y}) \cdot \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)] \right].$$

and its surrogate risk on the noisy distribution, where the indicator function  $\mathbb{1}[x \geq y]$  is replaced by the Sigmoid function  $\tilde{\mathbb{1}}[x \geq y] = 1/(1 + \exp(-(x - y)/\text{temp}))$ , as:

$$\tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f) = -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} \left[ M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)] \right].$$

Recall that:

$$\widehat{\mathcal{L}}_{\text{cr}}(f) = \frac{1}{n} \sum_{i=1}^n -M_f(X_i, \tilde{Y}_i) \cdot \tilde{\mathbb{1}}[M_f(X_i, \tilde{Y}_i) \geq \widehat{\tau}_\alpha(f)],$$

$$\widehat{\mathcal{L}}_{\text{cr}}^s(f) = \mathbb{E}_{\widehat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} \left[ \sum_{i=1}^n -M_f(X_i, \tilde{Y}_i) \cdot \tilde{\mathbb{1}}[M_f(X_i, \tilde{Y}_i) \geq \widehat{\tau}_\alpha^s(f)] \right],$$

and

$$\mathcal{L}_{\text{cr}}(f) = -\mathbb{E}_{(X, Y) \sim \mathcal{P}} [M_f(X, Y) | M_f(X, Y) \geq \tau_\alpha(f)].$$

Then we show the following technical lemmas:

**Lemma 1.** *(Immediate results from Theorem 3.3 in (Mohri et al., 2018)) For any  $f \in \mathcal{F}$ , the following inequality holds with high probability:*

$$\mathbb{E}_{\widehat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} \left[ -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} \left[ M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \widehat{\tau}_\alpha^s(f)] \right] \right] \leq \tilde{\mathcal{L}}_{\text{cr}}^s(f) + \tilde{O} \left( \widehat{\mathfrak{R}}_n(\mathcal{F}) + \frac{1}{\sqrt{n}} \right).$$

**Lemma 2.** *Suppose the assumptions in Proposition 1 hold. Then, the following inequality holds:*

$$\left| \tilde{\mathcal{L}}_{cr}^{noisy}(f) - \mathbb{E}_{\hat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} \left[ -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{noisy}} \left[ M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{I}}[M_f(X, \tilde{Y}) \geq \hat{\tau}_\alpha^s(f)] \right] \right] \right| \leq \tilde{O}(1/\sqrt{s}).$$

**Lemma 3.** *Suppose the assumptions in Proposition 1 hold. Then, the following inequality holds:*

$$|\mathcal{L}_{cr}^{noisy}(f) - \tilde{\mathcal{L}}_{cr}^{noisy}(f)| \leq O(\mathbf{temp}).$$

**Lemma 4.** *Define  $\delta_w$  as the average Wasserstein-1 distance between posteriors under the training and test domains, where  $\delta_w = \mathbb{E}_{X \sim \mathcal{P}(X)} \left[ W_1(\mathbb{P}_{\mathcal{P}_{noisy}}(\cdot | X), \mathbb{P}_{\mathcal{P}}(\cdot | X)) \right]$ . Then the following inequality holds:*

$$|\mathcal{L}_{cr}(f) - \mathcal{L}_{cr}^{noisy}(f)| \leq 2\delta_w + \alpha.$$

Now we begin to prove Theorem 1:

$$\begin{aligned} & \mathcal{L}_{cr}(f) - \hat{\mathcal{L}}_{cr}^s(f) \\ &= \underbrace{\mathcal{L}_{cr}(f) - \mathcal{L}_{cr}^{noisy}(f)}_{2\delta_w + \alpha, \text{ Lemma 4}} + \underbrace{\mathcal{L}_{cr}^{noisy}(f) - \tilde{\mathcal{L}}_{cr}^{noisy}(f)}_{O(\mathbf{temp}), \text{ Lemma 3}} \\ & \quad + \underbrace{\tilde{\mathcal{L}}_{cr}^{noisy}(f) - \mathbb{E}_{\hat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} \left[ -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{noisy}} \left[ M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{I}}[M_f(X, \tilde{Y}) \geq \hat{\tau}_\alpha^s(f)] \right] \right]}_{\tilde{O}(1/\sqrt{s}), \text{ Lemma 2}} \\ & \quad + \underbrace{\mathbb{E}_{\hat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} \left[ -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{noisy}} \left[ M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{I}}[M_f(X, \tilde{Y}) \geq \hat{\tau}_\alpha^s(f)] \right] \right] - \hat{\mathcal{L}}_{cr}^s(f)}_{\tilde{O}\left(\hat{\mathfrak{R}}_n(\mathcal{F}) + \frac{1}{\sqrt{n}}\right), \text{ Lemma 1}} \\ & \leq 2\delta_w + \alpha + O(\mathbf{temp}) + \tilde{O}(1/\sqrt{s}) + \tilde{O}\left(\hat{\mathfrak{R}}_n(\mathcal{F}) + \frac{1}{\sqrt{n}}\right) \\ & \leq \tilde{O}\left(\hat{\mathfrak{R}}_n(\mathcal{F}) + \frac{1}{\sqrt{s}} + \delta_w + \alpha + \mathbf{temp}\right), \end{aligned}$$

where the first inequality is due to Lemma 1, 2, 3, 4, and the second inequality is due to  $s \leq n$ .

Therefore, we have:

$$\mathcal{L}_{cr}(f) - \hat{\mathcal{L}}_{cr}^s(f) \leq \tilde{O}\left(\hat{\mathfrak{R}}_n(\mathcal{F}) + \frac{1}{\sqrt{s}} + \delta_w + \alpha + \mathbf{temp}\right).$$

□

## A.3 Proofs for Technical Lemmas

### A.3.1 Proofs for Lemma 2

**Lemma (Re) 1.** *(Lemma 2 restated.) Suppose the assumptions in Proposition 1 hold. Then, the following inequality holds:*

$$\left| \tilde{\mathcal{L}}_{cr}^{noisy}(f) - \mathbb{E}_{\hat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} \left[ -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{noisy}} \left[ M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{I}}[M_f(X, \tilde{Y}) \geq \hat{\tau}_\alpha^s(f)] \right] \right] \right| \leq \tilde{O}(1/\sqrt{s}).$$

*Proof.* of Lemma 2.

Before proving Lemma 2, we first show the following technical lemma:

**Lemma 5.** (Lipschitz continuity of  $\tilde{\mathcal{L}}_{cr}^{noisy}$  in  $\tau$ ) Define the smoothed conformal margin risk on the noisy distribution as a function of the threshold  $\tau$ :

$$\tilde{\mathcal{L}}_{cr}^{noisy}(f; \tau) := -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{noisy}} [M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau]],$$

where the softened indicator is the Sigmoid function  $\tilde{\mathbb{1}}(M_f \geq \tau) := \sigma\left(\frac{M_f(X, \tilde{Y}) - \tau}{\mathbf{temp}}\right)$ , and  $\sigma(u) = \frac{1}{1 + \exp(-u)}$ . Assume  $|M_f(X, \tilde{Y})| \leq 1$  almost surely and  $\mathbf{temp} > 0$ . Then, for any fixed  $f$ , the function  $\tau \rightarrow \tilde{\mathcal{L}}_{cr}^{noisy}(f; \tau)$  is  $L$ -Lipschitz continuous with  $L = 1/(4\mathbf{temp})$ .

The proof of Lemma 5 is deferred to the end of this proof.

Now we begin to prove Lemma 2.

$$\begin{aligned} & \left| \tilde{\mathcal{L}}_{cr}^{noisy}(f) - \mathbb{E}_{\hat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} \left[ -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{noisy}} [M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \hat{\tau}_\alpha^s(f)]] \right] \right| \\ &= \left| -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{noisy}} [M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] - \mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{noisy}} [M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \hat{\tau}_\alpha^s(f)]] \right| \\ &\leq \frac{1}{4\mathbf{temp}} \cdot \left| \tau_\alpha(f) - \mathbb{E}_{\hat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} [\hat{\tau}_\alpha^s(f)] \right| \\ &\leq \frac{1}{4\mathbf{temp}} \cdot \tilde{O}(1/\sqrt{s}) \\ &\leq \tilde{O}(1/\sqrt{s}), \end{aligned}$$

where the first inequality is due to Lemma 5, the second inequality is due to Proposition 1, and the last inequality is due to the setting  $\mathbf{temp} = 1.0$ .

Therefore, we have:

$$\left| \tilde{\mathcal{L}}_{cr}^{noisy}(f) - \mathbb{E}_{\hat{\tau}_\alpha^s(f) \sim \mathcal{T}_\alpha(f)} \left[ -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{noisy}} [M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \hat{\tau}_\alpha^s(f)]] \right] \right| \leq \tilde{O}(1/\sqrt{s}).$$

□

Now we begin to prove Lemma 5.

*Proof.* (of Lemma 5)

Let  $Z := M_f(X, \tilde{Y})$ ,  $u := \frac{Z - \tau}{\mathbf{temp}}$ .

We can rewrite the interested function as  $\tilde{\mathcal{L}}_{cr}^{noisy}(f; \tau) = -\mathbb{E}[Z \cdot \sigma(u)]$ , where the expectation is taken over  $(X, \tilde{Y}) \sim \mathcal{P}_{noisy}$ .

First, we develop the differentiation of  $\tilde{\mathcal{L}}_{cr}^{noisy}$  w.r.t.  $\tau$  by applying the chain rule as follows

$$\frac{\partial \tilde{\mathcal{L}}_{cr}^{noisy}(f; \tau)}{\partial \tau} = -\mathbb{E}[Z \cdot \frac{\partial \sigma(u)}{\partial \tau}]. \quad (11)$$

Since  $u = (Z - \tau)/\mathbf{temp}$ , we have  $\frac{\partial u}{\partial \tau} = -1/\mathbf{temp}$ .

The derivative of the Sigmoid function is

$$\frac{\partial \sigma(u)}{\partial u} = \frac{\exp(-u)}{(1 + \exp(-u))^2}.$$

Thus, by using the chain rule, we have

$$\frac{\partial \sigma(u)}{\partial \tau} = \frac{\partial \sigma(u)}{\partial u} \cdot \frac{\partial u}{\partial \tau} = -\frac{1}{\mathbf{temp}} \cdot \frac{\exp(-u)}{(1 + \exp(-u))^2}.$$

Substituting the above equality back into Equation (11), we have

$$\frac{\partial \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}}{\partial \tau}(f; \tau) = -\mathbb{E}\left[Z \cdot \left(-\frac{1}{\text{temp}} \cdot \frac{\exp(-u)}{(1 + \exp(-u))^2}\right)\right] = \frac{1}{\text{temp}} \cdot \mathbb{E}\left[Z \frac{\exp(-u)}{(1 + \exp(-u))^2}\right]. \quad (12)$$

Then, we would like to bound the derivative. Taking absolute value of (12) and using  $|Z| \leq 1$ , we have

$$\left| \frac{\partial \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f; \tau)}{\partial \tau} \right| \leq \frac{1}{\text{temp}} \cdot \mathbb{E}\left[\frac{\exp(-u)}{(1 + \exp(-u))^2}\right].$$

Recall that  $h(u) := \frac{\exp(-u)}{(1 + \exp(-u))^2} = \sigma(u) \cdot (1 - \sigma(u))$ , which is the standard logistic hat function. It satisfies

$$0 \leq h(u) \leq 1/4, \text{ with } h(0) = 1/4.$$

Therefore, we have  $\mathbb{E}[h(u)] \leq 1/4$ . Hence, we conclude  $\left| \frac{\partial \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f; \tau)}{\partial \tau} \right| \leq 1/(4\text{temp})$ .

Finally, to determine the Lipschitz parameter, for  $\tau_1, \tau_2$ , we have

$$\left| \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f; \tau_1) - \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f; \tau_2) \right| \leq \sup_{\tau} \left| \frac{\partial \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f; \tau)}{\partial \tau} \right| \cdot |\tau_1 - \tau_2| \leq \frac{1}{4\text{temp}} |\tau_1 - \tau_2|.$$

Thus, it shows  $\tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f; \tau)$  is Lipschitz in  $\tau$  with constant  $L = 1/(4\text{temp})$ .  $\square$

### A.3.2 Proof for Lemma 3

**Lemma (Re) 2.** (Lemma 3 restated.) *Suppose the assumptions in Proposition 1 hold. Then, for any  $f$ , there exists a constant  $C > 0$  independent from  $f$  and  $\text{temp}$  such that:*

$$|\mathcal{L}_{\text{cr}}^{\text{noisy}}(f) - \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f)| \leq C \text{temp}.$$

In particular, with Big-O notation, we have

$$|\mathcal{L}_{\text{cr}}^{\text{noisy}}(f) - \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f)| \leq O(\text{temp}).$$

*Proof.* (of Lemma 3.)

Before proving Lemma 3, we first present the following technical lemma:

**Lemma 6.** (Upper bounding the gap between hard and soft indicator functions) *Let*

$$\Gamma(\Delta) := \mathbb{1}[\Delta \geq 0] - \tilde{\mathbb{1}}[\Delta \geq 0], \text{ and } \tilde{\mathbb{1}}[\Delta \geq 0] = 1/(1 + \exp(-\Delta/\text{temp})),$$

where  $\text{temp} > 0$  is the temperature parameter. Then, for every  $\Delta \in \mathbb{R}$ , the following inequality holds:

$$|\Gamma(\Delta)| \leq \exp(-|\Delta|/\text{temp}).$$

Now we begin to prove Lemma 3.

Recall the definitions of the noisy conformal margin risk and its smoothed variant:

$$\begin{aligned} \mathcal{L}_{\text{cr}}^{\text{noisy}}(f) &= -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \cdot \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_{\alpha}(f)]], \\ \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f) &= -\mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau_{\alpha}(f)]]. \end{aligned}$$

where the soft smoothed indicator function is  $\tilde{\mathbb{1}}[M \geq \tau] = \frac{1}{1 + \exp(-(M - \tau)/\text{temp})}$ .

For a fixed  $f$ , define the margin-threshold difference as

$$\Delta_f(X, \tilde{Y}) := M_f(X, \tilde{Y}) - \tau_{\alpha}(f),$$

and define the difference between the hard and soft smoothed indicators as

$$\Gamma(\Delta_f(X, \tilde{Y})) := \mathbb{1}[\Delta_f(X, \tilde{Y}) \geq 0] - \tilde{\mathbb{1}}[\Delta_f(X, \tilde{Y}) \geq 0].$$

By assumption,  $M_f(X, \tilde{Y}) \in (-1, 1)$ . Since  $\tau_\alpha(f)$  is the  $\alpha$ -quantile of  $M_f(X, \tilde{Y})$  under  $\mathcal{P}_{\text{noisy}}$ , we also have  $\tau_\alpha(f) \in (-1, 1)$ . Hence we have

$$\Delta_f(X, \tilde{Y}) \in (-2, 2) \text{ almost surely.}$$

In what follows, we use three steps to prove the desired result.

**Step 1: Reducing to bounding  $\mathbb{E}[|\Gamma(\Delta_f)|]$ .** We first write the difference between the two risks explicitly as follows:

$$\begin{aligned} \mathcal{L}_{\text{cr}}^{\text{noisy}}(f) - \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f) &= -\mathbb{E}_{\text{noisy}}[M_f(X, \tilde{Y}) \cdot \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] + \mathbb{E}_{\text{noisy}}[M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] \\ &= \mathbb{E}_{\text{noisy}}[M_f(X, \tilde{Y}) \cdot (\tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)] - \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)])] \\ &= -\mathbb{E}_{\text{noisy}}[M_f(X, \tilde{Y}) \cdot \Gamma(\Delta_f(X, \tilde{Y}))]. \end{aligned}$$

Taking the absolute values and using  $|M_f(X, \tilde{Y})| \leq B$ , we have

$$\begin{aligned} |\mathcal{L}_{\text{cr}}^{\text{noisy}}(f) - \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f)| &= |\mathbb{E}_{\text{noisy}}[M_f(X, \tilde{Y}) \cdot \Gamma(\Delta_f(X, \tilde{Y}))]| \leq \mathbb{E}_{\text{noisy}}[|M_f(X, \tilde{Y})| \cdot |\Gamma(\Delta_f(X, \tilde{Y}))|] \\ &\leq \mathbb{E}_{\text{noisy}}[B|\Gamma(\Delta_f(X, \tilde{Y}))|]. \end{aligned}$$

Therefore, it suffices to upper bound  $\mathbb{E}_{\text{noisy}}[|\Gamma(\Delta_f)|]$ .

**Step 2: Use the density of  $\Delta_f$  and Lemma 6.** Let  $G(t)$  and  $g(t) = G'(t)$  denote the CDF and the density of  $M_f(X, \tilde{Y})$  under  $\mathcal{P}_{\text{noisy}}$ , respectively. By the assumptions in Proposition 1,  $g$  is continuous and bounded, i.e.,

$$g_\infty := \sup_{t \in \mathbb{R}} g(t) < \infty$$

For fixed  $f$ ,  $\tau_\alpha(f)$  is a constant, so the random variable  $\Delta_f = M_f - \tau_\alpha(f)$  has density which is also continuous and bounded with  $\|g_{\Delta_f}\|_\infty \leq g_\infty$

Therefore, we can write

$$\mathbb{E}_{\text{noisy}}[|\Gamma(\Delta_f)|] = \int_{-2}^2 |\Gamma(\delta)| \cdot g_{\Delta_f}(\delta) d\delta,$$

where we use the fact that  $\Delta_f \in (-2, 2)$  almost surely.

By Lemma 6, for all  $\delta \in \mathbb{R}$ , we have

$$|\Gamma(\delta)| \leq \exp(-|\delta|/\text{temp}).$$

Combining this with the bound on  $g_{\Delta_f}$ , we derive

$$\mathbb{E}_{\text{noisy}}[|\Gamma(\Delta_f)|] = \int_{-2}^2 |\Gamma(\delta)| g_{\Delta_f}(\delta) d\delta \leq g_\infty \cdot \int_{-2}^2 \exp(-|\delta|/\text{temp}) d\delta.$$

The integral can be computed explicitly:

$$\int_{-2}^2 \exp\left(-\frac{|\delta|}{\text{temp}}\right) d\delta = 2 \int_0^2 \exp(-\delta/\text{temp}) d\delta = 2\text{temp}(1 - \exp(-2/\text{temp})) \leq 2\text{temp},$$

where we used  $1 - \exp(-2/\text{temp}) \leq 1$ .

Thus, we have  $\mathbb{E}_{\text{noisy}}[|\Gamma(\Delta_f)|] \leq 2g_\infty \text{temp}$ .

**Step 3: Conclude the upper bound.** Putting everything together, we derive

$$\|\cdot\| \leq \mathbb{E}_{\text{noisy}}[|\Gamma(\Delta_f)|] \leq 2g_\infty \mathbf{temp}.$$

Hence the difference between the hard formal margin risk and its smoothed counter is  $O(\mathbf{temp})$ .

Recall that the  $M_f(X, \tilde{Y})$  is the confidence gap between the confidence score of the observed label and the highest confidence among other candidate labels. Thus,  $M_f(X, \tilde{Y}) \in (-1, 1)$ . As  $\tau_\alpha(f)$  is the  $\alpha$ -quantile of  $M_f(X, \tilde{Y})$ ,  $\tau_\alpha(f) \in (-1, 1)$  and  $\Delta_f(X, Y) \in (-2, 2)$ .

We also recall that  $G(t)$  is the CDF of  $M_f(X, \tilde{Y})$  under  $\mathcal{P}_{\text{noisy}}$ , and  $g(t) = G'(t)$  is its density, which is assumed to be continuous and bounded. Denote  $g_\infty = \sup_t g(t) < +\infty$ . Since  $\tau_\alpha(f)$  is a constant for fixed  $f$ , the density of  $\Delta_f$  is  $g_{\Delta_f}(t) = g(t + \tau_\alpha(f))$ , which satisfies  $\|g_{\Delta_f}\|_\infty = g_\infty$ .

Then, we have:

$$\begin{aligned} & |\mathcal{L}_{\text{cr}}^{\text{noisy}}(f) - \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f)| \\ &= \left| \mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \cdot \tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] - \mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \cdot \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] \right| \\ &= \left| \mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \cdot (\tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)] - \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)])] \right| \\ &\leq \mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} [|M_f(X, \tilde{Y})|] \cdot \left| \mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} [\tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)] - \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] \right| \\ &\leq \mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} [|\tilde{\mathbb{1}}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)] - \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]|] \\ &= \mathbb{E}_{(X, \tilde{Y}) \sim \mathcal{P}_{\text{noisy}}} |\tilde{\mathbb{1}}[\Delta_f(X, Y) \geq 0] - \mathbb{1}[\Delta_f(X, Y) \geq 0]| \\ &= \mathbb{E}[\Gamma(\Delta_f(X, Y))] \\ &= \int_{-2}^2 \Gamma(\Delta_f(X, Y)) \partial_{\Delta_f(X, Y)} \Gamma(\Delta_f(X, Y)) d\Delta_f(X, Y) \\ &\leq \int_{-2}^2 \exp\left(-\frac{|\Delta_f(X, Y)|}{\mathbf{temp}}\right) \partial_{\Delta_f(X, Y)} \Gamma(\Delta_f(X, Y)) d\Delta_f(X, Y) \\ &\leq g_\infty \cdot \int_{-2}^2 \exp\left(-\frac{|\Delta_f(X, Y)|}{\mathbf{temp}}\right) d\Delta_f(X, Y) \\ &= g_\infty \cdot 2\mathbf{temp}(1 - \exp(-2/\mathbf{temp})) \\ &\leq O(\mathbf{temp}), \end{aligned}$$

where the first inequality is due to the Hölder's inequality, the second inequality is due to the Jensen's inequality and expectation is convex function, the third inequality is due to  $M_f(X, \tilde{Y}) \in (-1, 1)$ , the fourth inequality is due to Lemma 6, the fifth inequality is due to the upper bound for the gradient of  $\Delta_f(X, Y)$  is  $g_\infty$ , and the last inequality is due to  $1 - \exp(-2/\mathbf{temp}) \leq 1$ .

Therefore, we have:

$$|\mathcal{L}_{\text{cr}}^{\text{noisy}}(f) - \tilde{\mathcal{L}}_{\text{cr}}^{\text{noisy}}(f)| \leq O(\mathbf{temp}).$$

□

Now we begin to prove Lemma 6.

*Proof.* (of Lemma 6)

We analyze the two cases, i.e.,  $\Delta \geq 0$  and  $\Delta < 0$  separately below.

**Case 1:**  $\Delta \geq 0$

In this case, by using the condition  $\Delta \geq 0$ , we first reformulate  $\Gamma(\Delta)$  as follows

$$\Gamma(\Delta) = 1 - \frac{1}{1 + \exp(-\Delta/\text{temp})} = \frac{\exp(-\Delta/\text{temp})}{1 + \exp(-\Delta/\text{temp})}.$$

Since the denominator is at least 1 (due to  $\exp(\cdot) > 0$ ), we have the following inequalities

$$0 < \Gamma(\Delta) = |\Gamma(\Delta)| \leq \exp(-\Delta/\text{temp}) = \exp(-|\Delta|/\text{temp}), \quad (13)$$

where the last equality is due to  $\Delta \geq 0$ , the condition for this **case 1**. Thus, the desired inequality holds for all  $\Delta > 0$ .

**Case 2:**  $\Delta < 0$

In this case, by using the condition  $\Delta < 0$ , we reformulate  $\Gamma(\Delta)$  as follows

$$\Gamma(\Delta) = 0 - \frac{1}{1 + \exp(-\Delta/\text{temp})} = -\frac{1}{1 + \exp(-\Delta/\text{temp})} < 0.$$

Then, with the condition  $\Delta < 0$ , we have  $|\Delta| = -\Delta$ . To upper bound  $|\Gamma(\Delta)|$  in this case, it suffices to show:

$$|\Gamma(\Delta)| = -\Gamma(\Delta) = \frac{1}{1 + \exp(-\Delta/\text{temp})} \leq \exp(\Delta/\text{temp}) = \exp(-|\Delta|/\text{temp}), \quad (14)$$

where the last inequality is due to  $\exp(a) \leq 1 + \exp(a)$  for any  $a \in \mathbb{R}$ .

Finally, combining (13) and (14) implies that for any  $\Delta \in \mathbb{R}$ , the following inequality holds:

$$|\Gamma(\Delta)| \leq \exp(-|\Delta|/\text{temp}).$$

□

### A.3.3 Proof for Lemma 4

**Lemma (Re) 3.** (*Lemma 4 restated.*) Define  $\delta_w$  as the average Wasserstein-1 distance between posteriors under the training and test domains, where  $\delta_w = \mathbb{E}_{X \sim \mathcal{P}(X)} \left[ W_1(\mathbb{P}_{\mathcal{P}_{\text{noisy}}}(\cdot | X), \mathbb{P}_{\mathcal{P}}(\cdot | X)) \right]$ . Then the following inequality holds:

$$|\mathcal{L}_{cr}(f) - \mathcal{L}_{cr}^{\text{noisy}}(f)| \leq 2\delta_w + \alpha.$$

*Proof.* of Lemma 4.

Before proving Lemma 4, we first show the following technical lemmas:

**Lemma 7** (Range of Probability). Define  $\delta_w$  as the average Wasserstein-1 distance between posteriors under the training and test domains, where  $\delta_w = \mathbb{E}_{X \sim \mathcal{P}(X)} \left[ W_1(\mathbb{P}_{\mathcal{P}_{\text{noisy}}}(\cdot | X), \mathbb{P}_{\mathcal{P}}(\cdot | X)) \right]$ . Then the probability  $\mathbb{P}_{\mathcal{P}}(M_f(X, y) \geq \tau_\alpha(f)) \in [1 - \alpha - \delta_w, 1 - \alpha + \delta_w]$ .

**Lemma 8** (Wasserstein-1 distance bound of expected difference (Dudley, 2018; Arjovsky et al., 2017)). For any bounded measurable function  $g : X \rightarrow \mathbb{R}$ , we have:

$$\left| \mathbb{E}_{\mathcal{P}_{\text{noisy}}}[g] - \mathbb{E}_{\mathcal{P}}[g] \right| \leq 2\|g\|_\infty W_1(\mathbb{P}_{\mathcal{P}_{\text{noisy}}}, \mathbb{P}_{\mathcal{P}}).$$

The proof of Lemma 7 is deferred to the end of this proof. Now we begin to prove Lemma 4.

Recall that  $M_f(X, \tilde{Y})$  as the gap of confidence scores, where  $|M_f(X, y)| \leq 1$  for any  $(X, \tilde{Y})$ . We also recall that  $\mathcal{L}_{cr}(f) = -\mathbb{E}_{X \sim \mathcal{P}} [M_f(X, Y) | M_f(X, Y) \geq \tau_\alpha(f)]$ .

Define  $P = \mathbb{P}_{\mathcal{P}}(M_f(X, Y) \geq \tau_\alpha(f))$ . Thus, we have:

$$\begin{aligned}
 & |\mathcal{L}_{\text{cr}}(f) - \mathcal{L}_{\text{cr}}^{\text{noisy}}(f)| \\
 &= \left| \mathbb{E}_{\mathcal{P}} [M_f(X, Y) \mid M_f(X, Y) \geq \tau_\alpha(f)] - \mathbb{E}_{\mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \cdot \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] \right| \\
 &= \left| \frac{\mathbb{E}_{\mathcal{P}} [M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]]}{\mathbb{P}_{\mathcal{P}}(M_f(X, Y) \geq \tau_\alpha(f))} - \mathbb{E}_{\mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \cdot \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] \right| \\
 &= \left| \frac{1}{P} \mathbb{E}_{\mathcal{P}} [M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]] - \mathbb{E}_{\mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] \right| \\
 &= \left| \left( \frac{1}{P} - 1 \right) \mathbb{E}_{\mathcal{P}} [M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]] + \mathbb{E}_{\mathcal{P}} [M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]] \right. \\
 &\quad \left. - \mathbb{E}_{\mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] \right| \\
 &\leq \underbrace{\left| \left( \frac{1}{P} - 1 \right) \mathbb{E}_{\mathcal{P}} [M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]] \right|}_A + \\
 &\quad \underbrace{\left| \mathbb{E}_{\mathcal{P}} [M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]] - \mathbb{E}_{\mathcal{P}_{\text{noisy}}} [M_f(X, \tilde{Y}) \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)]] \right|}_B, \tag{15}
 \end{aligned}$$

where the first inequality is due to the triangle inequality.

Then we analyze the upper bound of  $A$  and  $B$ , respectively.

**Part I: Upper bound of  $A$**

$$\begin{aligned}
 & \left| \left( \frac{1}{P} - 1 \right) \mathbb{E}_{\mathcal{P}} [M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]] \right| \\
 &\leq \left( \frac{1}{P} - 1 \right) \left| \mathbb{E}_{\mathcal{P}} [M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]] \right| \\
 &= \left( \frac{1}{P} - 1 \right) \left| \mathbb{E}_{\mathcal{P}} [M_f(X, Y)] \right| \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)] \\
 &\leq \left( \frac{1}{P} - 1 \right) \left| \mathbb{E}_{\mathcal{P}} \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)] \right| \\
 &= \left( \frac{1}{P} - 1 \right) P \\
 &= \frac{|1 - P|}{P} P \\
 &\leq \alpha + \delta_w, \tag{16}
 \end{aligned}$$

where the first inequality is due to the triangle inequality, the second inequality is due to  $|M_f(X, Y)| \leq 1$ , and the last inequality is due to the Lemma 7.

**Part II: Upper bound of  $B$**

For any fixed  $X$ , we define  $\phi_X(Y) = M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]$ . Then, we try to bound  $\|\phi_X(y)\|_\infty$  as:

$$\|\phi_X(Y)\|_\infty = \|M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)]\|_\infty \leq \|M_f(X, Y)\|_\infty \leq 1,$$

where the first inequality is due to  $\mathbb{1}[\cdot] \in [0, 1]$ , and the second inequality is due to  $|M_f(X, Y)| \leq 1$ .

Then, we rewrite  $B$  as:

$$\begin{aligned}
 & \left| \mathbb{E}_{\mathcal{P}} \left[ M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)] \right] - \mathbb{E}_{\mathcal{P}_{\text{noisy}}} \left[ M_f(X, \tilde{Y}) \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)] \right] \right| \\
 &= \left| \mathbb{E}_X \left[ \mathbb{E}_{y \sim \mathbb{P}(\cdot | X)} \left[ M_f(X, Y) \mathbb{1}[M_f(X, Y) \geq \tau_\alpha(f)] \right] - \mathbb{E}_{y \sim \mathbb{P}_{\text{noisy}}(\cdot | X)} \left[ M_f(X, \tilde{Y}) \mathbb{1}[M_f(X, \tilde{Y}) \geq \tau_\alpha(f)] \right] \right] \right| \\
 &= \left| \mathbb{E}_X \left[ \mathbb{E}_{y \sim \mathbb{P}(\cdot | X)} \left[ \phi_X(Y) \right] - \mathbb{E}_{y \sim \mathbb{P}_{\text{noisy}}(\cdot | X)} \left[ \phi_X(\tilde{Y}) \right] \right] \right| \\
 &\leq \left| \mathbb{E}_X \left[ 2W_1(\mathbb{P}_{\mathcal{P}_{\text{noisy}}}(\cdot | X), \mathbb{P}_{\mathcal{P}}(\cdot | X)) \right] \right| \\
 &= 2 \left| \mathbb{E}_X \left[ W_1(\mathbb{P}_{\mathcal{P}_{\text{noisy}}}(\cdot | X), \mathbb{P}_{\mathcal{P}}(\cdot | X)) \right] \right| \\
 &= 2\delta_w, \tag{17}
 \end{aligned}$$

where the first inequality is due to Lemma 8.

Combining inequalities (15),(16), and (17), we have that:

$$|\mathcal{L}_{\text{cr}}(f) - \mathcal{L}_{\text{cr}}^{\text{noisy}}(f)| \leq 2\delta_w + \alpha.$$

□

*Proof.* (of Lemma 7)

Recall that  $\mathcal{Y}$  denotes the label space equipped with the 0–1 distance  $d(y_1, y_2) = \mathbb{I}\{y_1 \neq y_2\}$ . For each fixed  $X$ , define  $\Phi_X(Y) = \mathbb{1}(M_f(X, Y) \geq \tau_\alpha(f))$ . Under this metric,  $\Phi_X(Y)$  is 1-Lipschitz because  $|\Phi_X(y_1) - \Phi_X(y_2)| \leq d(y_1, y_2)$ .

Denote by  $\mu_X = \mathbb{P}_{\mathcal{P}_{\text{noisy}}}(\cdot | X)$  and  $\nu_X = \mathbb{P}_{\mathcal{P}}(\cdot | X)$  the label posteriors under the noisy and clean domains, respectively. Then, we have:

$$\begin{aligned}
 & \left| \mathbb{P}_{\mathcal{P}_{\text{noisy}}}(M_f(X, \tilde{Y}) \geq \tau_\alpha(f)) - \mathbb{P}_{\mathcal{P}}(M_f(X, Y) \geq \tau_\alpha(f)) \right| \\
 &= \left| \mathbb{E}_{X \sim \mathcal{P}(X)} \left[ \mathbb{E}_{\tilde{Y} \sim \mu_X} \Phi_X(\tilde{Y}) - \mathbb{E}_{Y \sim \nu_X} \Phi_X(Y) \right] \right| \\
 &\leq \mathbb{E}_{X \sim \mathcal{P}(X)} \left[ \left| \mathbb{E}_{\tilde{Y} \sim \mu_X} \Phi_X(\tilde{Y}) - \mathbb{E}_{Y \sim \nu_X} \Phi_X(Y) \right| \right] \\
 &= \mathbb{E}_{X \sim \mathcal{P}(X)} \left[ \left| \sum_{y \in \mathcal{Y}} \Phi_X(y) (\mu_X(y) - \nu_X(y)) \right| \right] \\
 &\leq \mathbb{E}_{X \sim \mathcal{P}_X} \left[ W_1(\mu_X, \nu_X) \right] \\
 &= \mathbb{E}_{X \sim \mathcal{P}_X} \left[ W_1(\mathbb{P}_{\mathcal{P}_{\text{noisy}}}(\cdot | X), \mathbb{P}_{\mathcal{P}}(\cdot | X)) \right] = \delta_w,
 \end{aligned}$$

where the first inequality is due to Jensen’s inequality, and the last inequality is due to Kantorovich–Rubinstein duality for Wasserstein-1.

Due to  $\mathbb{P}_{\mathcal{P}_{\text{noisy}}}(M_f(X, \tilde{Y}) \geq \tau_\alpha(f)) = 1 - \alpha$ , we have:  $\left| \mathbb{P}_{\mathcal{P}}(M_f(X, Y) \geq \tau_\alpha(f)) - (1 - \alpha) \right| = \delta_w$ . □

## B ADDITIONAL EXPERIMENTS FOR MULTI-CLASS CLASSIFICATION

### B.1 Additional Experimental Setup Details

**Datasets.** To evaluate model robustness under noisy supervision, we construct asymmetric label noise on both CIFAR-100 and mini-ImageNet, simulating realistic mislabeling patterns.

**CIFAR-100** consists of 100 fine-grained classes organized into 20 coarse-grained superclasses (e.g., aquatic mammals, large carnivores). To inject structured label noise, we first build a mapping from each fine label to its corresponding coarse superclass. For noise rates  $\rho \in \{0\%, 5\%, 10\%, 20\%, 30\%, 40\%\}$ , we randomly select  $\rho n$  training samples, where  $n$  is the dataset size. Each selected label is replaced by a randomly chosen different label from the same superclass, ensuring that the corrupted label remains semantically similar to the original. The indices of noisy samples are recorded to facilitate evaluation and ablation.

**mini-ImageNet** is a 100-class image classification dataset. To inject asymmetric noise, we randomly select a fraction  $\rho = 20\%$  of training samples and replace each label  $Y$  with  $(Y + 1) \bmod 100$ , introducing a deterministic and minimal perturbation.

**FOOD101** is a 101-class image classification dataset. To inject asymmetric noise, we randomly select a fraction  $\rho = 20\%$  of training samples and replace each label  $Y$  with  $(Y + 1) \bmod 100$ , introducing a deterministic and minimal perturbation.

This form of circular asymmetric noise preserves the class index structure and ensures the new label is different from the original. The transformation is applied only to training labels, leaving the validation and test sets clean for evaluation.

**Hyperparameters for training.** We set datasets, base loss, batch size, training epochs, training parameters (learning rate, learning schedule, momentum, gamma, and weight decay),  $\lambda$  and  $\alpha$  as hyperparameter choices. We search for hyperparameters on batch size  $\in \{64, 128, 256, 512\}$ , epochs  $\in \{50, 100\}$ , learning rate ( $\eta$ )  $\in \{0.1, 0.07, 0.05, 0.03, 0.01\}$ , learning rate schedule  $\in \{[10], [10, 30], [60], [60, 80]\}$ , Momentum = 0.9, weight decay = 0.0002,  $\gamma = 0.01$ ,  $\lambda = \{0.05, 0.1, 0.15, 0.2, 0.25\}$ , and  $\alpha = \{0.05, 0.1, 0.15, 0.2, 0.25\}$  to select the best combination of hyperparameters of each methods. For GCE, we additionally scale  $\lambda$  by a factor of 0.1, resulting in  $\lambda \in \{0.005, 0.01, 0.015, 0.02, 0.025\}$ . The hyperparameters employed to get the results presented in the main paper are summarized in Table 5.

Data	Loss	Batch size	Epochs	$\eta$	lr schedule	Momentum	$\gamma$	weight decay	$\lambda$	$\alpha$
CIFAR-100	CE+CMRM	128	50	0.05	[10]	0.9	0.01	0.0002	0.1	0.15
	Focal+CMRM	128	50	0.05	[10]	0.9	0.01	0.0002	0.15	0.1
	LDAM+CMRM	128	100	0.05	[60, 80]	0.9	0.01	0.0002	0.1	0.15
	GCE+CMRM	128	50	0.05	[10]	0.9	0.01	0.0002	0.005	0.05
mini-ImageNet	CE+CMRM	512	50	0.05	[10]	0.9	0.01	0.0002	0.15	0.2
	Focal+CMRM	512	50	0.05	[10]	0.9	0.01	0.0002	0.15	0.15
	LDAM+CMRM	512	100	0.05	[60, 80]	0.9	0.01	0.0002	0.2	0.1
	GCE+CMRM	512	50	0.05	[10]	0.9	0.01	0.0002	0.0005	0.15
FOOD101	CE+CMRM	512	50	0.05	[10]	0.9	0.01	0.0002	0.15	0.15
	Focal+CMRM	512	50	0.05	[10]	0.9	0.01	0.0002	0.15	0.1
	LDAM+CMRM	512	100	0.05	[60, 80]	0.9	0.01	0.0002	0.05	0.15
	GCE+CMRM	512	50	0.05	[10]	0.9	0.01	0.0002	0.0001	0.1

Table 5: **The details we used to train our models for multi-class classification corrupted by synthetic noise with noise rate 20%.** We reported the hyperparameters that give the best accuracy. We employed SGD optimizer for all training unless specified.

## B.2 Additional Experimental Results

### Result: CMRM improves accuracy and reduces uncertainty under different types of noise.

Table 6 summarizes results on CIFAR-100 with synthetic label noise at rates  $\{0\%, 5\%, 10\%, 20\%, 30\%, 40\%\}$ , where 0% corresponds to the clean-label setting. For each objective (CE, Focal, LDAM, and GCE), we compare the Base model with its +CMRM variant. On average across all objectives and noise levels, CMRM improves accuracy by 1.34 and reduces M.APSS by 3.89%. The most notable gains occur under moderate to high noise. For example, CE and Focal achieve up to +2.14 and 2.22 accuracy improvements, while GCE shows the largest uncertainty reduction (up to -18.44% in M.APSS). Even when combined with LDAM, which already promotes margin separation, CMRM consistently provides additional accuracy improvements without increasing predictive uncertainty.

Table 7 reports the mean  $\pm$  standard deviation of Top-1 accuracy across multiple random seeds on CIFAR-100 with 20% synthetic label noise. For each objective (CE, Focal, LDAM, and GCE), we compare the Base model

Noise Rates	Metric	CE		Focal		LDAM		GCE	
		Base	+CMRM	Base	+CMRM	Base	+CMRM	Base	+CMRM
0%	ACC (%) ↑	68.82	<b>69.38</b> (+0.56)	68.32	<b>69.27</b> (+0.95)	62.94	<b>64.15</b> (+1.19)	<b>62.80</b>	<b>62.80</b> (+0)
	M.APSS ↓	<b>3.22</b>	3.37 (+4.52%)	<b>3.04</b>	3.26 (+7.24%)	<b>10.89</b>	11.06 (+1.56%)	<b>5.95</b>	<b>5.95</b> (+0%)
5%	ACC (%) ↑	67.48	<b>68.02</b> (+0.54)	66.76	<b>67.66</b> (+0.90)	61.63	<b>63.00</b> (+1.37)	63.88	<b>65.29</b> (+1.41)
	M.APSS ↓	<b>3.79</b>	4.13 (+8.97%)	<b>3.66</b>	3.82 (+4.37%)	<b>13.03</b>	13.28 (+1.92%)	<b>6.27</b>	6.30 (+0.48%)
10%	ACC (%) ↑	66.29	<b>67.32</b> (+1.03)	66.52	<b>66.82</b> (+0.30)	61.06	<b>62.28</b> (+1.22)	62.80	<b>64.73</b> (+1.93)
	M.APSS ↓	<b>4.31</b>	4.57 (+6.03%)	<b>4.28</b>	4.53 (+5.84%)	15.57	<b>15.32</b> (-1.61%)	6.55	<b>5.58</b> (-14.81%)
20%	ACC (%) ↑	65.16	<b>66.32</b> (+1.16)	64.42	<b>65.39</b> (+0.97)	59.63	<b>61.12</b> (+1.49)	62.17	<b>63.65</b> (+1.48)
	M.APSS ↓	6.67	<b>6.52</b> (-2.25%)	6.89	<b>6.61</b> (-4.06%)	17.85	<b>17.67</b> (-0.78%)	7.70	<b>6.28</b> (-18.44%)
30%	ACC (%) ↑	64.02	<b>65.22</b> (+1.2)	63.46	<b>65.33</b> (+1.87)	58.27	<b>59.56</b> (+1.29)	61.49	<b>63.23</b> (+1.74)
	M.APSS ↓	8.81	<b>7.88</b> (-10.56%)	9.12	<b>7.66</b> (-16.01%)	20.44	<b>18.63</b> (-8.86%)	6.77	<b>6.18</b> (-8.71%)
40%	ACC (%) ↑	62.63	<b>64.77</b> (+2.14)	62.65	<b>64.87</b> (+2.22)	58.25	<b>59.27</b> (+1.02)	59.34	<b>60.85</b> (+1.51)
	M.APSS ↓	10.31	<b>8.99</b> (-12.80%)	10.55	<b>9.12</b> (-13.55%)	<b>19.54</b>	20.95 (+7.22%)	7.76	<b>7.74</b> (-0.26%)

Table 6: **Top-1 accuracy (%) and marginal average prediction set size (M.APSS ↓) on CIFAR-100 datasets corrupted by synthetic noise with noise rates {0%, 5%, 10%, 20%, 30%, 40%}**. Each Base objective is paired with its +CMRM counterpart; the better value within each pair is in **bold**. Numbers in parentheses indicate the relative change (%): + denotes accuracy improvement, and - denotes M.APSS reduction compared to the corresponding Base objective. On average across all datasets and objectives, CMRM improves accuracy by 1.34 and reduces M.APSS by 3.89%.

Metric	CE		Focal		LDAM		GCE	
	Base	+CMRM	Base	+CMRM	Base	+CMRM	Base	+CMRM
ACC (%) ↑	64.74 ± 0.50	<b>65.95 ± 0.30</b> (+1.21)	64.33 ± 0.15	<b>65.26 ± 0.21</b> (+0.93)	59.64 ± 0.34	<b>60.76 ± 0.26</b> (+1.12)	62.46 ± 0.70	<b>63.60 ± 0.24</b> (+1.14)

Table 7: **Mean ± standard deviation (std) of Top-1 accuracy (%) across multiple random seeds on CIFAR-100 with synthetic label noise (noise rate 20%)**. Each base objective is paired with its +CMRM variant. The better value within each pair is shown in **bold**. Numbers in parentheses denote the absolute accuracy improvement over the corresponding base objective.

Method	CIFAR-10N (Aggre)		CIFAR-10N (Rand1)		CIFAR-10N (Rand2)		CIFAR-10N (Rand3)		CIFAR-10N (Worst)		CIFAR-100N	
	ACC(%)	M.APSS	ACC(%)	M.APSS	ACC(%)	M.APSS	ACC(%)	M.APSS	ACC(%)	M.APSS	ACC(%)	M.APSS
NI-ERM	98.69	0.904	98.80	0.902	98.65	<b>0.903</b>	98.67	0.904	95.71	0.93	83.17	1.49
NI-ERM+CMRM	<b>98.81</b>	<b>0.903</b>	<b>99.03</b>	<b>0.901</b>	<b>98.95</b>	<b>0.903</b>	<b>98.88</b>	<b>0.899</b>	<b>97.19</b>	<b>0.91</b>	<b>83.95</b>	<b>1.29</b>
	(+0.12)	(-0.11%)	(+0.23)	(-0.11%)	(+0.30)	(0%)	(+0.21)	(-0.55%)	(+1.48)	(-2.15%)	(+0.78)	(-13.42%)

Table 8: **Top-1 accuracy (%) and marginal average prediction set size (M.APSS ↓) on CIFAR-10N and CIFAR-100N corrupted by human annotation noise**. Numbers in parentheses indicate the relative change: + denotes accuracy improvement and -% denotes M.APSS reduction. CMRM consistently improves accuracy and reduces uncertainty across CIFAR-N variants, with the largest gains observed on CIFAR-10N and CIFAR-100N. On average across all datasets, CMRM improves accuracy by 0.52 and reduces M.APSS by 2.72%.

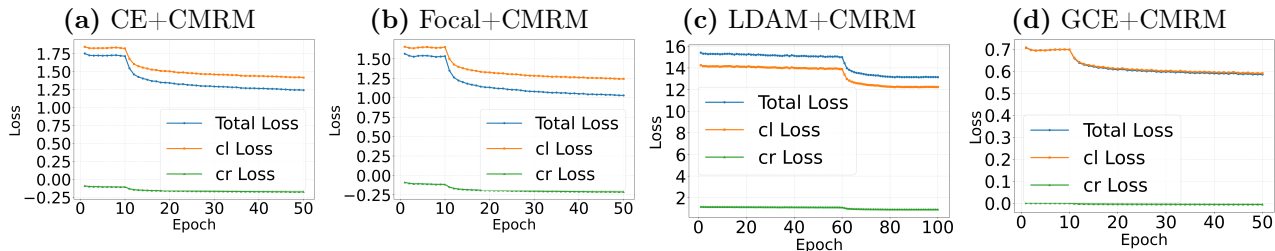


Figure 5: **training dynamics of total loss (Total), classification loss (cl), and CMRM loss (cr) for all base losses over epochs on CIFAR-100 with 20% synthetic label noise**. Subfigure (a) CE+CMRM; Subfigure (b) Focal+CMRM; Subfigure (c) LDAM+CMRM; Subfigure (d) GCE+CMRM; CMRM exhibits stable and monotonic convergence alongside standard loss components.

with its +CMRM variant. Across all objectives, CMRM consistently improves accuracy under the multi-seed evaluation. In particular, CMRM yields accuracy gains of +1.21, +0.93, +1.12, and +1.14 for CE, Focal, LDAM, and GCE, respectively. These results demonstrate that the performance improvements brought by CMRM remain stable across different random seeds.

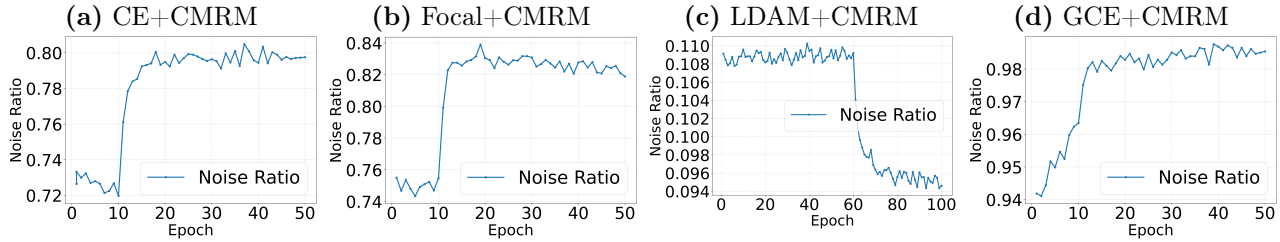


Figure 6: **The ratio of noisy samples among those filtered out by CMRM at each epoch on CIFAR-100 with 20% synthetical label noise.** Subfigure (a) CE+CMRM ( $\alpha = 0.15$ ); Subfigure (b) Focal+CMRM ( $\alpha = 0.1$ ); Subfigure (c) LDAM+CMRM ( $\alpha = 0.15$ ); Subfigure (d) GCE+CMRM ( $\alpha = 0.05$ ); CMRM consistently suppresses noisy examples by excluding low-margin samples during training.

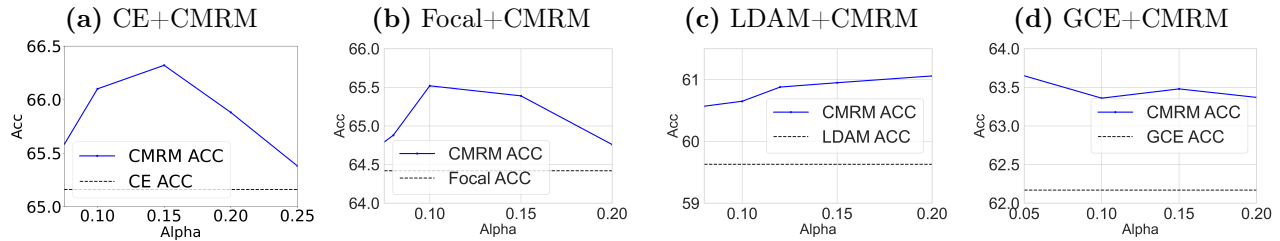


Figure 7: **The sensitivity of  $\alpha$  of different base losses on CIFAR-100 with 20% synthetical label noise.** Subfigure (a) CE+CMRM; Subfigure (b) Focal+CMRM; Subfigure (c) LDAM+CMRM; Subfigure (d) GCE+CMRM; CMRM maintains higher accuracy than CE across a range of  $\alpha$  values, indicating robustness to hyperparameter  $\alpha$ .

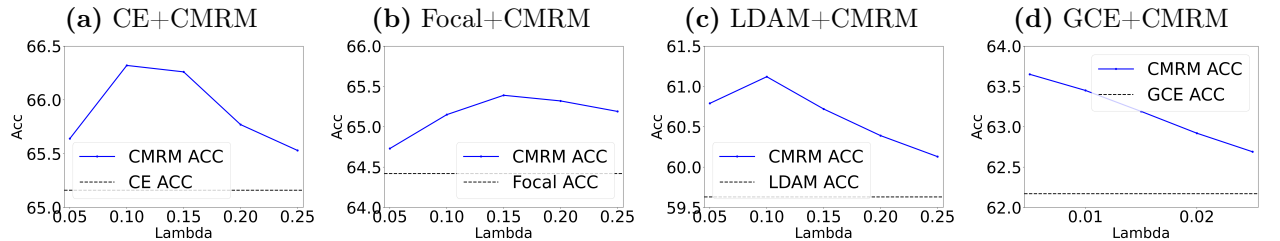


Figure 8: **The sensitivity of  $\lambda$  of different base losses on CIFAR-100 with 20% synthetical label noise.** Subfigure (a) CE+CMRM; Subfigure (b) Focal+CMRM; Subfigure (c) LDAM+CMRM; Subfigure (d) GCE+CMRM; CMRM maintains higher accuracy than CE across a range of  $\lambda$  values, indicating robustness to hyperparameter  $\lambda$ .

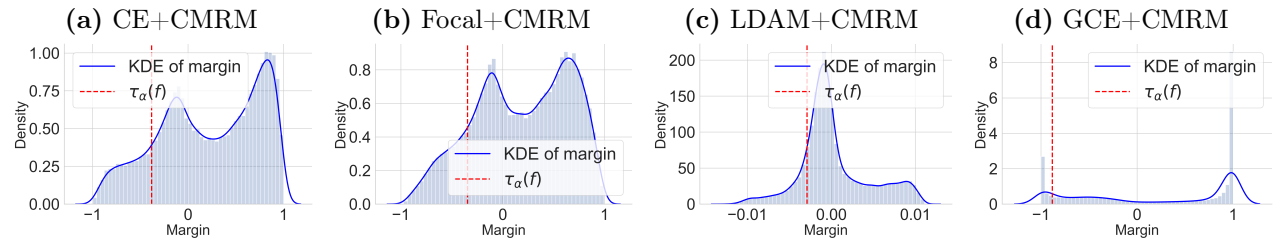


Figure 9: **the kernel density estimate (KDE) of the margin distribution of different base losses on CIFAR-100 with 20% synthetical label noise.** Subfigure (a) CE+CMRM ( $\alpha = 0.15$ ); Subfigure (b) Focal+CMRM ( $\alpha = 0.1$ ); Subfigure (c) LDAM+CMRM ( $\alpha = 0.15$ ); Subfigure (d) GCE+CMRM ( $\alpha = 0.05$ ); The vertical dashed line indicating the estimated  $\tau_\alpha(f)$ . The density curve is smooth and strictly positive around  $\tau_\alpha(f)$ , supporting the differentiability and positive-density assumption in Proposition 1.

Table 8 summarizes results on CIFAR-10N and CIFAR-100N with human-annotated label noise. Across all datasets, CMRM consistently improves both accuracy and uncertainty compared to NI-ERM. On average, CMRM increases accuracy by 0.52 and reduces M.APSS by 2.72%, indicating more confident and reliable predictions under noisy supervision. The improvements are particularly pronounced on the most challenging settings, i.e., CIFAR-10N (Worst) and CIFAR-100N, where accuracy rises by up to 1.48 and predictive uncertainty decreases

by as much as 13.42%. These consistent gains demonstrate that CMRM effectively enhances both robustness and calibration under real-world label noise.

**Result: CMRM loss convergences.** Figure 5 shows the training dynamics of the classification loss and the CMRM regularization loss. Both components decrease steadily and stabilize as training progresses, indicating smooth joint optimization. The CMRM term integrates with standard objectives and does not introduce instability or slowing down of convergence, demonstrating that CMRM can be efficiently optimized.

**Result: CMRM filters out noisy samples during training.** Figure 6 shows the fraction of noisy samples among those excluded by CMRM at each epoch. This proportion rapidly increases during the early training phase and stabilizes above 78%, indicating that CMRM consistently identifies and filters out mislabeled examples via its margin-based thresholding mechanism.

**Result: CMRM is robust to the choices of hyperparameter  $\alpha$  and  $\lambda$ .** Figure 7 and 8 examine the sensitivity of CMRM to the hyperparameter  $\alpha$  and  $\lambda$ , respectively. Across a range of  $\alpha$  and  $\lambda$  values, CMRM consistently achieves higher accuracy than CE, indicating that its performance is robust to the choice of  $\alpha$  and  $\lambda$ , and does not rely on careful hyperparameter tuning.

**Result: Assumptions in Proposition 1 are empirically valid.** Figure 9 presents the kernel density estimate (KDE) of the margin distribution, with the vertical dashed line indicating the estimated  $\tau_\alpha(f)$ . The density curve is smooth and strictly positive in the neighborhood of  $\tau_\alpha(f)$ , supporting the differentiability and positive-density assumption in Proposition 1.

## C ADDITIONAL EXPERIMENTS FOR BINARY CLASSIFICATION

### C.1 Additional Experimental Setup Details

**Datasets.** We evaluate our methods on four standard binary classification benchmarks, three of which are sourced from the UCI Machine Learning Repository:

- **EMAIL:** This dataset contains features extracted from emails and aims to classify whether an email is spam or not. We treat spam emails as the positive class (label 1) and non-spam emails as the negative class (label 0). Following prior work, we invert the original labels in the dataset to conform to this definition.
- **ADULT:** This dataset involves predicting whether an individual’s income exceeds 50K based on demographic and employment features. We define the positive class as individuals earning  $\leq 50K$  (label 1) and the negative class as those earning  $> 50K$  (label 0), effectively focusing on identifying lower-income individuals.
- **CREDIT:** This dataset contains information on credit card clients and whether they defaulted on their payment in the following month. We define the positive class as non-defaulting clients (label 1) and the negative class as clients who defaulted (label 0), inverting the original label to emphasize reliable borrowers.

**Hyperparameters for training.** We set datasets, base losses, base models, learning rate,  $\lambda^+$ ,  $\lambda^-$ ,  $\alpha^+$ , and  $\alpha^-$  as hyperparameter choices. We search for hyperparameters on learning rate ( $\eta$ )  $\in \{0.001, 0.005, 0.01\}$ ,  $\lambda^+ \in \{0, 0.05, 0.1, 0.15, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1.1, 1.2, 1.3, 1.4, 1.5\}$ ,  $\lambda^- \in \{0, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1.1, 1.2, 1.3, 1.4, 1.5\}$ ,  $\alpha^+ \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ , and  $\alpha^- \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$  to select the best combination of hyperparameters of each methods. The hyperparameters employed to get the results presented in the main paper are summarized in Table 10.

### C.2 Additional Experimental Results

**Result: CMRM improves robustness for binary classification.** Table 9 reports results on the Adult, Email, and Credit datasets under a 20% label noise setting. Across all base models (LR, Focal, SVM, GCE), CMRM consistently improves ranking (AUROC, AUPRC) and classification (ACC, FPR) performance, while maintaining comparable FNR. It also reduces predictive uncertainty, as shown by lower M.APSS, PC APSS, and NC APSS values. These results demonstrate that CMRM enhances model robustness and uncertainty estimation under noisy supervision without requiring any noise priors.

**Conformal Margin Risk Minimization: An Envelope Framework for Robust Learning under Label Noise**

Dataset	Method	Measure							
		AUROC ( $\uparrow$ )	AUPRC ( $\uparrow$ )	FNR ( $\downarrow$ )	FPR ( $\downarrow$ )	ACC ( $\uparrow$ )	M.APSS ( $\downarrow$ )	PC APSS ( $\downarrow$ )	NC APSS ( $\downarrow$ )
Adult	LR	0.784	0.885	<b>0.073</b>	0.571	0.802	1.223	1.154	1.432
	LR + CMRM	<b>0.852</b>	<b>0.925</b>	0.082	<b>0.422</b>	<b>0.833</b>	<b>1.209</b>	<b>1.109</b>	<b>1.308</b>
	Focal	0.809	0.890	0.136	0.388	0.801	1.257	1.224	1.356
	Focal + CMRM	<b>0.872</b>	<b>0.942</b>	<b>0.128</b>	<b>0.324</b>	<b>0.823</b>	<b>1.221</b>	<b>1.148</b>	<b>1.295</b>
	SVM	0.808	0.925	<b>0.029</b>	0.807	0.776	1.276	1.370	1.512
	SVM + CMRM	<b>0.847</b>	<b>0.937</b>	0.048	<b>0.585</b>	<b>0.817</b>	<b>1.199</b>	<b>1.322</b>	<b>1.343</b>
	GCE	0.819	0.904	<b>0.119</b>	0.424	<b>0.804</b>	1.286	<b>1.176</b>	1.396
	GCE + CMRM	<b>0.846</b>	<b>0.928</b>	0.172	<b>0.286</b>	0.800	<b>1.273</b>	1.207	<b>1.340</b>
Email	LR	0.831	0.869	<b>0.246</b>	0.200	0.773	1.405	1.415	1.392
	LR + CMRM	<b>0.875</b>	<b>0.910</b>	0.259	<b>0.129</b>	<b>0.793</b>	<b>1.335</b>	<b>1.247</b>	<b>1.291</b>
	Focal	0.833	0.858	<b>0.228</b>	0.208	0.780	1.404	1.449	1.345
	Focal + CMRM	<b>0.907</b>	<b>0.916</b>	0.246	<b>0.080</b>	<b>0.822</b>	<b>1.235</b>	<b>1.152</b>	<b>1.202</b>
	SVM	0.952	0.964	<b>0.043</b>	0.219	0.885	1.029	1.003	1.001
	SVM + CMRM	<b>0.954</b>	<b>0.967</b>	0.060	<b>0.125</b>	<b>0.913</b>	<b>0.975</b>	<b>0.996</b>	<b>0.993</b>
	GCE	0.931	0.925	<b>0.069</b>	0.101	0.918	<b>0.976</b>	<b>0.977</b>	<b>0.975</b>
	GCE + CMRM	<b>0.938</b>	<b>0.933</b>	0.074	<b>0.082</b>	<b>0.922</b>	0.984	0.982	0.986
Credit	LR	0.690	0.866	0.134	0.600	0.764	1.468	1.414	1.522
	LR + CMRM	<b>0.714</b>	<b>0.877</b>	<b>0.121</b>	<b>0.565</b>	<b>0.782</b>	<b>1.400</b>	<b>1.365</b>	<b>1.436</b>
	Focal	0.673	0.862	<b>0.169</b>	0.588	0.739	<b>1.507</b>	<b>1.467</b>	1.547
	Focal + CMRM	<b>0.677</b>	<b>0.861</b>	0.176	<b>0.547</b>	<b>0.743</b>	1.527	1.543	<b>1.511</b>
	SVM	0.688	0.856	<b>0.079</b>	0.618	0.803	<b>1.387</b>	<b>1.335</b>	1.438
	SVM + CMRM	<b>0.711</b>	<b>0.869</b>	0.084	<b>0.553</b>	<b>0.813</b>	1.421	1.437	<b>1.405</b>
	GCE	0.671	0.861	<b>0.134</b>	0.845	0.751	1.446	1.366	1.526
	GCE + CMRM	<b>0.701</b>	<b>0.875</b>	<b>0.115</b>	<b>0.627</b>	<b>0.773</b>	<b>1.417</b>	1.339	<b>1.495</b>

Table 9: **Performance comparison on binary classification.** We evaluate three datasets (Adult, Email, and Credit) under a 20% label noise setting. Models are assessed across accuracy (ACC), ranking (AUROC, AUPRC), calibration (FNR, FPR), and uncertainty metrics: marginal average prediction set size (M. APSS), positive-class prediction set size (PC APSS), and negative-class prediction set size (NC APSS).  $\uparrow$  means higher is better;  $\downarrow$  lower is better. Best results for each base model (LR, Focal, SVM, GCE) are in **bold**. Our method consistently improves uncertainty estimation (M.APSS, PC APSS, NC APSS), ranking metrics (AUROC, AUPRC), and accuracy (ACC, FPR), while slightly increasing FNR.

Data	loss	Architecture	Epochs	$\eta$	$\lambda^+$	$\lambda^-$	$\alpha^+$	$\alpha^-$
Adult	LR+CMRM	MLP	150	0.01	0.4	0.5	0.3	0.1
	Focal+CMRM	MLP	150	0.01	0.1	0.4	0.3	0.5
	SVM+CMRM	SVM	150	0.01	0.15	1.5	0.2	0.5
	GCE+CMRM	MLP	150	0.001	0.0	0.2	0.1	0.5
Email	LR+CMRM	MLP	150	0.01	0.7	0.2	0.2	0.2
	Focal+CMRM	MLP	150	0.01	0.4	0.6	0.1	0.1
	SVM+CMRM	SVM	150	0.01	0.3	1.0	0.3	0.3
	GCE+CMRM	MLP	150	0.01	0.9	0.7	0.4	0.4
Credit	LR+CMRM	MLP	150	0.005	0.0	0.3	0.1	0.1
	Focal+CMRM	MLP	150	0.001	0.05	0.3	0.3	0.1
	SVM+CMRM	SVM	150	0.001	0.0	1.0	0.1	0.5
	GCE+CMRM	MLP	150	0.0005	0.0	0.3	0.1	0.1

Table 10: **The details we used to train our models for binary classification.** We reported the hyperparameters that give the best combination of metrics.

**CMRM Loss convergences.** To verify the stability and trainability of our method, we monitor the learning dynamics of CMRM during optimization. As shown in the top row of Figure 10, both the classification loss and CMRM loss steadily decrease and stabilize, indicating that the joint objective converges smoothly. The CMRM regularization term integrates seamlessly with standard classification training and does not introduce optimization instability or slow down convergence. This confirms that CMRM can be efficiently optimized using standard gradient-based methods and is compatible with commonly used loss functions.

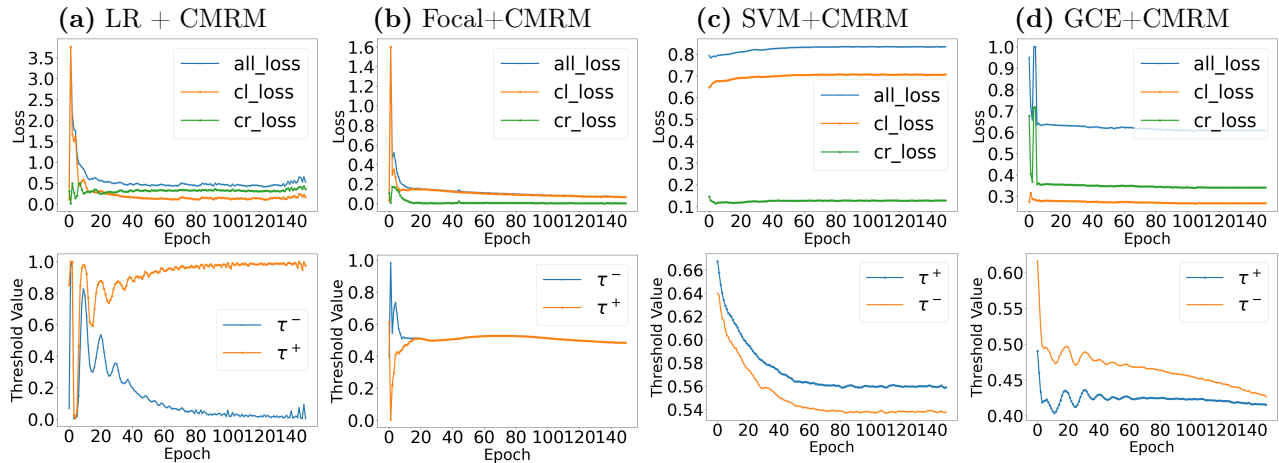


Figure 10: **Dynamics of losses (Top Row) and two thresholds ( $\tau^+$  and  $\tau^-$ , Bottom Row) for binary classification of all base loss (LR, Focal, SVM, GCE) with CMRM on Email datasets. Top row shows the training dynamics of total loss (all loss), classification loss (cl loss), and CMRM regularization loss (CMRM loss) over epochs. CMRM exhibits stable and monotonic convergence alongside standard loss components. Bottom row shows  $\tau^+$  (negative class threshold) and  $\tau^-$  (positive class threshold) of LR+CMRM during training. The separation between the thresholds increases, indicating that CMRM actively maximizes the margin between favorable and unfavorable classes.**

**Margin  $\tau^- - \tau^+$  of CMRM grows over Time.** The bottom row of Figure 10 illustrates the evolution of the class-conditional thresholds  $\tau^+$  and  $\tau^-$  over training. These thresholds define the CMRM margin region, with  $\tau^-$  indicating the lower bound for confident positives and  $\tau^+$  the upper bound for confident negatives. As training progresses, we observe that  $\tau^-$  increases while  $\tau^+$  decreases, leading to an expanding margin between the two thresholds. This demonstrates that CMRM successfully enforces separation between favorable and unfavorable classes in high-confidence regions, which is critical for robustness under posterior shift.