

CONFIDENCE DIFFERENCE REFLECTS VARIOUS SUPERVISED SIGNALS IN CONFIDENCE-DIFFERENCE CLASSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Training a precise binary classifier with limited supervision in weakly supervised learning scenarios holds considerable research significance in practical settings. Leveraging pairwise unlabeled data with confidence differences has been demonstrated to outperform learning from pointwise unlabeled data. We theoretically analyze the various supervisory signals reflected by confidence differences in confidence difference (ConfDiff) classification and identify challenges arising from noisy signals when confidence differences are small. To address this, we partition the dataset into two subsets with distinct supervisory signals and propose a consistency regularization-based risk estimator to encourage similar outputs for similar instances, mitigating the impact of noisy supervision. We further derive and analyze its estimation error bounds theoretically. Extensive experiments on benchmark and UCI datasets demonstrate the effectiveness of our method. Additionally, to effectively capture the influence of real-world noise on the confidence difference, we artificially perturb the confidence difference distribution and demonstrate the robustness of our method under noisy conditions through comprehensive experiments.

1 INTRODUCTION

Weakly supervised learning is an essential research field in machine learning, focusing on training accurate predictive models under conditions of low supervision or imprecise labeling. Due to the difficulty of obtaining precise supervision in real-world scenarios, weakly supervised learning holds significant research value and significance for effectively leveraging limited available supervision information. Consequently, the field of weakly supervised learning has increasingly attracted attention from experts and scholars in recent years, leading to the emergence of many typical weakly supervised learning methods, such as multi-instance learning [32; 30; 24; 19], positive and unlabeled (PU) learning [10; 5; 31; 16; 23], and others.

A prevalent idea in weakly supervised classification involves maximizing the utilization of pointwise weakly supervised information [4], thereby prompting the development of various techniques based on soft labels [18; 26], mixup [28; 21; 27; 9; 7; 13], and others. Nevertheless, it is undeniable that annotating pointwise information in real-world classification problems is a complex and laborious task, further compounded by the personal biases of annotators which frequently exacerbate the probability of inaccuracies. In such scenarios, pairwise comparison information between data points may be more readily obtainable in real-world settings than pointwise information, and it often exhibits greater resistance to biases compared to pointwise semi-supervised information [1]. For instance, in medical diagnosis, accurately determining whether a patient has a disease solely based on their presented symptoms is challenging. However, comparing the symptoms of this patient with those of others provides more accessible information and reduces the probability of misdiagnosis. Extensive research has been conducted on pairwise analysis in numerous binary classification problems, leading to the development of risk minimization functions capable of inducing binary classifiers across various combinations of pairwise similarities, dissimilarities, and unlabeled data [1; 20; 14; 15; 22].

In recent work, pairwise comparison (Pcomp) classification has shown that in tackling difficult point labeling tasks, people can more easily gather comparative information between two examples,

constituting a form of weakly supervised information [4]. However, in real-world application scenarios, individuals may not only distinguish which of two examples is more likely to be classified as positive over the other but also gauge the extent of the disparity in their confidence levels regarding positivity. In light of this framework, Wang et al. introduced a new pairwise weakly supervised classification problem called confidence-difference (ConfDiff) classification, and proposed the corresponding ConfDiff method [22]. To establish confidence difference, the ConfDiff method first utilizes binary-labeled data to train a probability classifier. Subsequently, unlabeled data pairs are fed into the classifier to generate posterior probabilities, from which confidence difference are computed based on the differences between these posterior probabilities. However, through the analysis of the various supervised signals in the ConfDiff method, we identify that ConfDiff method encourages unlabeled data pairs to predict opposite classes from both experimental and theoretical perspectives. This prediction direction is valid when the confidence difference is large. However, when the confidence difference is small, the instances may belong to the same or different classes, and such a predictive tendency may lead to samples from the same class being incorrectly classified as belonging to different classes, thereby introducing noisy supervisory signals.

To handle this problem, in this paper, we concentrate on mitigating the impact of inaccurate predictions when confidence differences are small. Specifically, we analyze the different supervised signals induced by varying confidence differences in the ConfDiff method. We find that pairwise instances with small confidence differences tend to introduce noisy supervised signals, while those with larger confidence differences provide more reliable supervision. Based on this observation, we propose a ConfDiff classification method that incorporates consistency regularization. By partitioning the dataset based on the accuracy of predictive information, we introduce a consistency regularization term for the subset with relatively precise predictions, encouraging the model to produce similar outputs for pairs with small confidence differences. Meanwhile, for the subset with relatively imprecise predictions, we preserve the benefit of reliable supervised signals. Experimental results demonstrate that our method outperforms existing baselines in most cases and exhibits strong robustness even under artificial noise interference.

In summary, this paper’s key contributions can be outlined as follows:

- We introduce a method for ConfDiff classification which aims to enhance the accuracy of weakly supervised classification by constructing risk estimator through **C**onsistency **R**isk and **C**onsistency **R**egularization (CRCR).
- We theoretically analyze various supervised signals reflected by different confidence differences in ConfDiff classification. Additionally, we theoretically estimate the error bounds of our proposed method.
- We validate the effectiveness of our method through experiments by comparing it with existing baselines on datasets of varying scales. In addition, the robustness of our method is further validated under the influence of artificially added noise.

2 PRELIMINARIES

In this section, we briefly review the problem definitions of binary classification, binary classification with soft labels, and ConfDiff classification.

Formulation of binary classification Binary classification is a typical task in the field of supervised learning, where the goal is to induce a classifier to partition the data space into two categories. Formally, let $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \{-1, +1\}$ be the d -dimensional feature space and label space, respectively. The dataset $\mathcal{D}_{\text{BC}} = \mathcal{D}_{\text{BC}}^p \cup \mathcal{D}_{\text{BC}}^n$ for binary classification consists of a positive dataset $\mathcal{D}_{\text{BC}}^p$ and a negative dataset $\mathcal{D}_{\text{BC}}^n$:

$$\begin{aligned} \mathcal{D}_{\text{BC}}^p &= \{(\mathbf{x}_i^p \in \mathcal{X}, y_i^p = +1)\}_{i=1}^{n_p}, \mathbf{x}_i^p \stackrel{i.i.d.}{\sim} p(\mathbf{x}|y = +1), \\ \mathcal{D}_{\text{BC}}^n &= \{(\mathbf{x}_i^n \in \mathcal{X}, y_i^n = -1)\}_{i=1}^{n_n}, \mathbf{x}_i^n \stackrel{i.i.d.}{\sim} p(\mathbf{x}|y = -1), \end{aligned}$$

where n_p and n_n denote the number of positive and negative instances, respectively. Let π denotes the class prior $p(y = +1)$ and $\ell : \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ denotes a binary loss function. Then binary classification induces a classifier $g : \mathcal{X} \rightarrow \mathbb{R}$ from \mathcal{D}_{BC} by minimizing the following classification risk:

$$R(g) = \pi \mathbb{E}_{p(\mathbf{x}|y=+1)}[\ell(g(\mathbf{x}), +1)] + (1 - \pi) \mathbb{E}_{p(\mathbf{x}|y=-1)}[\ell(g(\mathbf{x}), -1)]. \quad (1)$$

Formulation of binary classification with soft labels In binary classification, soft labels typically represent the confidence of each sample belonging to the positive class. Moreover, several studies have shown that using soft labels rather than hard labels can more accurately reflect the data distribution, thus enhancing the accuracy of training binary classifiers. Formally, let q_i denotes the positive confidence of \mathbf{x}_i , the dataset $\mathcal{D}_{\text{BC-soft}}$ for binary classification can be defined as follows:

$$\mathcal{D}_{\text{BC-soft}} = \{(\mathbf{x}_i, q_i)\}_{i=1}^n, \mathbf{x}_i \stackrel{i.i.d.}{\sim} p(\mathbf{x}), q_i = p(y_i = +1|\mathbf{x}_i),$$

where $p(\mathbf{x}) = \pi p(\mathbf{x}|y = +1) + (1 - \pi)p(\mathbf{x}|y = -1)$. Subsequently, the risk minimization objective function for binary classification with soft labels can be reformulated into the following form:

$$R_{\text{BC-soft}}(g) = \mathbb{E}_{p(\mathbf{x})}[q\ell(g(\mathbf{x}), +1) + (1 - q)\ell(g(\mathbf{x}), -1)]. \quad (2)$$

Formulation of confidence-difference (ConfDiff) classification Given that pairwise supervision is typically more accessible than pointwise supervision and it's feasible to not only determine which sample in an unlabeled data pair is more likely positive but also quantify the confidence difference between them in practical scenarios, ConfDiff classification precisely serves as a weakly supervised classification tailored to address this scenario. It specifically deals with weakly supervised classification problems where training data comprises only pairwise unlabeled data and the confidence difference associated with each pair. Formally, let $c_i = c(\mathbf{x}_i, \mathbf{x}'_i) = p(y'_i = +1|\mathbf{x}'_i) - p(y_i = +1|\mathbf{x}_i)$ be the confidence difference between pairwise unlabeled data $(\mathbf{x}_i, \mathbf{x}'_i)$ drawn from an independent identically distribution probability density $p(\mathbf{x}, \mathbf{x}') = p(\mathbf{x})p(\mathbf{x}')$. Considering a pairwise dataset \mathcal{D} drawn from the pairwise unlabeled data and the confidence differences between them:

$$\mathcal{D}_{\text{CD}} = \{(\mathbf{x}_i, \mathbf{x}'_i, c_i)\}_{i=1}^n, \mathbf{x}_i \stackrel{i.i.d.}{\sim} p(\mathbf{x}), \mathbf{x}'_i \stackrel{i.i.d.}{\sim} p(\mathbf{x}').$$

In a recent study, Wang et al. tackled the ConfDiff classification problem in such challenging scenarios [22]. They deduced an unbiased risk estimator for confidence-difference classification from Eq. 1 and trained a binary classifier solely utilizing unlabeled data and confidence differences by minimizing it. The classification risk can be expressed as:

$$R_{\text{CD}}(g) = \frac{1}{2}\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[\mathcal{L}(\mathbf{x}, \mathbf{x}') + \mathcal{L}(\mathbf{x}', \mathbf{x})], \quad (3)$$

where $\mathcal{L}(\mathbf{x}, \mathbf{x}') = (\pi - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), +1) + (1 - \pi - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}'), -1)$. Then Eq. 3 can be refined as follows:

$$R_{\text{CD}}(g) = \frac{1}{2}\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')}[(\pi - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), +1) + (1 - \pi - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}'), -1) \\ + (\pi + c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}'), +1) + (1 - \pi + c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), -1)]. \quad (4)$$

3 THE PROPOSED METHOD

In this section, we introduce the proposed noisy ConfDiff method named CRCR.

3.1 ANALYSIS OF THE CONFDIFF METHOD

In the ConfDiff method, pairwise instances with confidence differences smaller than 0.5 are prone to introducing noise, while those with larger confidence differences (greater than 0.5) are considered to provide stronger and more reliable supervised signals. To explain this, we consider the general form of many commonly used losses for the prediction function $g(x)$ and target y [29]:

$$\mathcal{L} = \{\ell(g(\mathbf{x}), y) | \ell(g(\mathbf{x}), y) = h(g(\mathbf{x})) - yg(\mathbf{x}) \text{ for some function } h\}, \quad (5)$$

Substituting the form of the loss function from Eq.5 into Eq.4, then the classification risk of ConfDiff method can be rewritten as follows and the proof details are presented in the Appendix B:

$$R_{\text{CD}}(g) = \frac{1}{2}\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[\left(\frac{1}{2} - c(\mathbf{x}, \mathbf{x}') \right) \ell(g(\mathbf{x}), +1) + \left(\frac{1}{2} + c(\mathbf{x}, \mathbf{x}') \right) \ell(g(\mathbf{x}'), +1) \right] \\ + \frac{1}{2}\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[\left(\frac{1}{2} + c(\mathbf{x}, \mathbf{x}') \right) \ell(g(\mathbf{x}), -1) + \left(\frac{1}{2} - c(\mathbf{x}, \mathbf{x}') \right) \ell(g(\mathbf{x}'), -1) \right] \\ + \frac{1}{2}\mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[(1 - 2\pi)(g(\mathbf{x}) + g(\mathbf{x}')) \right]. \quad (6)$$

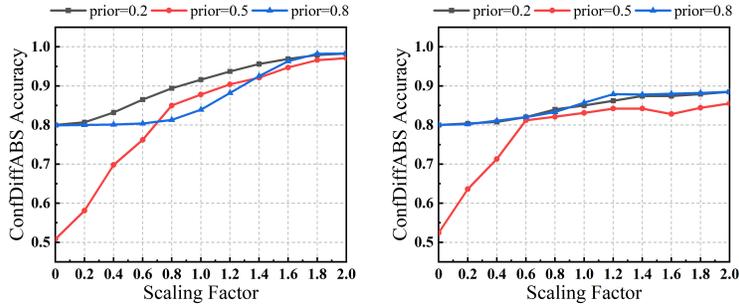


Figure 1: The Accuracy for the binary classifier concerning different proportion of pairwise data with $|c(\mathbf{x}, \mathbf{x}')| > 0.5$ on two benchmark datasets MNIST (left) and CIFAR-10 (right). (The value of the x-axis values $\times \min(\pi, 1 - \pi)$ denotes the proportion of pairwise instances with $|c(\mathbf{x}, \mathbf{x}')| > 0.5$.)

where the first and second terms denote the pairwise instance $(\mathbf{x}, \mathbf{x}')$ contrastive losses for positive and negative class predictions, respectively; and the third term serves as a regularization. We first analyze the critical components of the first term, where the weights $\frac{1}{2} - c(\mathbf{x}, \mathbf{x}')$ and $\frac{1}{2} + c(\mathbf{x}, \mathbf{x}')$ determine the contributions of \mathbf{x} and \mathbf{x}' to the positive class prediction loss, respectively. These weights exhibit an inherent balance, as their sum equals 1, indicating that $\frac{1}{2}$ serves as a boundary distinguishing the prediction directions. The weights lie on opposite sides of this boundary, ensuring that one of \mathbf{x} or \mathbf{x}' is encouraged to predict more strongly as the positive class, while the other is encouraged to weaken its positive class tendency (i.e., predict as the negative class). In other words, the first loss term ensures \mathbf{x} and \mathbf{x}' to adjust their predictions in opposite directions, thereby emphasizing the predictive divergence of pairwise instances in the positive class predictions. Similarly, the second loss term forces to diverge in their predictions for the negative class.

Referring to the definition of $c(\mathbf{x}, \mathbf{x}')$, if $|c(\mathbf{x}, \mathbf{x}')| > 0.5$, \mathbf{x} and \mathbf{x}' must belong to different classes; and if $|c(\mathbf{x}, \mathbf{x}')| \leq 0.5$, \mathbf{x} and \mathbf{x}' can belong to the same class or different classes, as the posterior difference is insufficient to surpass the classification threshold. So the prediction trend encouraged by R_{CD} holds correctly for pairwise instances with $|c(\mathbf{x}, \mathbf{x}')| > 0.5$. However, when $|c(\mathbf{x}, \mathbf{x}')| \leq 0.5$, the prediction trend may lead to samples from the same class being predicted as belonging to different classes, introducing erroneous supervisory signals. Accordingly, we consider that the pairwise instances whose confidence difference are greater than 0.5 contain more supervised signals, but the other ones may result in noisy signals in the existing ConfDiff method.

To further validate this perspective, we conduct experiments on the MNIST and CIFAR-10 by varying the proportion of the pairwise instances with $|c(\mathbf{x}, \mathbf{x}')| > 0.5$. The empirical results (see in Figure 1) illustrate the accuracy of the binary classifier under different proportions of the pairwise instances with $|c(\mathbf{x}, \mathbf{x}')| > 0.5$. We observe a positive correlation between classification accuracy and the proportion value. Notably, when the proportion is 0, the classifier accuracy is approximately 0.5, indicating that the classifier performs nearly at random. These findings demonstrate that the pairwise instances with $|c(\mathbf{x}, \mathbf{x}')| > 0.5$ provide stronger and more reliable supervised signals and dominate the contribution to R_{CD} .

3.2 CRCR METHOD

Based on the discussion in Section 3.1, it is demonstrated that noise signals is introduced when $|c(\mathbf{x}, \mathbf{x}')| \leq 0.5$, while it remains more supervised signals when $|c(\mathbf{x}, \mathbf{x}')| > 0.5$. To address it, we propose setting a threshold θ to partition the dataset into two subsets: one with relatively precise predictive information (denoted as D^S) and the other with comparatively imprecise predictive information (denoted as D^C). For D^C , we aim to provide additional information to guide the predictions of pairwise instances toward the correct direction. Specifically, for pairwise instances with small confidence differences, we encourage the model to produce more similar outputs for these pairs. To achieve this, we introduce a consistency regularization term that encourages alignment between the confidence difference and the model’s outputs. Meanwhile, for D^S , we retain the original strategy to preserve the accuracy of predictions driven by this strong guidance. Our objective is to

induce a classifier $g: \mathbb{R}^d \rightarrow \mathcal{Y}$ from \mathcal{D} by minimizing the expected risk with respect to the data distribution:

$$R_{\text{CRCR}}(g) = \frac{1}{2} \mathbb{E}_{p_{\mathcal{D}^S}(\mathbf{x}, \mathbf{x}')} [(\pi - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), +1) + (1 - \pi - c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}'), -1) \\ + (\pi + c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}'), +1) + (1 - \pi + c(\mathbf{x}, \mathbf{x}'))\ell(g(\mathbf{x}), -1)] \\ + \alpha \mathbb{E}_{p_{\mathcal{D}^C}(\mathbf{x}, \mathbf{x}')} \left[\left(\frac{1}{\log(|c(\mathbf{x}, \mathbf{x}')| + \varepsilon)} \right) \cdot \|g(\mathbf{x}) - g(\mathbf{x}')\|_2 \right], \quad (7)$$

where α denotes the parameter of the consistency regularization term, and $\varepsilon = 1.1$ is a smoothing parameter introduced to mitigate numerical issues when $|c(\mathbf{x}, \mathbf{x}')|$ approaches or equals zero. Let $|\mathcal{D}^S| = n_1$ and $|\mathcal{D}^C| = n_2$. Then the risk estimator can be expressed as follows:

$$\hat{R}_{\text{CRCR}}(g) = \frac{1}{2n_1} \sum_{i=1}^{n_1} \left((\pi - c_i)\ell(g(\mathbf{x}_i), +1) \right) + (1 - \pi - c_i)\ell(g(\mathbf{x}'_i), -1) + (\pi + c_i)\ell(g(\mathbf{x}'_i), +1) \\ + (1 - \pi + c_i)\ell(g(\mathbf{x}_i), -1) \Big) + \frac{\alpha}{n_2} \sum_{i=1}^{n_2} \left(\frac{1}{\log(|c_i| + \varepsilon)} \cdot \|g(\mathbf{x}_i) - g(\mathbf{x}'_i)\|_2 \right). \quad (8)$$

3.3 ANALYSIS OF ERROR BOUND

Assuming there exists a constant C_g such that $\sup_{g \in G} \|G\|_\infty \leq C_g$, and another constant C_ℓ such that $\sup_{|z|} \ell(z, y) \leq C_g$ and $\ell(z, y) \leq C_\ell$. Additionally, we presume the binary loss function $\ell(z, y)$ to be Lipschitz continuous with respect to both z and y , and to have a Lipschitz constant denoted by L_ℓ . $\mathfrak{R}_{n_1}(\mathcal{G})$ and $\mathfrak{R}_{n_2}(\mathcal{G})$ denote the Rademacher complexity of unlabeled data \mathcal{G} with size n_1 and n_2 , respectively.

Theorem 1. *Let $g^* = \arg \min_{g \in \mathcal{G}} R(g)$ is the minimizer of the true classification risk in Eq.1 and $\hat{g}_{\text{CRCR}} = \arg \min_{g \in \mathcal{G}} \hat{R}_{\text{CRCR}}(g)$ denotes the minimizer of the risk form in Eq.8. Then for any $\delta > 0$, we believe that the following expression holds with a probability at least $1 - \delta$:*

$$R(\hat{g}_{\text{CRCR}}) - R(g^*) \leq 8L_\ell \mathfrak{R}_{n_1}(\mathcal{G}) + \frac{4\alpha}{\log(\varepsilon)} \mathfrak{R}_{n_2}(\mathcal{G}) \\ + \left(\frac{4C_\ell}{n_1} + \left| \frac{1}{\log(\varepsilon)} - \frac{1}{\log(\theta + \varepsilon)} \right| \frac{4\alpha C_g}{n_2} \right) \sqrt{2n \ln(2/\delta)}. \quad (9)$$

Due to the space limitation, the proof details are presented in the Appendix A. As $n_1, n_2 \rightarrow \infty$, the Rademacher complexities $\mathfrak{R}_{n_1}(\mathcal{G})$ and $\mathfrak{R}_{n_2}(\mathcal{G})$ decrease to zero, and the third term involving \sqrt{n}/n_1 and \sqrt{n}/n_2 also diminishes. Furthermore, the convergence rates of $\mathfrak{R}_{n_1}(\mathcal{G})$ and $\mathfrak{R}_{n_2}(\mathcal{G})$ are $O(1/\sqrt{n_1})$ and $O(1/\sqrt{n_2})$, while the third term's rate is dominated by $O(\sqrt{n}/n_1)$ and $O(\sqrt{n}/n_2)$. Consequently, as $n \rightarrow \infty$, $R(\hat{g}_{\text{CRCR}}) \rightarrow R(g^*)$, and the overall convergence rate is characterized by $O(\max(\sqrt{n}/n_1, \sqrt{n}/n_2))$.

3.4 EMPIRICAL RISK CORRECTION

It can potentially lead to severe overfitting problems when the empirical risk becomes negative due to the application of a revised unbiased form. Fortunately, risk correction functions $f(\cdot)$ can be utilized to mitigate this issue. Examples include the absolute value function or the rectified linear unit (ReLU) function. Consequently, the corrected risk estimator can be expressed as follows:

$$\tilde{R}_{\text{CRCR}}(g) = \frac{1}{2n_1} f \left(\sum_{i=1}^{n_1} (\pi - c_i)\ell(g(\mathbf{x}_i), +1) \right) + \frac{1}{2n_1} f \left(\sum_{i=1}^{n_1} (1 - \pi - c_i)\ell(g(\mathbf{x}'_i), -1) \right) \\ + \frac{1}{2n_1} f \left(\sum_{i=1}^{n_1} (\pi + c_i)\ell(g(\mathbf{x}'_i), +1) \right) + \frac{1}{2n_1} f \left(\sum_{i=1}^{n_1} (1 - \pi + c_i)\ell(g(\mathbf{x}_i), -1) \right) \\ + \alpha \frac{1}{n_2} f \left(\sum_{i=1}^{n_2} \left(\frac{1}{\log(|c_i| + \varepsilon)} \cdot \|g(\mathbf{x}_i) - g(\mathbf{x}'_i)\|_2 \right) \right). \quad (10)$$

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

Table 1: Detailed characteristics of datasets.

Dataset	#Instance	#Trainset	#Testset	#Fea	Pos Class	Neg Class	Backbone
MNIST	70,000	15,000	5,000	28×28	0,2,4,6,8	1,3,5,7,9	3-layer MLP
F-MNIST	70,000	15,000	5,000	28×28	0,2,4,6,8	1,3,5,7,9	3-layer MLP
K-MNIST	70,000	15,000	5,000	28×28	0,2,4,6,8	1,3,5,7,9	3-layer MLP
CIFAR-10	60,000	10,000	5,000	$3 \times 32 \times 32$	2,3,4,5,6,7	0,1,8,9	ResNet-34
Optdigits	5,620	1,200	1,125	62	0,2,4,6,8	1,3,5,7,9	Linear
Pendigits	10,992	2,500	2,199	16	0,2,4,6,8	1,3,5,7,9	Linear

Additionally, we report corresponding versions in the experiments that utilized absolute risk correction function (CRCR-ABS) and ReLU risk correction function (CRCR-ReLU).

4 EXPERIMENTS

In this section, we empirically evaluate the proposed CRCR method.

4.1 EXPERIMENTAL SETTINGS

Datasets For comprehensive experimentation, we employ four popular benchmark datasets, including MNIST [12], Kuzushiji-MNIST (K-MNIST)[3], Fashion-MNIST (F-MNIST)[25] and CIFAR-10[11]. Additionally, experiments are conducted on two UCI datasets[2], including Optdigits and Pendigits. These datasets encompass more than just two labels, therefore, we categorize the class labels into positive and negative classes, effectively transforming them into binary classification datasets. Furthermore, for each dataset, we randomly selected $m\% \times n$ instances to add noise, where the noise ratio m is varied over $[0, 50, 75, 100]$. As a result, in our experiments, we generate 24 synthetic datasets in total.

Furthermore, we choose different models as backbones based on the varying feature dimensions of each dataset. Specifically, for MNIST, K-MNIST and F-MNIST, we use a 3-layer multilayer perceptron (MLP) with three hidden layers of width 300 equipped with the ReLU [17] activation function and batch normalization [8]. For CIFAR-10, we train a ResNet-34 model [6] as the backbone. For all UCI datasets, we use a linear model for training. The detailed information for each dataset is presented in Table 1.

Baseline methods We employ seven state-of-the-art algorithms for comparison, including four Pcomp methods (*i.e.*, PcompTeacher, PcompABS, PcompReLU and PcompUnbiased) and three ConfDiff methods (*i.e.*, ConfDiffABS, ConfDiffReLU and ConfDiffUnbiased). Details of baselines are described as follows:

- Pointwise Binary Classification with **Pairwise Confidence Comparisons (Pcomp)** [4]: A weakly supervised learning method that trains a binary classifier using pairwise comparison data, composed of unlabeled data pairs where one is more likely to be positive, instead of using pointwise data. Pcomp comprises four versions: PcompTeacher, PcompABS, PcompReLU, and PcompUnbiased. We use the code provided by its authors ¹.
- Binary Classification with **Confidence Difference (ConfDiff)** [22]: A weakly supervised learning method that trains a binary classifier using pairwise comparison data, which consists of pairwise unlabeled data where the difference in the probabilities of being positive (confidence difference) is known. ConfDiff comprises three versions: ConfDiff-ABS, ConfDiff-ReLU, and ConfDiff-Unbiased. We utilize the publicly available code online ².

Implementation details For each comparison method under every experimental configuration, we execute the code five times, employing the logistic loss function and Adam optimizer consistently. Specifically, during the training phase, each run is independently performed for 200 epochs with a batch size of 256. In balanced scenarios (*i.e.*, $\pi = 0.5$), the learning rate is set to 10^{-3} across all

¹<https://lfeng1995.github.io/codedata.html>

²<https://github.com/wwangwitsel/ConfDiff>

324 datasets, with weight decay parameters set to 10^{-5} for MNIST, K-MNIST, F-MNIST, and Pendigits,
 325 10^{-4} for Optdigits, and 10^{-3} for Pendigits. In imbalanced scenarios (*i.e.*, $\pi = 0.2$), the learning rate
 326 is set to 10^{-4} for MNIST and K-MNIST, and 10^{-3} for the remaining datasets, with weight decay
 327 parameters set to 10^{-4} for K-MNIST and Optdigits, and 10^{-5} for the remaining datasets. During the
 328 pretraining phase, each run is independently executed for 20 epochs with a batch size of 256. The
 329 learning rate and weight decay remain consistent with those in the training phase. All experiments
 330 are conducted on a server equipped with two Nvidia RTX 4090 GPUs.

332 4.2 CONSTRUCTION OF THE CONFIDENCE DIFFERENCES

333 In this subsection, we present the confidence differences construction method to address the challenge
 334 of fitting scenarios where precise posterior probabilities are difficult to obtain, along with a noise
 335 generation method to validate the robustness of our method under noisy conditions.

336 **The confidence differences construction method.** The ConfDiff method generates class posterior
 337 probabilities using a logistic regression-based probabilistic classifier trained on labeled data and
 338 calculates the confidence difference according to its definition. Although this generation method
 339 benefits comprehensive experimental analysis, it fails to accurately reflect the posterior probability
 340 distribution derived from manual annotations in real-world scenarios. Inspired by this, we incorporate
 341 an a posterior probability construction method based on outlier detection into the probabilistic
 342 classifier and computed confidence differences according to its definition to achieve a more uniform
 343 and realistic distribution. Specifically, we apply Gaussian kernel-based probability density estimation
 344 method to discrete posterior probabilities.

$$345 \hat{d}(x_i) = \frac{1}{nh\sqrt{2\pi}} \sum_{j=1}^n \exp\left(-\frac{(x_i - x_j)^2}{2h^2}\right), \quad (11)$$

346 where $\hat{d}(x_i)$ represents the estimated probability density function at instance x_i and $\exp\left(-\frac{(x_i - x_j)^2}{2h^2}\right)$
 347 is the standard Gaussian kernel function. Furthermore, h denotes the kernel bandwidth, which controls
 348 the degree of smoothing. This parameter is adaptively set based on the standard deviation of the
 349 probability distributions used in our work. We identify instances with densities below the threshold
 350 o as outliers. (Notably, o is also adaptively determined based on different probability density
 351 distributions. In our work, it is set at the 2nd percentile of the probability density to avoid filtering
 352 out too many instances.) The posterior probabilities of remaining non-outlier instances, are then
 353 rescaled to ensure a more uniform distribution within the range $[0, 1]$.

$$354 p(y_i = +1|\mathbf{x}_i) = \begin{cases} \text{Scaling}(p(y_i = +1|\mathbf{x}_i)), & \text{if } \hat{d}(x_i) \leq o \\ p(y_i = +1|\mathbf{x}_i), & \text{otherwise} \end{cases} \quad (12)$$

355 where $\text{Scaling}(\cdot)$ denotes a scaling function as:

$$356 \text{Scaling}(p(y_i = +1|\mathbf{x}_i)) = \begin{cases} \log(p(y_i = +1|\mathbf{x}_i) + \vartheta), & \text{if } p(y_i = +1|\mathbf{x}_i) \leq 0.5 \\ \log(1 - p(y_i = +1|\mathbf{x}_i) + \vartheta), & \text{otherwise} \end{cases} \quad (13)$$

357 where $\vartheta = e^{-10}$ is a smoothing parameter. Then, the confidence difference can be calculated
 358 according to its definition $c(\mathbf{x}_i, \mathbf{x}'_i) = p(y'_i = +1|\mathbf{x}'_i) - p(y_i = +1|\mathbf{x}_i)$.

359 **The noise generation method.** One straightforward method is to add noise directly to c . However,
 360 this method overlooks the intrinsic logic behind the original construction of c . We might be more
 361 interested in observing how the noise impacts the posterior probability distribution, thereby further
 362 influencing c indirectly. Then, we focus on adding noise to the posterior probabilities generated by
 363 the probabilistic classifier, thereby indirectly adding noise to c . In the real world, individuals tend to
 364 exhibit smaller judgment biases towards more similar sample pairs, while generating larger biases
 365 towards samples with lower similarity. Therefore, White Gaussian Noise (WGN) is introduced into
 366 the posterior probabilities $p(y_i = +1|\mathbf{x}_i)$ and $p(y'_i = +1|\mathbf{x}'_i)$ provided by the probabilistic classifier
 367 for the instance pair $(\mathbf{x}_i, \mathbf{x}'_i)$. Then, the noisy posterior probabilities are used to generate the label
 368 confidence difference, *i.e.*, $\tilde{c}_i = \tilde{c}(\mathbf{x}_i, \mathbf{x}'_i) = \tilde{p}(y'_i = +1|\mathbf{x}'_i) - \tilde{p}(y_i = +1|\mathbf{x}_i)$, where

$$369 \tilde{p}(y'_i = +1|\mathbf{x}'_i) = p(y'_i = +1|\mathbf{x}'_i) + \zeta'_i, \quad \zeta'_i \sim N(0, \sigma^2)$$

$$370 \tilde{p}(y_i = +1|\mathbf{x}_i) = p(y_i = +1|\mathbf{x}_i) + \zeta_i, \quad \zeta_i \sim N(0, \sigma^2), \quad (14)$$

371 where ζ'_i and ζ_i represent the noise offsets which follow a standard Gaussian distribution $N(0, \sigma^2)$.
 372 In our experiments, we set $\sigma = 1/3$.

Table 2: Classification accuracy of each comparing method on six datasets (mean±std) when $\pi = 0.5$, where the best performance is shown in boldface.

m	Method	MNIST	K-MNIST	F-MNIST	CIFAR-10	Pendigits	Opltdigits
0	PcompUnbiased	0.815±0.007	0.588±0.087	0.813±0.066	0.752±0.005	0.775±0.018	0.795±0.020
	PcompReLU	0.719±0.108	0.692±0.012	0.614±0.132	0.794±0.009	0.746±0.014	0.766±0.038
	PcompABS	0.830±0.005	0.727±0.015	0.837±0.010	0.828±0.006	0.645±0.059	0.722±0.027
	PcompTeacher	0.882±0.024	0.708±0.008	0.887±0.012	0.812±0.010	0.496±0.016	0.507±0.067
	ConfDiffUnbiased	0.723±0.072	0.576±0.029	0.771±0.085	0.848±0.014	0.675±0.071	0.799±0.023
	ConfDiffReLU	0.929±0.003	0.771±0.025	0.912±0.020	0.848±0.014	0.675±0.071	0.799±0.023
	ConfDiffABS	0.944±0.003	0.825±0.011	0.952±0.004	0.848±0.014	0.675±0.071	0.799±0.023
	CRCR_Unbiased	0.777±0.034	0.769±0.004	0.921±0.009	0.869±0.009	0.756±0.006	0.823±0.023
	CRCR_ReLU	0.919±0.019	0.685±0.080	0.925±0.017	0.869±0.009	0.753±0.007	0.823±0.023
CRCR_ABS	0.962±0.006	0.848±0.013	0.955±0.002	0.869±0.009	0.753±0.009	0.823±0.023	
50	PcompUnbiased	0.814±0.050	0.606±0.086	0.855±0.061	0.733±0.010	0.760±0.020	0.793±0.022
	PcompReLU	0.849±0.008	0.722±0.003	0.833±0.063	0.810±0.008	0.756±0.036	0.772±0.017
	PcompABS	0.853±0.016	0.730±0.013	0.876±0.015	0.833±0.005	0.676±0.069	0.736±0.017
	PcompTeacher	0.898±0.019	0.723±0.018	0.907±0.021	0.812±0.007	0.495±0.017	0.503±0.068
	ConfDiffUnbiased	0.678±0.046	0.602±0.021	0.794±0.034	0.833±0.013	0.675±0.073	0.792±0.021
	ConfDiffReLU	0.933±0.002	0.766±0.020	0.933±0.012	0.836±0.014	0.675±0.073	0.792±0.021
	ConfDiffABS	0.937±0.004	0.819±0.007	0.953±0.007	0.834±0.013	0.675±0.073	0.792±0.021
	CRCR_Unbiased	0.845±0.043	0.779±0.008	0.928±0.001	0.859±0.003	0.759±0.029	0.821±0.022
	CRCR_ReLU	0.923±0.023	0.793±0.019	0.936±0.007	0.860±0.003	0.757±0.030	0.821±0.022
CRCR_ABS	0.961±0.005	0.851±0.010	0.956±0.005	0.860±0.003	0.762±0.033	0.821±0.022	
75	PcompUnbiased	0.849±0.010	0.596±0.086	0.832±0.129	0.716±0.006	0.754±0.028	0.794±0.021
	PcompReLU	0.858±0.006	0.728±0.013	0.880±0.012	0.820±0.008	0.743±0.038	0.783±0.018
	PcompABS	0.865±0.008	0.734±0.017	0.874±0.011	0.836±0.003	0.688±0.060	0.743±0.020
	PcompTeacher	0.908±0.010	0.735±0.013	0.920±0.018	0.813±0.008	0.495±0.018	0.501±0.069
	ConfDiffUnbiased	0.620±0.084	0.560±0.025	0.650±0.051	0.844±0.008	0.674±0.073	0.795±0.018
	ConfDiffReLU	0.922±0.019	0.778±0.008	0.931±0.015	0.843±0.009	0.674±0.073	0.795±0.018
	ConfDiffABS	0.933±0.006	0.817±0.009	0.954±0.004	0.844±0.009	0.674±0.073	0.795±0.018
	CRCR_Unbiased	0.797±0.075	0.791±0.010	0.926±0.010	0.858±0.003	0.723±0.033	0.819±0.022
	CRCR_ReLU	0.938±0.006	0.792±0.010	0.942±0.005	0.858±0.003	0.721±0.035	0.819±0.022
CRCR_ABS	0.962±0.003	0.851±0.006	0.959±0.001	0.858±0.003	0.756±0.009	0.819±0.022	
100	PcompUnbiased	0.832±0.051	0.631±0.079	0.897±0.013	0.708±0.014	0.735±0.024	0.796±0.015
	PcompReLU	0.862±0.015	0.726±0.012	0.883±0.017	0.827±0.004	0.725±0.035	0.787±0.019
	PcompABS	0.865±0.014	0.735±0.009	0.886±0.009	0.837±0.006	0.688±0.059	0.766±0.020
	PcompTeacher	0.914±0.011	0.738±0.020	0.921±0.011	0.812±0.010	0.495±0.018	0.499±0.070
	ConfDiffUnbiased	0.631±0.056	0.548±0.022	0.573±0.060	0.835±0.012	0.669±0.070	0.791±0.021
	ConfDiffReLU	0.920±0.014	0.769±0.008	0.923±0.032	0.834±0.012	0.669±0.070	0.791±0.021
	ConfDiffABS	0.934±0.006	0.812±0.004	0.953±0.005	0.835±0.012	0.669±0.070	0.791±0.021
	CRCR_Unbiased	0.860±0.081	0.804±0.009	0.910±0.030	0.851±0.007	0.751±0.008	0.815±0.019
	CRCR_ReLU	0.939±0.006	0.797±0.006	0.941±0.006	0.851±0.007	0.752±0.008	0.815±0.019
CRCR_ABS	0.960±0.002	0.856±0.008	0.960±0.002	0.851±0.007	0.752±0.008	0.815±0.019	

4.3 RESULT ANALYSIS

Table 2 and Table 3 present the results of all baselines on four benchmark datasets and two UCI datasets for class-balanced (*i.e.*, prior = 0.5) and class-imbalanced scenarios (*i.e.*, prior = 0.2), respectively. Accuracy is chosen as the evaluation metric, and experiments are conducted five times on all datasets, with average and variance results recorded. Overall, our method performs nearly optimally across all scenarios compared to the baseline methods, consistently achieving nearly the best results using the ABS risk correction function.

In scenarios with balanced classes, our method outperforms Pcomp by improving accuracy from 0.02 to 0.341 and surpasses ConfDiff from 0.01 to 0.387, as observed from a baseline perspective. CRCR_ABS outperforms nearly all baselines, with the only observed exception being the results of PcompUnbiased on the Pendigits dataset when no noise is added. This may be due to the fact that the Pcomp method leverages only the information that one instance is more likely to be positive than another, without requiring knowledge of the exact difference between them. The posterior probability distribution is simply reconstructed in the absence of noise, and this reconstruction function preserves the monotonic increasing relationship of the posterior probabilities, without altering the relative likelihood of positivity between instances. Moreover, compared to Pcomp and ConfDiff, our method demonstrates increasingly stable and consistent accuracy as the noise ratio increases, with notable improvements in both accuracy and standard deviation, especially when the noise ratio reaches 100%. This indicates its ability to produce more competitive results in the presence of noise interference.

In scenarios with imbalanced classes, PcompReLU and ConfDiffReLU tend to exhibit random outcomes when confronted with imbalanced data augmented with noise. This phenomenon may be attributed to the introduced noise, which significantly increases the likelihood of predictions where

Table 3: Classification accuracy of each comparing method on six datasets (mean \pm std) when $\pi = 0.2$, where the best performance is shown in boldface.

m	Method	MNIST	K-MNIST	F-MNIST	CIFAR-10	Pendigits	Optdigits
0	PcompUnbiased	0.744 \pm 0.037	0.555 \pm 0.076	0.748 \pm 0.047	0.634 \pm 0.021	0.820 \pm 0.025	0.813 \pm 0.024
	PcompReLU	0.800 \pm 0.000	0.800 \pm 0.000	0.800 \pm 0.000	0.802 \pm 0.003	0.819 \pm 0.020	0.816 \pm 0.007
	PcompABS	0.804 \pm 0.009	0.800 \pm 0.000	0.801 \pm 0.001	0.833 \pm 0.004	0.797 \pm 0.023	0.805 \pm 0.006
	PcompTeacher	0.788 \pm 0.074	0.695 \pm 0.046	0.883 \pm 0.026	0.813 \pm 0.020	0.482 \pm 0.212	0.684 \pm 0.097
	ConfDiffUnbiased	0.743 \pm 0.033	0.622 \pm 0.077	0.724 \pm 0.025	0.812 \pm 0.004	0.797 \pm 0.028	0.830 \pm 0.016
	ConfDiffReLU	0.800 \pm 0.000	0.800 \pm 0.000	0.846 \pm 0.064	0.800 \pm 0.000	0.797 \pm 0.028	0.830 \pm 0.016
	ConfDiffABS	0.910 \pm 0.015	0.841 \pm 0.014	0.940\pm0.010	0.800 \pm 0.001	0.797 \pm 0.028	0.830 \pm 0.016
	CRCR_Unbiased	0.816 \pm 0.043	0.597 \pm 0.055	0.886 \pm 0.009	0.841\pm0.012	0.823\pm0.005	0.838\pm0.017
	CRCR_ReLU	0.929 \pm 0.055	0.814 \pm 0.031	0.930 \pm 0.049	0.801 \pm 0.001	0.817 \pm 0.012	0.830 \pm 0.008
	CRCR_ABS	0.916\pm0.022	0.856\pm0.006	0.922 \pm 0.007	0.812 \pm 0.017	0.784 \pm 0.024	0.825 \pm 0.005
50	PcompUnbiased	0.742 \pm 0.015	0.547 \pm 0.038	0.768 \pm 0.070	0.623 \pm 0.017	0.818 \pm 0.025	0.810 \pm 0.027
	PcompReLU	0.800 \pm 0.000	0.801 \pm 0.002	0.800 \pm 0.000	0.801 \pm 0.003	0.806 \pm 0.023	0.821 \pm 0.007
	PcompABS	0.824 \pm 0.029	0.800 \pm 0.000	0.809 \pm 0.006	0.833 \pm 0.006	0.801 \pm 0.030	0.811 \pm 0.010
	PcompTeacher	0.822 \pm 0.061	0.707 \pm 0.062	0.902 \pm 0.014	0.797 \pm 0.033	0.483 \pm 0.211	0.682 \pm 0.096
	ConfDiffUnbiased	0.694 \pm 0.030	0.640 \pm 0.043	0.711 \pm 0.018	0.805 \pm 0.006	0.797 \pm 0.029	0.834 \pm 0.015
	ConfDiffReLU	0.800 \pm 0.000	0.800 \pm 0.000	0.821 \pm 0.046	0.800 \pm 0.001	0.797 \pm 0.029	0.834 \pm 0.015
	ConfDiffABS	0.891 \pm 0.025	0.818 \pm 0.010	0.938 \pm 0.014	0.801 \pm 0.002	0.797 \pm 0.029	0.834 \pm 0.015
	CRCR_Unbiased	0.794 \pm 0.043	0.623 \pm 0.079	0.880 \pm 0.016	0.789 \pm 0.035	0.795 \pm 0.025	0.843\pm0.023
	CRCR_ReLU	0.908 \pm 0.063	0.815 \pm 0.015	0.936 \pm 0.043	0.811 \pm 0.025	0.808 \pm 0.021	0.838 \pm 0.014
	CRCR_ABS	0.916\pm0.013	0.830\pm0.029	0.950\pm0.011	0.850\pm0.017	0.822\pm0.019	0.835 \pm 0.011
75	PcompUnbiased	0.753 \pm 0.031	0.535 \pm 0.048	0.775 \pm 0.059	0.616 \pm 0.038	0.817 \pm 0.017	0.813 \pm 0.030
	PcompReLU	0.804 \pm 0.009	0.804 \pm 0.007	0.800 \pm 0.000	0.805 \pm 0.012	0.822 \pm 0.020	0.827 \pm 0.008
	PcompABS	0.863 \pm 0.014	0.800 \pm 0.000	0.828 \pm 0.016	0.832 \pm 0.005	0.803 \pm 0.038	0.813 \pm 0.009
	PcompTeacher	0.840 \pm 0.061	0.714 \pm 0.055	0.908 \pm 0.019	0.793 \pm 0.044	0.482 \pm 0.211	0.680 \pm 0.096
	ConfDiffUnbiased	0.704 \pm 0.058	0.630 \pm 0.026	0.745 \pm 0.088	0.804 \pm 0.004	0.796 \pm 0.031	0.830 \pm 0.016
	ConfDiffReLU	0.800 \pm 0.000	0.800 \pm 0.000	0.800 \pm 0.000	0.800 \pm 0.000	0.796 \pm 0.031	0.830 \pm 0.016
	ConfDiffABS	0.862 \pm 0.030	0.806 \pm 0.005	0.921 \pm 0.023	0.800 \pm 0.001	0.796 \pm 0.031	0.830 \pm 0.016
	CRCR_Unbiased	0.792 \pm 0.047	0.640 \pm 0.028	0.861 \pm 0.033	0.772 \pm 0.020	0.811 \pm 0.030	0.836 \pm 0.022
	CRCR_ReLU	0.901 \pm 0.055	0.817 \pm 0.024	0.801 \pm 0.001	0.828 \pm 0.024	0.827\pm0.008	0.836 \pm 0.017
	CRCR_ABS	0.914\pm0.008	0.819\pm0.022	0.947\pm0.006	0.853\pm0.004	0.819 \pm 0.011	0.839\pm0.012
100	PcompUnbiased	0.752 \pm 0.021	0.540 \pm 0.069	0.834 \pm 0.034	0.643 \pm 0.053	0.805 \pm 0.024	0.817 \pm 0.027
	PcompReLU	0.845 \pm 0.040	0.808 \pm 0.010	0.814 \pm 0.019	0.806 \pm 0.005	0.808 \pm 0.020	0.834 \pm 0.008
	PcompABS	0.871 \pm 0.006	0.801 \pm 0.001	0.844 \pm 0.015	0.835 \pm 0.003	0.803 \pm 0.029	0.823 \pm 0.012
	PcompTeacher	0.869 \pm 0.068	0.711 \pm 0.062	0.922 \pm 0.011	0.787 \pm 0.033	0.482 \pm 0.211	0.681 \pm 0.096
	ConfDiffUnbiased	0.772 \pm 0.056	0.693 \pm 0.028	0.748 \pm 0.101	0.810 \pm 0.007	0.796 \pm 0.028	0.831 \pm 0.017
	ConfDiffReLU	0.800 \pm 0.000	0.800 \pm 0.000	0.800 \pm 0.000	0.800 \pm 0.001	0.796 \pm 0.028	0.831 \pm 0.016
	ConfDiffABS	0.814 \pm 0.006	0.801 \pm 0.002	0.870 \pm 0.043	0.801 \pm 0.001	0.796 \pm 0.028	0.831 \pm 0.016
	CRCR_Unbiased	0.790 \pm 0.036	0.639 \pm 0.059	0.838 \pm 0.056	0.780 \pm 0.014	0.797 \pm 0.018	0.837 \pm 0.026
	CRCR_ReLU	0.905 \pm 0.059	0.808 \pm 0.007	0.903 \pm 0.039	0.800 \pm 0.001	0.800 \pm 0.014	0.838 \pm 0.020
	CRCR_ABS	0.926\pm0.010	0.828\pm0.025	0.958\pm0.009	0.841\pm0.005	0.810\pm0.006	0.839\pm0.019

one instance in a pair is incorrectly predicted to be more likely positive than the other, contrary to the actual scenario. This contradiction becomes significantly more pronounced as class imbalance and noise ratio increase. For other baselines, we observe advantages in both accuracy mean and variance. From the dataset perspective, CRCR_ABS significantly outperforms other methods on the MNIST, K-MNIST, F-MNIST, and CIFAR-10 datasets in the presence of noise, while maintaining strong competitiveness on the Pendigits and Optdigits datasets. CRCR_Unbiased shows promising results without noise; however, the experiments clearly demonstrate that its training challenges on complex and noisy datasets often lead to a notable decline in performance. This further underscores the effectiveness of CRCR_ABS in maintaining robust performance when dealing with complex datasets.

4.4 PARAMETER SENSITIVITY

In this subsection, we conduct experiments with different thresholds θ for partitioning subsets and the parameter α for the consistency term, and the results are shown in Figure 2.

About different threshold θ To evaluate the sensitivity of the threshold θ , we vary its value within the range $\{0.1, 0.2, \dots, 1\}$ and examine its influence on four distinct benchmark datasets (*i.e.*, MNIST, K-MNIST, F-MNIST and CIFAR-10). The results reveal that the accuracy score peaks for the four benchmark datasets when $\theta = 0.4$ with $\pi = 0.5$, and when $\theta = 0.2$ or 0.3 with $\pi = 0.2$. This observation may be attributed to the distribution of confidence differences resembling a waveform akin to a normal distribution. A low threshold results in numerous inaccurate predictions within the subset D^S utilized for risk consistency, while a high threshold leads to a scarcity of samples within

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

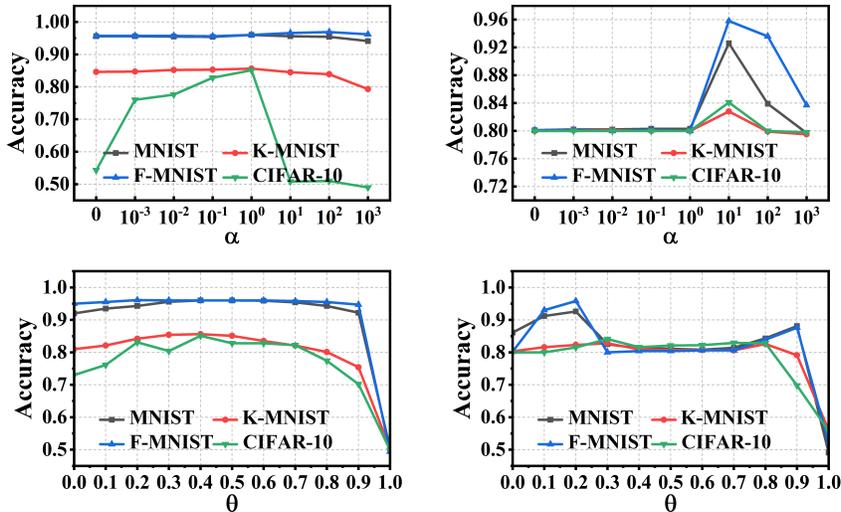


Figure 2: Sensitivity analysis of parameters α (top) and θ (bottom) on four benchmark datasets when $\pi = 0.5$ (left) and $\pi = 0.2$ (right).

D^S , thus diminishing the available supervisory information. Therefore, we empirically recommend setting the threshold at $\theta = 0.4$ when $\pi = 0.5$, and $\theta = 0.2$ or 0.3 when $\pi = 0.2$.

About different parameter α To assess the sensitivity of the parameter α , we vary its values across the range $\{10^i | i = -3, \dots, +3\}$ and observe its effects on four benchmark datasets. Our analysis reveals that α shows increased sensitivity on the larger-scale CIFAR-10 dataset when $\pi = 0.5$, while maintaining relatively stable performance on the smaller-scale datasets. Moreover, α leads to a consistent trend in accuracy variation across the four datasets when $\pi = 0.2$. Notably, it achieves relatively optimal results when $\alpha = 1$ with $\pi = 0.5$, and $\alpha = 10^1$ with $\pi = 0.2$. Thus, we recommend setting $\alpha = 1$ or 10^1 in experimental setups.

4.5 ABLATION STUDY

In this subsection, we conduct ablation studies on various strategies by setting corresponding parameters to zero. Specifically, setting $\{\alpha = 0, \theta = 0\}$ represent versions without consistency strategy and without subset segmentation strategy, respectively. The experimental results, presented in Figure 1, demonstrate that our proposed subset segmentation strategy and consistency term contribute to performance improvement to some extent in the context of noisy confidence difference classification.

5 CONCLUSION

In this paper, we propose a novel ConfDiff classification method based on consistency risk and consistency regularization to address the challenge of noisy supervised signals in ConfDiff classification. We conduct a theoretical analysis of various supervised signals associated with different confidence differences. Based on this analysis, the ConfDiff dataset is partitioned into two subsets according to the reliability of the supervised information. For the subset with more reliable supervision, we employ a consistency risk to preserve precise supervised information. Conversely, for the subset with less reliable supervision, we leverage consistency regularization to mitigate the impact of erroneous predictions. Extensive experimental results demonstrate that the proposed CRCR method outperforms state-of-the-art baselines and exhibits strong robustness, even under artificially induced noise.

REFERENCES

[1] Han Bao, Gang Niu, and Masashi Sugiyama. Classification from pairwise similarity and unlabeled data. In *International Conference on Machine Learning*, pp. 452–461. PMLR, 2018.

- 540 [2] Catherine L Blake. Uci repository of machine learning databases. [http://www.ics.uci.edu/~](http://www.ics.uci.edu/~mlearn/MLRepository.html)
541 [mlearn/MLRepository.html](http://www.ics.uci.edu/~mlearn/MLRepository.html), 1998.
- 542 [3] Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and
543 David Ha. Deep learning for classical japanese literature. *arXiv preprint arXiv:1812.01718*,
544 2018.
- 545 [4] Lei Feng, Senlin Shu, Nan Lu, Bo Han, Miao Xu, Gang Niu, Bo An, and Masashi Sugiyama.
546 Pointwise binary classification with pairwise confidence comparisons. In *International Confer-*
547 *ence on Machine Learning*, pp. 3252–3262. PMLR, 2021.
- 548 [5] Zayd Hammoudeh and Daniel Lowd. Learning from positive and unlabeled data with arbitrary
549 positive shift. In *Neural Information Processing Systems*, 2020.
- 550 [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
551 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
552 pp. 770–778, 2016.
- 553 [7] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi-
554 narayanan. Augmix: A simple data processing method to improve robustness and uncertainty.
555 *arXiv preprint arXiv:1912.02781*, 2019.
- 556 [8] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by
557 reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456.
558 pmlr, 2015.
- 559 [9] Jang-Hyun Kim, Wonho Choo, and Hyun Oh Song. Puzzle mix: Exploiting saliency and local
560 statistics for optimal mixup. In *International Conference on Machine Learning*, pp. 5275–5285.
561 PMLR, 2020.
- 562 [10] Ryuichi Kiryo, Gang Niu, Marthinus Christoffel du Plessis, and Masashi Sugiyama. Positive-
563 unlabeled learning with non-negative risk estimator. In *Neural Information Processing Systems*,
564 pp. 1675–1685, 2017.
- 565 [11] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
566 2009. URL <https://api.semanticscholar.org/CorpusID:18268744>.
- 567 [12] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning
568 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 569 [13] Changchun Li, Ximing Li, Lei Feng, and Jihong Ouyang. Who is your right mixup partner
570 in positive and unlabeled learning. In *International Conference on Learning Representations*,
571 2021.
- 572 [14] Nan Lu, Gang Niu, Aditya Krishna Menon, and Masashi Sugiyama. On the minimal supervision
573 for training any binary classifier from only unlabeled data. *arXiv preprint arXiv:1808.10585*,
574 2018.
- 575 [15] Nan Lu, Tianyi Zhang, Gang Niu, and Masashi Sugiyama. Mitigating overfitting in super-
576 vised classification from two unlabeled datasets: A consistent risk correction approach. In
577 *International Conference on Artificial Intelligence and Statistics*, pp. 1115–1125. PMLR, 2020.
- 578 [16] Chuan Luo, Pu Zhao, Chen Chen, Bo Qiao, Chao Du, Hongyu Zhang, Wei Wu, Shaowei Cai,
579 Bing He, Saravanakumar Rajmohan, and Qingwei Lin. PULNS: positive-unlabeled learning
580 with effective negative sample selector. In *AAAI Conference on Artificial Intelligence*, pp.
581 8784–8792, 2021.
- 582 [17] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines.
583 In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–
584 814, 2010.
- 585 [18] Quang Nguyen, Hamed Valizadegan, and Milos Hauskrecht. Learning classification models
586 with soft-label information. *Journal of the American Medical Informatics Association*, 21(3):
587 501–508, 2014.

- 594 [19] Xiaoshuang Shi, Fuyong Xing, Yuanpu Xie, Zizhao Zhang, Lei Cui, and Lin Yang. Loss-based
595 attention for deep multiple instance learning. In *Proceedings of the AAAI conference on artificial*
596 *intelligence*, volume 34, pp. 5742–5749, 2020.
- 597 [20] Takuya Shimada, Han Bao, Issei Sato, and Masashi Sugiyama. Classification from pairwise sim-
598 ilarities/dissimilarities and unlabeled data via empirical risk minimization. *Neural Computation*,
599 33(5):1234–1268, 2021.
- 600 [21] Vikas Verma, Alex Lamb, Christopher Beckham, Amir Najafi, Ioannis Mitliagkas, David Lopez-
601 Paz, and Yoshua Bengio. Manifold mixup: Better representations by interpolating hidden states.
602 In *International conference on machine learning*, pp. 6438–6447. PMLR, 2019.
- 603 [22] Wei Wang, Lei Feng, Yuchen Jiang, Gang Niu, Min-Ling Zhang, and Masashi Sugiyama. Binary
604 classification with confidence difference. *Advances in Neural Information Processing Systems*,
605 36, 2024.
- 606 [23] Xinrui Wang, Wenhai Wan, Chuanxing Geng, Shaoyuan Li, and Songcan Chen. Beyond
607 myopia: Learning from positive and unlabeled data through holistic predictive trends. In *Neural*
608 *Information Processing Systems*, 2023.
- 609 [24] Jia Wu, Shirui Pan, Xingquan Zhu, Chengqi Zhang, and Xindong Wu. Multi-instance learning
610 with discriminative bag mapping. *IEEE Transactions on Knowledge and Data Engineering*, 30
611 (6):1065–1080, 2018.
- 612 [25] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
613 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 614 [26] Yanbing Xue and Milos Hauskrecht. Learning of classification models from noisy soft-labels.
615 In *Proceedings of the Twenty-second European Conference on Artificial Intelligence*, pp. 1618–
616 1619, 2016.
- 617 [27] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon
618 Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In
619 *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6023–6032,
620 2019.
- 621 [28] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond
622 empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- 623 [29] Linjun Zhang, Zhun Deng, Kenji Kawaguchi, Amirata Ghorbani, and James Zou. How does
624 mixup help with robustness and generalization? In *International Conference on Learning*
625 *Representations*, 2021.
- 626 [30] Min-Ling Zhang and Zhi-Hua Zhou. Multi-instance clustering with applications to multi-
627 instance prediction. *Applied intelligence*, 31:47–68, 2009.
- 628 [31] Yunrui Zhao, Qianqian Xu, Yangbangyan Jiang, Peisong Wen, and Qingming Huang. Dist-pu:
629 Positive-unlabeled learning from a label distribution perspective. In *IEEE/CVF Conference on*
630 *Computer Vision and Pattern Recognition*, pp. 14441–14450, 2022.
- 631 [32] Zhi-Hua Zhou, Yu-Yin Sun, and Yu-Feng Li. Multi-instance learning by treating instances
632 as non-iid samples. In *Proceedings of the 26th annual international conference on machine*
633 *learning*, pp. 1249–1256, 2009.
- 634
635
636
637
638
639
640
641
642
643
644
645
646
647

648 A PROOF OF THEOREM 1

649 In this appendix, we provide the proof of the Theorem 1 and the corresponding technical lemmas.

650 **Lemma 1.** *The Rademacher complexity $\bar{\mathfrak{R}}_n(\mathcal{L}_{\text{CRCR}} \circ \mathcal{G})$ on \mathcal{D} for ConfDiff data with noise of size n*
 651 *can be defined as follows:*

$$652 \bar{\mathfrak{R}}_n(\mathcal{L}_{\text{CRCR}} \circ \mathcal{G}) \leq 2L_\ell \mathfrak{R}_{n_1}(\mathcal{G}) + \frac{\alpha}{\log(\varepsilon)} \mathfrak{R}_{n_2}(\mathcal{G}) \quad (15)$$

653 The proof of Lemma 1:

$$654 \begin{aligned} 655 \bar{\mathfrak{R}}_n(\mathcal{L}_{\text{CRCR}} \circ \mathcal{G}) &= \mathbb{E}_{\mathcal{D}_{n_1}} \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \frac{1}{n_1} \sum_{i=1}^{n_1} \sigma_i \mathcal{L}_{\text{CRCR}}^S(g; \mathbf{x}_i, \mathbf{x}'_i) \right] \\ 656 &+ \mathbb{E}_{\mathcal{D}_{n_2}} \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \frac{1}{n_2} \sum_{i=1}^{n_2} \sigma_i \mathcal{L}_{\text{CRCR}}^C(g; \mathbf{x}_i, \mathbf{x}'_i) \right] \\ 657 &= \mathbb{E}_{\mathcal{D}_{n_1}} \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \frac{1}{n_1} \sum_{i=1}^{n_1} \frac{1}{2} \sigma_i \left((\pi - c_i) \ell(g(\mathbf{x}_i), +1) + (1 - \pi - c_i) \ell(g(\mathbf{x}'_i), -1) \right. \right. \\ 658 &\quad \left. \left. + (\pi + c_i) \ell(g(\mathbf{x}'_i), +1) + (1 - \pi + c_i) \ell(g(\mathbf{x}_i), -1) \right) \right] \\ 659 &+ \mathbb{E}_{\mathcal{D}_{n_2}} \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \frac{1}{n_2} \sum_{i=1}^{n_2} \alpha \sigma_i \frac{1}{\log(|\tilde{c}_i| + \varepsilon)} \cdot \|(g(\mathbf{x}_i) - g(\mathbf{x}'_i))\|_2 \right] \\ 660 &= \mathbb{E}_{\mathcal{D}_{n_1}} \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \frac{1}{n_1} \sum_{i=1}^{n_1} \sigma_i \|\nabla \mathcal{L}_{\text{CD}}^S(g; \mathbf{x}_i, \mathbf{x}'_i)\|_2 g(\mathbf{x}_i) \right] \\ 661 &+ \mathbb{E}_{\mathcal{D}_{n_2}} \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \frac{1}{n_2} \sum_{i=1}^{n_2} \sigma_i \|\nabla \mathcal{L}_{\text{CD}}^C(g; \mathbf{x}_i, \mathbf{x}'_i)\|_2 g(\mathbf{x}_i) \right] \end{aligned} \quad (16)$$

662 where

$$663 \begin{aligned} 664 &\|\nabla \mathcal{L}_{\text{CRCR}}^S(g; \mathbf{x}_i, \mathbf{x}'_i)\|_2 \\ 665 &= \frac{1}{2} \left\| \nabla \left((\pi - c_i) \ell(g(\mathbf{x}_i), +1) + (1 - \pi - c_i) \ell(g(\mathbf{x}'_i), -1) \right. \right. \\ 666 &\quad \left. \left. + (\pi + c_i) \ell(g(\mathbf{x}'_i), +1) + (1 - \pi + c_i) \ell(g(\mathbf{x}_i), -1) \right) \right\|_2 \\ 667 &\leq \frac{1}{2} \left(\|\nabla((\pi - c_i) \ell(g(\mathbf{x}_i), +1))\|_2 + \|\nabla((1 - \pi - c_i) \ell(g(\mathbf{x}'_i), -1))\|_2 \right. \\ 668 &\quad \left. + \|\nabla((\pi + c_i) \ell(g(\mathbf{x}'_i), +1))\|_2 + \|\nabla((1 - \pi + c_i) \ell(g(\mathbf{x}_i), -1))\|_2 \right) \\ 669 &\leq \frac{1}{2} |\pi - c_i| L_\ell + \frac{1}{2} |1 - \pi - c_i| L_\ell + \frac{1}{2} |\pi + c_i| L_\ell + \frac{1}{2} |1 - \pi + c_i| L_\ell \\ 670 &\leq 2L_\ell \end{aligned} \quad (17)$$

671 and,

$$672 \begin{aligned} 673 \|\nabla \mathcal{L}_{\text{CRCR}}^C(g; \mathbf{x}_i, \mathbf{x}'_i)\|_2 &= \alpha \left\| \nabla \frac{1}{\log(|\tilde{c}_i| + \varepsilon)} \cdot \|(g(\mathbf{x}_i) - g(\mathbf{x}'_i))\|_2 \right\|_2 \\ 674 &\leq \alpha \frac{1}{\log(|\tilde{c}_i| + \varepsilon)} \cdot \frac{g(\mathbf{x}_i) - g(\mathbf{x}'_i)}{\|(g(\mathbf{x}_i) - g(\mathbf{x}'_i))\|_2} \\ 675 &\leq \frac{\alpha}{\log(\varepsilon)} \end{aligned} \quad (18)$$

676 Replacing the corresponding term in Eq.16 with Eq.18 and Eq.19, we can prove the Lemma 1:

$$677 \begin{aligned} 678 \bar{\mathfrak{R}}_n(\mathcal{L}_{\text{CRCR}} \circ \mathcal{G}) &\leq 2L_\ell \mathbb{E}_{\mathcal{D}_{n_1}} \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \frac{1}{n_1} \sum_{i=1}^{n_1} \sigma_i g(\mathbf{x}_i) \right] + \frac{\alpha}{\log(\varepsilon)} \mathbb{E}_{\mathcal{D}_{n_2}} \mathbb{E}_\sigma \left[\sup_{g \in \mathcal{G}} \frac{1}{n_2} \sum_{i=1}^{n_2} \sigma_i g(\mathbf{x}_i) \right] \\ 679 &\leq 2L_\ell \mathfrak{R}_{n_1}(\mathcal{G}) + \frac{\alpha}{\log(\varepsilon)} \mathfrak{R}_{n_2}(\mathcal{G}) \end{aligned} \quad (19)$$

Lemma 2.

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left| R(g) - \hat{R}_{\text{CRCR}}(g) \right| &\leq 4L_\ell \mathfrak{R}_{n_1}(\mathcal{G}) + \frac{2\alpha}{\log(\varepsilon)} \mathfrak{R}_{n_2}(\mathcal{G}) \\ &\quad + \left(\frac{C_\ell}{n_1} + \left| \frac{1}{\log(\varepsilon)} - \frac{1}{\log(\theta + \varepsilon)} \right| \frac{4\alpha C_g^2}{n_2} \right) \sqrt{2n \ln(2/\delta)} \end{aligned} \quad (21)$$

The proof of Lemma 2: Let $\hat{R}_{\text{CRCR}}(g)$ and $\hat{\hat{R}}_{\text{CRCR}}(g)$ represent the empirical risks of two sets of training samples, each differing by exactly one point, denoted as $\{(\mathbf{x}_i, \mathbf{x}'_i), c(\mathbf{x}_i, \mathbf{x}'_i)\}$ and $\{(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i), c(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}'_i)\}$ respectively.

$$\begin{aligned} &\sup_{g \in \mathcal{G}} \left| (R(g) - \hat{R}_{\text{CRCR}}(g)) - (R(g) - \hat{\hat{R}}_{\text{CRCR}}(g)) \right| \\ &\leq \sup_{g \in \mathcal{G}} \left| \hat{R}_{\text{CRCR}}(g) - \hat{\hat{R}}_{\text{CRCR}}(g) \right| \\ &\leq \sup_{g \in \mathcal{G}} \left| \frac{1}{2n_1} (\pi - \tilde{c}_i) (\ell(g(\mathbf{x}_i), +1) - \ell(g(\bar{\mathbf{x}}_i), +1)) \right. \end{aligned} \quad (22)$$

$$\left. + (1 - \pi - \tilde{c}_i) (\ell(g(\mathbf{x}'_i), -1) - \ell(g(\bar{\mathbf{x}}'_i), -1)) \right) \quad (23)$$

$$\left. + (\pi + \tilde{c}_i) (\ell(g(\mathbf{x}'_i), +1) - \ell(g(\bar{\mathbf{x}}'_i), +1)) \right) \quad (24)$$

$$\left. + (1 - \pi + \tilde{c}_i) (\ell(g(\mathbf{x}_i), -1) - \ell(g(\bar{\mathbf{x}}_i), -1)) \right) \quad (25)$$

$$\begin{aligned} &+ \frac{\alpha}{n_2} \left(\frac{1}{\log(|\tilde{c}_i| + \varepsilon)} \cdot \|g(\mathbf{x}_i) - g(\mathbf{x}'_i)\|_2 - \frac{1}{\log(|\tilde{c}_i| + \varepsilon)} \cdot \|g(\bar{\mathbf{x}}_i) - g(\bar{\mathbf{x}}'_i)\|_2 \right) \Big| \\ &\leq \frac{2C_\ell}{n_1} + \left| \frac{1}{\log(\varepsilon)} - \frac{1}{\log(\theta + \varepsilon)} \right| \frac{2\alpha C_g}{n_2} \end{aligned} \quad (26)$$

Then according McDiarmid's inequality:

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left| R(g) - \hat{R}_{\text{CRCR}}(g) \right| &\leq \mathbb{E}_{\mathcal{D}_n} [\sup_{g \in \mathcal{G}} (R(g) - \hat{R}_{\text{CRCR}}(g))] \\ &\quad + \left(\frac{2C_\ell}{n_1} + \left| \frac{1}{\log(\varepsilon)} - \frac{1}{\log(\theta + \varepsilon)} \right| \frac{2\alpha C_g}{n_2} \right) \sqrt{2n \ln(2/\delta)} \\ &\leq 2\bar{\mathfrak{R}}_n(\mathcal{L}_{\text{CRCR}} \circ \mathcal{G}) \\ &\quad + \left(\frac{2C_\ell}{n_1} + \left| \frac{1}{\log(\varepsilon)} - \frac{1}{\log(\theta + \varepsilon)} \right| \frac{2\alpha C_g}{n_2} \right) \sqrt{2n \ln(2/\delta)} \\ &\leq 4L_\ell \mathfrak{R}_{n_1}(\mathcal{G}) + \frac{2\alpha}{\log(\varepsilon)} \mathfrak{R}_{n_2}(\mathcal{G}) \\ &\quad + \left(\frac{2C_\ell}{n_1} + \left| \frac{1}{\log(\varepsilon)} - \frac{1}{\log(\theta + \varepsilon)} \right| \frac{2\alpha C_g}{n_2} \right) \sqrt{2n \ln(2/\delta)} \end{aligned} \quad (27)$$

The proof of Theorem 1:

$$\begin{aligned} R(\hat{g}_{\text{CRCR}}) - R(g^*) &= (R(\hat{g}_{\text{CRCR}}) - \hat{R}_{\text{CRCR}}(\hat{g}_{\text{CRCR}})) + (\hat{R}_{\text{CRCR}}(\hat{g}_{\text{CRCR}}) - \hat{R}_{\text{CRCR}}(g^*)) \\ &\quad + (\hat{R}_{\text{CRCR}}(g^*) - R(g^*)) \\ &\leq (R(\hat{g}_{\text{CRCR}}) - \hat{R}_{\text{CRCR}}(\hat{g}_{\text{CRCR}})) + (\hat{R}_{\text{CRCR}}(g^*) - R(g^*)) \\ &\leq 2 \sup_{g \in \mathcal{G}} \left| R(g) - \hat{R}_{\text{CRCR}}(g) \right| \\ &\leq 8L_\ell \mathfrak{R}_{n_1}(\mathcal{G}) + \frac{4\alpha}{\log(\varepsilon)} \mathfrak{R}_{n_2}(\mathcal{G}) \\ &\quad + \left(\frac{4C_\ell}{n_1} + \left| \frac{1}{\log(\varepsilon)} - \frac{1}{\log(\theta + \varepsilon)} \right| \frac{4\alpha C_g}{n_2} \right) \sqrt{2n \ln(2/\delta)} \end{aligned} \quad (28)$$

756 **B PROOF OF EQ. 6**
 757
 758

759 In this appendix, we provide the proof of the Eq. 6.

760 Substituting the form of the loss function from Eq.5 into Eq.3, then we can obtain:
 761

$$\begin{aligned}
 762 R_{CD}(g) &= \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[(\pi - c(\mathbf{x}, \mathbf{x}')) \ell(g(\mathbf{x}), +1) + (1 - \pi - c(\mathbf{x}, \mathbf{x}')) \ell(g(\mathbf{x}'), -1) \right. \\
 763 &\quad \left. + (\pi + c(\mathbf{x}, \mathbf{x}')) \ell(g(\mathbf{x}'), +1) + (1 - \pi + c(\mathbf{x}, \mathbf{x}')) \ell(g(\mathbf{x}), -1) \right] \\
 764 &= \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[(\pi - c(\mathbf{x}, \mathbf{x}')) (h(g(\mathbf{x})) - g(\mathbf{x})) + (1 - \pi - c(\mathbf{x}, \mathbf{x}')) (h(g(\mathbf{x}')) + g(\mathbf{x}')) \right. \\
 765 &\quad \left. + (\pi + c(\mathbf{x}, \mathbf{x}')) (h(g(\mathbf{x}')) - g(\mathbf{x}')) + (1 - \pi + c(\mathbf{x}, \mathbf{x}')) (h(g(\mathbf{x})) + g(\mathbf{x})) \right] \\
 766 &= \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[h(g(\mathbf{x})) + (1 - 2\pi + 2c(\mathbf{x}, \mathbf{x}')) g(\mathbf{x}) \right. \\
 767 &\quad \left. + h(g(\mathbf{x}')) + (1 - 2\pi - 2c(\mathbf{x}, \mathbf{x}')) g(\mathbf{x}') \right] \\
 768 &= \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[h(g(\mathbf{x})) + 2c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}) + h(g(\mathbf{x}')) - 2c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') \right] \\
 769 &\quad + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[(1 - 2\pi) (g(\mathbf{x}) + g(\mathbf{x}')) \right] \\
 770 &= \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[h(g(\mathbf{x})) + 2c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}) + h(g(\mathbf{x}')) - 2c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') \right. \\
 771 &\quad \left. + c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x})) - c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x})) + \frac{1}{2} g(\mathbf{x}) - \frac{1}{2} g(\mathbf{x}) \right. \\
 772 &\quad \left. + c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x}')) - c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x}')) + \frac{1}{2} g(\mathbf{x}') - \frac{1}{2} g(\mathbf{x}') \right] \\
 773 &\quad + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[(1 - 2\pi) (g(\mathbf{x}) + g(\mathbf{x}')) \right] \\
 774 &= \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[\frac{1}{2} h(g(\mathbf{x})) - c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x})) - \frac{1}{2} g(\mathbf{x}) + c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}) \right. \\
 775 &\quad \left. + \frac{1}{2} h(g(\mathbf{x}')) + c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x}')) - \frac{1}{2} g(\mathbf{x}') - c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') \right] \\
 776 &\quad + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[(1 - 2\pi) (g(\mathbf{x}) + g(\mathbf{x}')) \right] \\
 777 &= \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[\frac{1}{2} h(g(\mathbf{x})) - c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x})) - \frac{1}{2} g(\mathbf{x}) + c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}) \right. \\
 778 &\quad \left. + \frac{1}{2} h(g(\mathbf{x}')) + c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x}')) - \frac{1}{2} g(\mathbf{x}') - c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') \right] \\
 779 &\quad + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[(1 - 2\pi) (g(\mathbf{x}) + g(\mathbf{x}')) \right] \\
 780 &= \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[\frac{1}{2} h(g(\mathbf{x})) - c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x})) - \frac{1}{2} g(\mathbf{x}) + c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}) \right. \\
 781 &\quad \left. + \frac{1}{2} h(g(\mathbf{x}')) + c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x}')) - \frac{1}{2} g(\mathbf{x}') - c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') \right] \\
 782 &\quad + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[\frac{1}{2} h(g(\mathbf{x})) + c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x})) + \frac{1}{2} g(\mathbf{x}) + c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}) \right. \\
 783 &\quad \left. + \frac{1}{2} h(g(\mathbf{x}')) - c(\mathbf{x}, \mathbf{x}') h(g(\mathbf{x}')) + \frac{1}{2} g(\mathbf{x}') - c(\mathbf{x}, \mathbf{x}') g(\mathbf{x}') \right] \\
 784 &\quad + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[(1 - 2\pi) (g(\mathbf{x}) + g(\mathbf{x}')) \right] \\
 785 &= \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[\left(\frac{1}{2} - c(\mathbf{x}, \mathbf{x}') \right) \ell(g(\mathbf{x}), +1) + \left(\frac{1}{2} + c(\mathbf{x}, \mathbf{x}') \right) \ell(g(\mathbf{x}'), +1) \right] \\
 786 &\quad + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[\left(\frac{1}{2} + c(\mathbf{x}, \mathbf{x}') \right) \ell(g(\mathbf{x}), -1) + \left(\frac{1}{2} - c(\mathbf{x}, \mathbf{x}') \right) \ell(g(\mathbf{x}'), -1) \right] \\
 787 &\quad + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}, \mathbf{x}')} \left[(1 - 2\pi) (g(\mathbf{x}) + g(\mathbf{x}')) \right]. \tag{29} \\
 788 & \\
 789 & \\
 790 & \\
 791 & \\
 792 & \\
 793 & \\
 794 & \\
 795 & \\
 796 & \\
 797 & \\
 798 & \\
 799 & \\
 800 & \\
 801 & \\
 802 & \\
 803 & \\
 804 & \\
 805 & \\
 806 & \\
 807 & \\
 808 & \\
 809 &
 \end{aligned}$$

Then Eq. 6 is proven.

C LIMITATIONS

The noise generation method we proposed primarily utilizes a Gaussian distribution to perturb confidence difference distributions originally concentrated around specific values, aiming to approximate the confidence difference distributions that may manifest in the real world. Consequently, artificial datasets are utilized. In the future, we may consider annotating pairwise confidence difference datasets derived from real-world scenarios. It would allow for experiments using authentic datasets rather than artificially constructed ones, offering substantial practical significance.

Additionally, the datasets used are actually multi-label datasets although we focus on binary classification problems in weakly supervised learning. Then the labels of these multi-label datasets are partitioned into two disjoint subsets, each serving as positive and negative classes, respectively, thereby converting them into binary classification datasets. In the future, we will consider expanding the problem scenario to multi-label classification.

D BROADER IMPACTS

The noise confidence difference classification proposed in this paper stands to notably improve decision accuracy in real-world settings. It addresses potential noise impacts present in real-world data and holds substantial practical significance as a plausible scenario in weakly supervised domains. Its applicability can be extended to various fields including medical diagnosis, rehabilitation assessment, and financial risk management.

However, it's important to acknowledge that the confidence difference utilized in our method within weakly supervised settings might be influenced by potential data biases inherent in the real world. Furthermore, we demonstrate the effectiveness of our approach in weakly supervised scenarios, there's a risk of excessive dependence on algorithms for decision-making, potentially overlooking the cultivation of individual decision-making capabilities and autonomy.