Ordinal Adaptive Correction: A Data-Centric Approach to Ordinal Image Classification with Noisy Labels

Alireza Sedighi Moghaddam, Mohammad Reza Mohammadi

Abstract—Labeled data is a fundamental component in training supervised deep learning models for computer vision tasks. However, the labeling process, especially for ordinal image classification where class boundaries are often ambiguous, is prone to error and noise. Such label noise can significantly degrade the performance and reliability of machine learning models. This paper addresses the problem of detecting and correcting label noise in ordinal image classification tasks. To this end, a novel data-centric method called ORDinal Adaptive Correction (ORDAC) is proposed for adaptive correction of noisy labels. The proposed approach leverages the capabilities of Label Distribution Learning (LDL) to model the inherent ambiguity and uncertainty present in ordinal labels. During training, ORDAC dynamically adjusts the mean and standard deviation of the label distribution for each sample. Rather than discarding potentially noisy samples, this approach aims to correct them and make optimal use of the entire training dataset. The effectiveness of the proposed method is evaluated on benchmark datasets for age estimation (Adience) and disease severity detection (Diabetic Retinopathy) under various asymmetric Gaussian noise scenarios. Results show that ORDAC and its extended versions (ORDAC_C and ORDAC_R) lead to significant improvements in model performance. For instance, on the Adience dataset with 40% noise, ORDAC_R reduced the mean absolute error from 0.86 to 0.62 and increased the recall metric from 0.37 to 0.49. The method also demonstrated its effectiveness in correcting intrinsic noise present in the original datasets. This research indicates that adaptive label correction using label distributions is an effective strategy to enhance the robustness and accuracy of ordinal classification models in the presence of noisy data.

Index Terms—Label Error Detection, Noisy Label Correction, Ordinal Classification, Label Noise, Label Distribution Learning, Data-centric Artificial Intelligence

I. INTRODUCTION

THE performance, fairness, and robustness of modern artificial intelligence (AI) systems are fundamentally shaped by the quality of the data they are trained on. This has given rise to a data-centric view of AI, which posits that systematic improvements in data quality are often more impactful than architectural innovations alone [1]. In computer vision, the success of deep neural networks (DNNs) has been fueled by large-scale annotated datasets like ImageNet [2]. However, creating such datasets is a costly and labor-intensive process. Consequently, researchers often rely on scalable but less controlled methods like crowdsourcing or web scraping, which inevitably introduce noisy or incorrect labels into the training data [3].

AS. Moghaddam, and MR. Mohammadi are with School of Computer Engineering, Iran University of Science and Technology, Islamic Republic of Iran (e-mail: a_sedighi77@comp.iust.ac.ir; mrmohammadi@iust.ac.ir), Corresponding author: MR. Mohammadi.

The prevalence of noisy labels in benchmark datasets such as ImageNet [2], CIFAR-10 [4], and MNIST [5] presents a significant challenge. DNNs trained on such data can exhibit severely degraded performance, as they may memorize the incorrect annotations. This raises a fundamental question: how can we effectively train models in the presence of label noise? This problem is particularly acute in ordinal classification (or ordinal regression), where the goal is to predict a label from a set of classes with an inherent order, such as estimating age, grading disease severity, or ranking customer satisfaction. In these tasks, the semantic closeness of adjacent classes makes annotation inherently ambiguous and increases the likelihood of label noise.

Existing approaches to mitigate label noise fall into two main categories. Model-centric methods aim to make the learning process itself robust, for example, by designing noise-tolerant loss functions or using regularization to prevent overfitting to incorrect labels [6]. In contrast, data-centric approaches focus on the data itself, with the dominant strategy being sample selection. Methods like CASSOR [7] and ICDF [8] identify and filter out samples that are likely mislabeled, training the model on a cleaner subset. While effective, the primary drawback of this approach is the outright removal of samples, which discards the potentially valuable information contained within the features of those instances. This reveals a significant research gap: a need for methods that can correct rather than discard noisy labels by adaptively modeling the uncertainty inherent in ordinal data.

To address this gap, this paper introduces a novel framework named **ORD**inal **A**daptive **C**orrection (**ORDAC**). Instead of discarding samples suspected of having noisy labels, ORDAC is designed to correct them by leveraging the expressive power of Label Distribution Learning (LDL) [9]. The core idea is to represent each label not as a single value but as a Gaussian distribution, where the mean represents the label value and the standard deviation quantifies the uncertainty. The ORDAC framework operates iteratively, using the model's own predictions in a cross-validation setup to dynamically update the mean and standard deviation of each training sample's label distribution. This adaptive correction mechanism allows the model to gradually learn from progressively cleaner and more reliable labels, enhancing its robustness and generalization.

The main contributions of this work are:

- We propose a novel framework, ORDAC, that shifts the paradigm for handling noisy ordinal labels from sample selection to adaptive label correction.
- We introduce a mechanism to dynamically update both the mean and standard deviation of label distributions,

allowing the model to explicitly represent and manage uncertainty during training.

 We demonstrate through extensive experiments on realworld datasets that our correction-based approach significantly outperforms both standard training methods and state-of-the-art sample selection techniques in the presence of label noise.

By focusing on improving data quality at a fundamental level, this research paves the way for more accurate and reliable ordinal classification, particularly in domains where clean data is scarce and label ambiguity is high. Our implementation is publicly available at https://github.com/AlirezaSM/ORDAC/.

The structure of the paper is as follows: Section 2 reviews related works. Section 3 describes methodology and the process of noisy label correction is explained. Section 4 is dedicated to evaluating the proposed method, and the experiments conducted to assess its performance are presented. Finally, Section 5 provides the conclusions and potential directions for future work.

II. RELATED WORK

Our work lies at the intersection of noisy-label learning, ordinal classification, and label distribution learning. We survey these areas to highlight the need for methods that go beyond sample filtering toward active label correction for ordinal data.

A. Learning from Noisy Labels

The paradigm of Data-Centric AI emphasizes that improving data quality is a cornerstone of building robust models [10]. A central challenge within this paradigm is learning from data with noisy labels. The strategies developed to address this can be broadly grouped into two families.

Model-centric approaches aim to make the learning algorithm itself resilient to noise. This includes designing robust loss functions like Generalized Cross-Entropy [11] that are less sensitive to large errors, or employing regularization techniques like Mixup [12] that discourage the model from memorizing incorrect labels. While often effective, these methods treat the dataset as a static, unchangeable entity.

Data-centric approaches, in contrast, focus on improving the dataset itself [3]. These methods generally fall into two categories: sample selection and label correction. The more common strategy has been sample selection, which aims to detect and reduce the effect of mislabeled samples. Many early works relied on the small-loss trick, based on the observation that deep networks tend to learn clean samples earlier. This idea underpins methods such as Co-teaching [13], where two networks exchange small-loss samples, and MentorNet [14], which learns a curriculum to prioritize easier, likely clean samples. Over time, this paradigm has evolved into label quality scoring frameworks. For example, Confident Learning [15] provides a principled, model-agnostic way to identify label errors by analyzing the relationship between noisy and predicted labels. Similarly, model-agnostic label quality scoring [16] introduces multiple scoring metrics to detect low-quality samples, and Dataset Cartography [17] visualizes learning dynamics to identify hard-to-learn samples that often correspond to mislabeled data. These approaches have proven effective not only in image classification but also in more complex tasks such as object detection [18] and semantic segmentation [19].

In contrast, label correction offers a more direct, datarestorative strategy. Rather than discarding noisy samples, these methods try to infer the correct labels, preserving the full dataset for training. Although promising, this line of work remains less explored compared to sample selection. Some methods adopt a meta-learning framework, where a metamodel learns a correction function for the noisy dataset, as in Meta Label Correction [20]. Others exploit the geometric structure of the feature space, for example, by using graphbased label propagation [21] to refine labels. Despite their strengths, most correction methods are designed for nominal classification tasks. Our work, ORDAC, advances this area by introducing a new label correction mechanism specifically tailored to the unique structural properties of ordinal data.

B. Ordinal Classification and Its Challenges

Ordinal classification, or ordinal regression, addresses supervised learning tasks where labels possess a natural order, such as age estimation or clinical severity grading [22]. Ignoring this intrinsic order and treating the problem as nominal classification is suboptimal, as the cost of misclassification is not uniform; for instance, a one-rank error is far less severe than a five-rank error.

To leverage this structure, a diverse set of methodologies has been developed. A prominent family of methods is based on ordinal binary decomposition, which reframes the K-rank problem into a series of simpler binary classification tasks (e.g., is the rank > k?). Early deep learning approaches like OR-CNN [23] applied this principle, but a key challenge was ensuring that the binary predictions were monotonically consistent. The CORAL framework [24] elegantly solved this by sharing weights across the binary classifiers, a technique further refined by its successor, CORN [25]. Other major paradigms include threshold models, which learn a mapping to a latent continuous score that is then partitioned by a set of learned thresholds [26], [27], and the direct design of ordinal loss functions that explicitly penalize predictions based on their rank distance from the true label, such as the Weighted Kappa Loss [28]. While powerful, these methods are typically designed for clean datasets and their performance can degrade significantly in the presence of label noise, as they lack an explicit mechanism to handle incorrect rank annotations.

C. Label Distribution Learning for Ordinal Tasks

Label Distribution Learning (LDL) has emerged as a particularly well-suited paradigm for ordinal tasks due to its ability to handle label ambiguity [9]. Instead of a single hard label, LDL assigns a probability distribution over all possible labels to each instance. For ordinal data, this is often a unimodal distribution (e.g., a Gaussian), which naturally captures the semantic closeness of adjacent classes and the decreasing likelihood of distant ranks.

LDL methods can be categorized as Fixed-form (FLDL) or Adaptive (ALDL) [29]. FLDL methods like DLDL-v2 [30] assume a static shape for the label distribution (e.g., a Gaussian with a fixed standard deviation). This is a strong assumption that limits the model's capacity to represent varying levels of uncertainty across different samples. ALDL methods, such as those using a Unimodal-Concentrated Loss [29] or modeling ordinal relationships explicitly [31], offer more flexibility. However, it is crucial to note that the reported gains of many specialized loss functions have been questioned by studies highlighting the disproportionate impact of evaluation protocols, such as inconsistent data splitting, on final performance [32]. While our framework initializes with a fixed-form distribution, characteristic of FLDL, its core novelty lies in a dynamic correction mechanism that adaptively modifies each sample's distribution, thus operating as an ALDL method.

D. Learning Ordinal Regression with Noisy Labels

The specific challenge of learning ordinal regression from noisy labels is an emerging research frontier. The few existing methods are primarily data-centric and have focused on sample selection. CASSOR [7], for example, proposes a class-aware selection strategy that estimates an "insufficiency score" to dynamically adjust the sampling rate for each class. Similarly, ICDF [8] introduces a filtering algorithm based on inter-class feature-space distances to identify and remove noisy samples.

While these state-of-the-art methods demonstrate strong performance, they share a fundamental limitation: they are designed to identify and discard samples. This strategy, by its nature, can lead to the loss of valuable feature information contained in the discarded instances, especially in data-scarce regimes. Furthermore, these methods do not leverage the unique capabilities of LDL to explicitly model label uncertainty as a means to correct, rather than simply remove, noisy labels. Our work is motivated by this clear and critical gap in the literature. We propose a method that uses the principles of LDL not for filtering, but for the adaptive correction of noisy ordinal labels, thereby preserving data while improving its quality throughout the training process.

III. METHODOLOGY

To address the challenge of label noise in ordinal classification, we propose a novel data-centric framework named **ORD**inal **A**daptive **C**orrection (**ORDAC**). Unlike existing methods that discard potentially noisy samples, ORDAC is designed to correct them by leveraging the principles of Label Distribution Learning (LDL). Our core idea is to represent each ordinal label not as a discrete value, but as a Gaussian distribution characterized by a mean (μ) and a standard deviation (σ). In this representation, μ corresponds to the label's value, while σ quantifies the model's uncertainty about that label. The framework iteratively refines both μ and σ for each training sample, using the model's own evolving knowledge to dynamically clean the dataset.

A. Framework Overview

The ORDAC framework operates using a cross-training strategy to prevent a model from being biased by its own confident but potentially incorrect predictions on data it has already seen. This ensures that label corrections are always guided by reliable, out-of-sample predictions. The overall workflow, depicted in Figure 1, consists of the following key stages:

3

- 1) Data Partitioning and Model Setup: The training dataset D is split into K folds. We then create K distinct training configurations. For each configuration $k \in \{1, \ldots, K\}$, one fold is designated as the validation set (D_k^{valid}) , and the remaining K-1 folds constitute the training set (D_k^{train}) . We initialize K identical Fixedform LDL (FLDL) models, M_1, \ldots, M_K , one for each configuration.
- 2) Initialization of Label Distributions: For a noisy dataset $D = \{(x_i, \tilde{y_i})\}_{i=1}^N$, we initialize a Gaussian label distribution $\mathcal{N}(\mu_i, \sigma_i^2)$ for each sample x_i . The initial mean μ_i is set to the provided noisy label $\tilde{y_i}$, and the standard deviation σ_i is initialized to a fixed, constant value for all samples, representing uniform initial uncertainty.
- 3) Warm-up Phase: All K models are trained concurrently on their respective training sets for a specified number of warm-up epochs ($E_{\rm corr}$). This allows the models to learn initial feature representations and converge to a stable state on the original noisy data before any corrections are made.
- 4) Iterative Correction Phase: After the warm-up phase, for each subsequent epoch, the adaptive correction mechanism is activated. The label distribution for each sample in the training set is updated based on the predictions from the corresponding model that held it out as a validation sample. For instance, the corrected labels from D_k^{valid} (generated by model M_k) are used to update the corresponding samples in all other training sets D_j^{train} where $j \neq k$. This process is detailed in Section III-B and Algorithm 1.

B. Adaptive Label Correction Mechanism

The core of ORDAC is its two-stage adaptive correction process, which is applied iteratively to each sample in the designated validation fold during the correction phase. This process first corrects for systematic model bias and then performs a sample-specific update of the label distribution.

1) Class-wise Prediction Debiasing: Ordinal regression models are prone to developing a systematic bias towards middle-rank classes. This phenomenon arises from two primary factors. First, the distance-sensitive loss functions commonly used in ordinal regression encourage a conservative model to predict middle-rank classes, as this strategy minimizes the upper bound of the loss in the worst-case scenario. Second, this tendency is often exacerbated by inherent data imbalance, where middle-rank classes frequently contain more samples than those at the extremes of the spectrum. To counteract this and prevent our iterative correction process

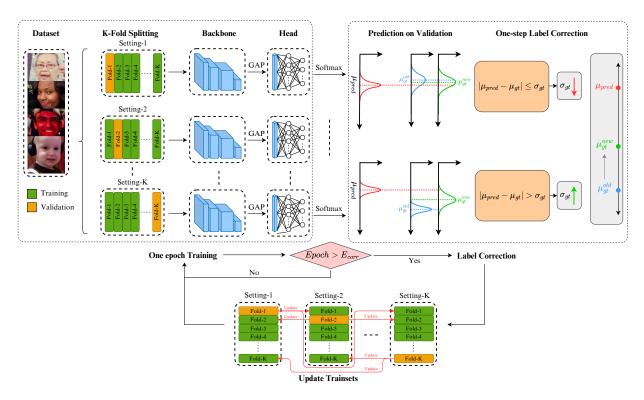


Fig. 1. An overview of the proposed ORDAC framework. Using a K-fold setup, models are trained on training folds and make predictions on a validation fold. These predictions are used to correct the label distributions (mean and standard deviation) of the validation samples, which are then propagated back to the training sets for subsequent epochs.

from collapsing towards a biased mean, we first perform a class-wise debiasing of the model's predictions. For each class c, we compute the mean prediction of the model across all samples currently labeled as belonging to that class:

$$\operatorname{mean}_{c} = \frac{1}{N_{c}} \sum_{i=1}^{N} \mathbf{1}_{1} [\mu_{i} = c] \hat{y}_{i}$$
 (1)

where N_c is the number of samples in class c, μ_i is the current mean of the label distribution for sample i, \hat{y}_i is the model's predicted mean for sample i, and $\mathbf{1}_1[\cdot]$ is the indicator function.

We then shift the prediction for each sample to re-center the class-wise predictions around the true class label. This assumes the label noise is unbiased and does not significantly alter the true mean of each class. The shifted prediction $\hat{y}_i^{\text{shifted}}$ for a sample i with label mean μ_i is:

$$\hat{y}_i^{\text{shifted}} = \hat{y}_i - (\text{mean}_{\mu_i} - \mu_i)$$
 (2)

2) Sample-wise Distribution Update: Following the debiasing step, we perform an adaptive update of the mean (μ_i) and standard deviation (σ_i) for each sample. The magnitude of this update is governed by a sample-specific correction coefficient, λ_i^{corr} , which balances model confidence with class frequency:

$$\lambda_i^{\text{corr}} = \left(\frac{\gamma_i}{1 - \log(\pi_{\mu_i} + \epsilon)}\right) \tag{3}$$

Here, γ_i is a measure of the model's confidence in its prediction for sample i. $\pi_{\mu_i} = N_{\mu_i}/N$ is the prior probability of the sample's current class, which ensures that samples from rare

classes are corrected more cautiously. ϵ is a small constant to prevent numerical instability.

This coefficient scales the base learning rates, α_{base} and β_{base} , to produce sample-specific update rates:

$$\alpha_i = \alpha_{\text{base}} \times \lambda_i^{\text{corr}}, \quad \beta_i = \beta_{\text{base}} \times \lambda_i^{\text{corr}}$$
 (4)

The prediction error, e_i , is calculated using the debiased prediction:

$$e_i = \hat{y}_i^{\text{shifted}} - \mu_i \tag{5}$$

The standard deviation is then updated based on the relationship between the prediction error and the current uncertainty:

$$\sigma_i^{\text{new}} = \sigma_i + \alpha_i \times (|e_i| - \sigma_i) \tag{6}$$

This rule intuitively increases the uncertainty (σ_i) if the model's error is larger than the current uncertainty, suggesting the label is likely noisy. Conversely, it decreases uncertainty if the error is small, indicating confidence in the label.

Finally, the mean of the label distribution is corrected by shifting it towards the model's prediction:

$$\mu_i^{\text{new}} = \mu_i + \beta_i \times e_i \tag{7}$$

This entire process is summarized in Algorithm 1.

C. Method Variants

To better understand the behavior of our framework, we introduce two variants designed to isolate specific effects of the correction process:

5

Algorithm 1 The ORDAC Algorithm for Ordinal Adaptive Correction

```
1: Input: Noisy dataset D = \{(x_i, \tilde{y_i})\}_{i=1}^N, folds K, epochs E_{\text{max}}, correction start epoch E_{\text{corr}}, standard deviation \sigma_i
 2: Output: Cleaned dataset D_{\text{clean}}, trained models \{M_k\}_{k=1}^K
 3: Initialize K data splits \{D_k\}_{k=1}^K and K models \{M_k\}_{k=1}^K.
4: Initialize label distributions \{\mathcal{N}(\mu_i, \sigma_i^2)\}_{i=1}^N with \mu_i = \tilde{y_i}.
 5: for e=1 to E_{\max} do
           for each configuration k = 1 to K do
 6:
                Train model M_k on D_k^{\text{train}} for one epoch.
 7:
 8:
           if e \geq E_{\rm corr} then
                Initialize a temporary set for corrected labels D_{\text{corr}} = \emptyset.
 9:
                for each configuration k = 1 to K do
10:
                      Get predictions \hat{Y}_k for validation set D_k^{\text{valid}} from model M_k.
11:
                      Compute class-wise means and get shifted predictions \hat{Y}_k^{\text{shifted}} using Eq. 2.
12:
                      for each sample (x_i, \mu_i, \sigma_i) \in D_k^{\text{valid}} do
13:
                           Compute correction coefficient \lambda_i^{\text{corr}}, update rates \alpha_i, \beta_i.
14:
                           Compute error e_i = \hat{y}_i^{\text{shifted}} - \mu_i.
15:
                           Update \sigma_i^{\text{new}} = \sigma_i + \alpha_i \times (|e_i| - \sigma_i).
16:
                           Update \mu_i^{\text{new}} = \mu_i + \beta_i \times e_i.
17:
                           Add (x_i, \mu_i^{\text{new}}, \sigma_i^{\text{new}}) to D_{\text{corr}}.
18:
                Update all training sets D_k^{\text{train}} with the corrected distributions from D_{\text{corr}}.
19:
20: D_{\text{clean}} = \text{union of all final corrected validation folds.}
21: Return D_{\text{clean}}, \{M_k\}_{k=1}^K.
```

- ORDAC_C (Correct): This variant is designed to demonstrate that our correction process genuinely improves the overall quality of the dataset. Here, the full iterative ORDAC process is run to generate a static, cleaned dataset. A new model is then trained from scratch on this corrected dataset without any further online corrections. Strong performance from this variant indicates that the corrected labels are of high quality.
- ORDAC_R (Remove): This variant explores the impact of removing samples that remain highly uncertain even after correction. It builds on ORDAC_C by first generating a corrected dataset. It then identifies and removes samples whose uncertainty (standard deviation) failed to decrease from its initial value (i.e., $\sigma_i^{\text{new}} \geq \sigma_i^{\text{initial}}$), treating them as outliers. A new model is then trained on this smaller, filtered-and-corrected dataset. This allows us to study the synergy between label correction and the removal of hard-to-correct samples.

IV. RESULTS

In this section, we present a comprehensive empirical evaluation of our proposed method, ORDAC. We assess its performance on two real-world ordinal classification datasets, one for age estimation and one for medical image analysis. We first detail the experimental setup, including the datasets, evaluation metrics, and noise injection protocol. We then analyze the robustness of ORDAC against controlled label noise, compare it with state-of-the-art methods, and provide in-depth ablation studies to validate our design choices.

A. Experimental Setup

1) Datasets: We use two challenging, publicly available datasets:

- Adience [33]: This dataset contains approximately 26K unfiltered face images from Flickr albums, annotated for age and gender. We focus on the age estimation task, which consists of 8 ordinal classes: (0-2), (4-6), (8-13), (15-20), (25-32), (38-43), (48-53), and (60+). The dataset's in-the-wild nature, with variations in lighting, pose, and resolution, makes it a challenging benchmark. Crucially, the labels are derived from user profiles and are known to contain inherent noise due to self-reporting errors or mismatches between the uploader and the person in the photo. We use the aligned version of the dataset and the official 5-fold split as recommended by [32].
- Diabetic Retinopathy (DR) [34]: This medical imaging dataset contains over 88K high-resolution retinal fundus images for detecting and grading diabetic retinopathy. The task is to classify each image into one of 5 ordinal severity levels: (0: No DR, 1: Mild, 2: Moderate, 3: Severe, 4: Proliferative DR). The dataset exhibits significant challenges, including variations in imaging conditions and image quality. We use the official training/test split and further divide the training data into 5 folds for our cross-validation setup.
- 2) Evaluation Metrics: Given the class imbalance common in ordinal datasets, we use macro-averaged metrics to ensure a fair evaluation across all classes. We report:
 - Macro-Averaged Mean Absolute Error (MAE): The primary metric for ordinal regression, measuring the average absolute difference between the predicted and true class ranks.
 - Macro-Averaged Recall (Accuracy): The average of per-class recall, measuring the classification accuracy without bias towards majority classes.

TABLE I

MAIN RESULTS COMPARING OUR PROPOSED METHODS (ORDAC AND ITS VARIANTS) AGAINST BASELINES ON THE ADIENCE AND DR DATASETS UNDER
DIFFERENT NOISE RATES (τ). WE REPORT MACRO-AVERAGED MAE AND RECALL (REC). LOWER MAE AND HIGHER REC ARE BETTER. BEST
RESULTS ARE IN **BOLD**, SECOND BEST ARE <u>UNDERLINED</u>.

Dataset	au	Metric	CORAL [24]	DLDL-v2 [30]	ORDAC	ORDAC _C	ORDACR
Adience	0.0	MAE REC	0.6640 ± 0.0287 0.4669 ± 0.0135	$\begin{array}{c} 0.6343 \pm 0.0505 \\ 0.5063 \pm 0.0221 \end{array}$	$\begin{array}{c} 0.5585 \pm 0.0606 \\ 0.5542 \pm 0.0353 \end{array}$	$\begin{array}{c} 0.4929 \pm 0.0296 \\ 0.5968 \pm 0.0193 \end{array}$	$\begin{array}{c} 0.5018 \pm 0.0274 \\ \hline 0.5954 \pm 0.0136 \end{array}$
	0.2	MAE REC	$\begin{array}{c} 0.8349 \pm 0.0289 \\ 0.3985 \pm 0.0102 \end{array}$	$\begin{array}{c} 0.7618 \pm 0.0482 \\ 0.4452 \pm 0.0211 \end{array}$	0.6463 ± 0.0394 0.4988 ± 0.0226	$\begin{array}{c} 0.5642 \pm 0.0400 \\ \hline 0.5460 \pm 0.0182 \end{array}$	$\begin{array}{c} \textbf{0.5366} \pm \textbf{0.0347} \\ \textbf{0.5624} \pm \textbf{0.0185} \end{array}$
	0.4	MAE REC	$\begin{array}{c} 1.0427 \pm 0.0676 \\ 0.3312 \pm 0.0245 \end{array}$	$\begin{array}{c} 0.8649 \pm 0.0382 \\ 0.3775 \pm 0.0290 \end{array}$	0.7192 ± 0.0393 0.4600 ± 0.0119	$\begin{array}{c} 0.6637 \pm 0.0326 \\ \hline 0.4813 \pm 0.0218 \end{array}$	$\begin{array}{c} 0.6283 \pm 0.0379 \\ 0.4950 \pm 0.0221 \end{array}$
DR	0.0	MAE REC	$\begin{array}{c} 0.7721 \pm 0.0173 \\ 0.4287 \pm 0.0052 \end{array}$	$\begin{array}{c} 0.7324 \pm 0.0268 \\ 0.4509 \pm 0.0136 \end{array}$	$\begin{array}{c} 0.6826 \pm 0.0122 \\ 0.4665 \pm 0.0103 \end{array}$	$\begin{array}{c} 0.7084 \pm 0.0103 \\ 0.4559 \pm 0.0058 \end{array}$	$\begin{array}{c} 0.6924 \pm 0.0172 \\ \hline 0.4599 \pm 0.0063 \end{array}$
	0.2	MAE REC	0.8403 ± 0.0347 0.3943 ± 0.0181	$\begin{array}{c} 0.8025 \pm 0.0260 \\ 0.3894 \pm 0.0176 \end{array}$	$\begin{array}{c} \textbf{0.7114} \pm \textbf{0.0202} \\ \textbf{0.4363} \pm \textbf{0.0088} \end{array}$	$\begin{array}{c} 0.7362 \pm 0.0113 \\ 0.4208 \pm 0.0033 \end{array}$	$\begin{array}{c} 0.7246 \pm 0.0060 \\ \hline 0.4217 \pm 0.0023 \end{array}$
	0.4	MAE REC	$\begin{array}{c} 0.8595 \pm 0.0436 \\ 0.3702 \pm 0.0189 \end{array}$	$\begin{array}{c} 0.8488 \pm 0.0493 \\ 0.3452 \pm 0.0236 \end{array}$	$\begin{array}{c} 0.8015 \pm 0.0620 \\ 0.3722 \pm 0.0323 \end{array}$	$\frac{0.7567 \pm 0.0168}{0.4168 \pm 0.0057}$	$\begin{array}{c} \textbf{0.7436} \pm \textbf{0.0111} \\ \textbf{0.4200} \pm \textbf{0.0059} \end{array}$

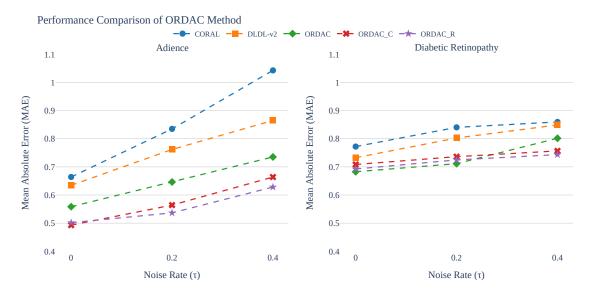


Fig. 2. MAE of the proposed methods and baselines on the Adience and DR datasets as a function of the injected noise rate (τ). Lower is better.

3) Noise Injection Protocol: Label noise in ordinal tasks is typically asymmetric, as mislabelings are more likely to occur between adjacent, semantically similar classes. To simulate this realistically, we inject **Gaussian Asymmetric Noise** into the training and validation sets, following the protocol used in prior works [35], [7], [8]. A noise transition matrix T, where $T_{ij} = P(\tilde{y} = j|y = i)$ is the probability of a true label i being flipped to a noisy label j, is generated. The off-diagonal elements are defined by a Gaussian function of the distance between ranks:

$$T_{ij} \propto \exp\left(-\frac{(i-j)^2}{2\sigma_n^2}\right) \quad \text{for } i \neq j$$
 (8)

where σ_n controls the spread of the noise (set to 3 in our experiments). The matrix is then normalized such that the overall noise rate corresponds to a target value τ . The test sets are always kept clean for evaluation.

4) Implementation Details: All experiments use a **ResNet-50** architecture, pre-trained on ImageNet, as the feature ex-

tractor backbone. For our proposed method, the model is trained for a maximum of 50 epochs ($E_{\rm max}$), with the adaptive correction mechanism activating after 10 warm-up periods ($E_{\rm corr}$). The base learning rates for correction were set to $\alpha_{\rm base}=0.2$ and $\beta_{\rm base}=0.8$ based on the hyperparameter analysis in IV-E2.

B. Performance under Label Noise

We first evaluate the robustness of ORDAC against synthetic label noise. We compare our three proposed variants (ORDAC, ORDAC, and ORDACR) with two strong baselines: **CORAL** [24], a standard ordinal regression method, and **DLDL-v2** [30], a fixed-form LDL method. The training and validation sets are corrupted with Gaussian asymmetric noise at rates $\tau \in \{0.2, 0.4\}$. The case $\tau = 0$ corresponds to training on the original, uncorrupted datasets.

The results are presented in Table I and visualized in Figure 2. Our proposed methods consistently and significantly



Fig. 3. Examples of successful (green box) and unsuccessful (red box) corrections on the Adience dataset with synthetic noise ($\tau = 0.4$).

outperform the baselines across both datasets and all noise levels. As expected, the performance of all methods degrades as the noise rate increases. However, the ORDAC framework demonstrates substantially greater resilience. For instance, on Adience with 40% noise, $ORDAC_R$ achieves an MAE of 0.6283, a dramatic improvement over CORAL (1.0427) and DLDL-v2 (0.8649).

Interestingly, even with no injected noise ($\tau=0$), our methods still improve performance over the baselines. This strongly suggests that the original datasets contain inherent label noise, which ORDAC successfully identifies and corrects.

C. Analysis of the Correction Mechanism

To see how our method works in practice, we show several examples of its label corrections in Figure 3. The figure highlights two main outcomes. Successful corrections, where a noisy label was fixed, are marked with green boxes. Apparent errors, where a clean label was changed, are marked with red boxes.

Crucially, some of these apparent errors (indicated by a small green box) are actually plausible corrections. This happens when the original "clean" label in the dataset was already incorrect. This result shows that our method can fix not only the noise we add for experiments but also the hidden errors that already exist in the dataset.

Quantitatively, we analyze the magnitude of corrections on the original Adience dataset ($\tau=0$). Figure 4 shows that most corrections are small (a shift of 1 class), with very few large changes, indicating that the model is making targeted, conservative adjustments rather than drastic, random changes.

We also measure the MAE and RMSE between the true labels and the noisy/corrected labels in Table II. While ORDAC sometimes increases the MAE of the labels themselves (e.g., for DR at $\tau=0.2$), it consistently improves the final model's performance on the test set (Table I). This indicates that simply measuring label correctness is insufficient; the ultimate goal

Histogram of Label Correction Differences on Adience with τ =0

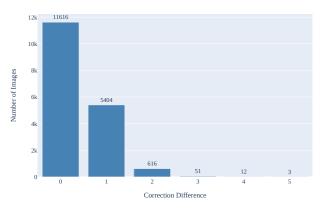


Fig. 4. Histogram of label changes made by ORDAC on the original (clean) Adience dataset.

TABLE II

MAE AND RMSE BETWEEN THE TRUE LABELS AND THE LABELS USED FOR TRAINING (NOISY FOR DLDL-V2, CORRECTED FOR ORDAC).

Dataset	au	Metric	DLDL-v2 (Noisy)	ORDAC (Corrected)
	0.0	MAE	0.0	0.4186
	0.0	RMSE	0.0	0.7116
Adience	0.2	MAE	0.4218	0.4933
	0.2	RMSE	1.1214	0.8070
	0.4	MAE	0.8727	0.6481
	0.4	RMSE	1.6493	0.9894
	0.0	MAE	0.0	0.6106
	0.0	RMSE	0.0	0.9265
DR	0.2	MAE	0.3508	0.6462
	0.2	RMSE	0.8845	0.9431
	0.4	MAE	0.6658	0.6688
	0.4	RMSE	1.2106	1.0096

TABLE III

MAE COMPARISON WITH THE CASSOR SAMPLE SELECTION METHOD. "NORMAL TRAINING" USES ALL NOISY DATA. "ORDAC $_{\rm C}$ + CASSOR" APPLIES CASSOR'S SELECTION TO OUR CORRECTED DATASET.

Dataset	au	Normal Training	CASSOR	ORDAC _C	ORDAC _C + CASSOR	ORDAC _R
Adience	0.0 0.2 0.4	0.5727 ± 0.0366 0.6969 ± 0.0427 0.8200 ± 0.0632	0.6158 ± 0.0611 0.5902 ± 0.0442 0.5908 ± 0.0605	$\begin{array}{c} 0.5038 \pm 0.0458 \\ \hline 0.5329 \pm 0.0348 \\ \hline 0.6321 \pm 0.0288 \end{array}$	$\begin{array}{c} 0.5671 \pm 0.0449 \\ 0.5655 \pm 0.0497 \\ \underline{0.5918 \pm 0.0523} \end{array}$	0.5033 ± 0.0353 0.5326 ± 0.0321 0.6159 ± 0.0316
DR	0.0 0.2 0.4	$\begin{array}{c} 0.6712 \pm 0.0040 \\ \hline 0.7447 \pm 0.0195 \\ 0.8597 \pm 0.0258 \end{array}$	0.6709 ± 0.0056 0.7930 ± 0.0543 0.8240 ± 0.0736	$\begin{array}{c} 0.6893 \pm 0.0107 \\ 0.7283 \pm 0.0100 \\ \hline 0.7532 \pm 0.0102 \end{array}$	0.7267 ± 0.0338 0.7453 ± 0.0237 0.7559 ± 0.0211	0.6828 ± 0.0102 0.7149 ± 0.0086 0.7435 ± 0.0120

is to produce a more generalizable model, which ORDAC achieves. The improved model performance suggests that even when a label is not corrected perfectly, it is often moved closer to the true value, or its associated uncertainty is adjusted appropriately, leading to better overall training dynamics.

D. Comparison with State-of-the-Art Sample Selection

We compare our correction-based approach with **CASSOR** [7], a state-of-the-art sample selection method for noisy ordinal regression. To ensure a fair comparison, we use their official implementation and train it with the same ResNet-50 backbone. We evaluate training with CASSOR's selected clean subset versus training with our corrected labels (ORDAC $_{\rm C}$ and ORDAC $_{\rm R}$).

The results in Table III show that our correction-based methods generally outperform sample selection, especially as the noise rate increases. On the DR dataset with 40% noise, ORDAC_R achieves an MAE of 0.7435, significantly better than CASSOR's 0.8240. This suggests that correcting noisy labels is a more data-efficient strategy than simply discarding them. An interesting case is Adience at $\tau=0.4$, where CASSOR performs best. This may indicate that at very high noise levels, aggressive filtering can be more effective than attempting to correct highly unreliable labels. We also test a hybrid approach (ORDAC_C + CASSOR), which yields strong results, suggesting that correction and selection are complementary.

E. Ablation Studies

- 1) Effect of Class-wise Prediction Debiasing: We evaluate the importance of the class-wise debiasing step described in Section III-B. Table IV shows that removing this step significantly degrades performance, especially at higher noise rates. As shown in Figure 5, without debiasing, the correction process develops a strong bias towards the majority middle class, drastically depleting the minority classes at the ends of the spectrum. The debiasing step successfully mitigates this, ensuring a more balanced and accurate correction process.
- 2) Hyperparameter Analysis: We analyzed the sensitivity to the base correction rates, $\alpha_{\rm base}$ and $\beta_{\rm base}$, on Adience with $\tau=0.4$. The results in Table V show that the best performance is achieved with a small $\alpha_{\rm base}$ (0.2) and a large $\beta_{\rm base}$ (0.8). This aligns with our intuition: the model should be aggressive in correcting the label's mean value (μ) but conservative and stable when updating its uncertainty (σ). We also found

TABLE IV
EFFECT OF REMOVING THE CLASS-WISE PREDICTION DEBIASING STEP ON
THE ADIENCE DATASET.

τ	DLDL-v2	ORDAC (w/o debiasing)	ORDAC (w/ debiasing)
0.0	0.6343 ± 0.0505	0.6033 ± 0.0362	0.5585 ± 0.0606
0.2	0.7618 ± 0.0482	0.6960 ± 0.0893	0.6463 ± 0.0394
0.4	0.8649 ± 0.0382	0.8366 ± 0.0516	$\textbf{0.7192}\pm\textbf{0.0393}$

that an initial standard deviation (std_{init}) of 0.75 yielded the best results, providing a good balance for initial uncertainty representation.

TABLE V Hyperparameter tuning results (MAE) for $\alpha_{\rm BASE}$ and $\beta_{\rm BASE}$ on Adience ($\tau=0.4$).

au	α_{base}	β_{base}	ORDAC
	0.2	0.2 0.5 0.8	0.8773 ± 0.0399 0.8269 ± 0.0572 0.7192 ± 0.0393
0.4	0.5	0.2 0.5 0.8	0.8807 ± 0.0331 0.7983 ± 0.0515 0.7330 ± 0.0581
	0.8	0.2 0.5 0.8	0.8760 ± 0.0350 0.8147 ± 0.0681 0.7463 ± 0.0897

V. CONCLUSION

In this paper, we addressed the critical challenge of label noise in ordinal classification, a problem that undermines model performance, particularly in domains with inherent annotation ambiguity. We argued that the dominant data-centric paradigm of sample selection, which discards potentially noisy instances, is suboptimal as it leads to the loss of valuable data. To overcome this limitation, we introduced ORDAC, a novel framework that shifts the focus from sample removal to adaptive label correction. By representing each label as a dynamic Gaussian distribution, our method successfully leverages the model's own evolving knowledge to iteratively refine both the value (mean) and uncertainty (standard deviation) of labels in the training set.

Our extensive experiments on two real-world datasets, Adience and Diabetic Retinopathy, demonstrated the effectiveness



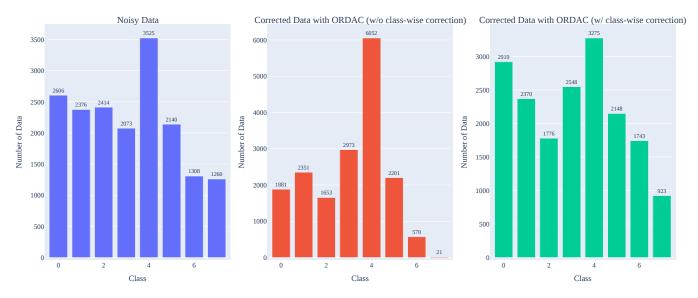


Fig. 5. Number of samples per class before and after correction on Adience ($\tau = 0.4$), with and without the class-wise debiasing step. Debiasing prevents a collapse into the majority class.

of this corrective approach. The results conclusively show that ORDAC and its variants significantly outperform standard ordinal regression baselines and state-of-the-art sample selection methods, especially under high levels of asymmetric label noise. Crucially, we also found that ORDAC improves performance even on the original, uncorrupted datasets, providing strong evidence of its ability to identify and correct inherent, real-world label errors. Our analysis confirmed that this performance gain is driven by a robust mechanism that not only moves noisy labels closer to their true values but also intelligently manages label uncertainty throughout the training process.

For future work, several promising avenues exist. While our method uses a Gaussian distribution, exploring more flexible, non-parametric distributions could allow for modeling more complex types of label ambiguity. Furthermore, developing a mechanism to automatically tune the correction rates based on an online estimation of the dataset's noise level would enhance the framework's autonomy. Finally, applying the principles of ORDAC to other tasks and domains presents an exciting direction for future research.

REFERENCES

- [1] N. Sambasivan, S. Kapania, H. Highfill, D. Akrong, P. Paritosh, and L. M. Aroyo, ""everyone wants to do the model work, not the data work": Data cascades in high-stakes ai," in proceedings of the 2021 CHI Conference on Human Factors in Computing Systems, pp. 1–15, 2021.
- [2] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., "Imagenet large scale visual recognition challenge," *International journal of computer* vision, vol. 115, pp. 211–252, 2015.
- [3] F. R. Cordeiro and G. Carneiro, "A survey on deep learning with noisy labels: How to train your model when you cannot trust on the annotations?," in 2020 33rd SIBGRAPI conference on graphics, patterns and images (SIBGRAPI), pp. 9–16, IEEE, 2020.
- [4] A. Krizhevsky, G. Hinton, et al., "Learning multiple layers of features from tiny images," 2009.

- [5] L. Deng, "The mnist database of handwritten digit images for machine learning research [best of the web]," *IEEE signal processing magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [6] H. Song, M. Kim, D. Park, Y. Shin, and J.-G. Lee, "Learning from noisy labels with deep neural networks: A survey," *IEEE transactions* on neural networks and learning systems, vol. 34, no. 11, pp. 8135– 8153, 2022.
- [7] Y. Yuan, S. Wan, C. Zhang, and C. Gong, "Cassor: Class-aware sample selection for ordinal regression with noisy labels," in *Pacific Rim International Conference on Artificial Intelligence*, pp. 117–123, Springer, 2023.
- [8] G. Jiang, F. Wang, and W. Wang, "Noise cleaning for nonuniform ordinal labels based on inter-class distance," *Applied Intelligence*, vol. 54, no. 11, pp. 6997–7011, 2024.
- [9] X. Geng, "Label distribution learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 7, pp. 1734–1748, 2016.
- [10] D. Zha, Z. P. Bhat, K.-H. Lai, F. Yang, and X. Hu, "Data-centric ai: Perspectives and challenges," in *Proceedings of the 2023 SIAM* international conference on data mining (SDM), pp. 945–948, SIAM, 2023
- [11] Z. Zhang and M. Sabuncu, "Generalized cross entropy loss for training deep neural networks with noisy labels," *Advances in neural information* processing systems, vol. 31, 2018.
- [12] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," arXiv preprint arXiv:1710.09412, 2017.
- [13] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, and M. Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," *Advances in neural information processing* systems, vol. 31, 2018.
- [14] L. Jiang, Z. Zhou, T. Leung, L.-J. Li, and L. Fei-Fei, "Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels," in *International conference on machine learning*, pp. 2304–2313, PMLR, 2018.
- [15] C. Northcutt, L. Jiang, and I. Chuang, "Confident learning: Estimating uncertainty in dataset labels," *Journal of Artificial Intelligence Research*, vol. 70, pp. 1373–1411, 2021.
- [16] J. Kuan and J. Mueller, "Model-agnostic label quality scoring to detect real-world label errors," in *ICML DataPerf Workshop*, 2022.
- [17] S. Swayamdipta, R. Schwartz, N. Lourie, Y. Wang, H. Hajishirzi, N. A. Smith, and Y. Choi, "Dataset cartography: Mapping and diagnosing datasets with training dynamics," arXiv preprint arXiv:2009.10795, 2020.
- [18] U. Tkachenko, A. Thyagarajan, and J. Mueller, "Objectlab: Automated diagnosis of mislabeled images in object detection data," arXiv preprint arXiv:2309.00832, 2023.

- [19] V. Lad and J. Mueller, "Estimating label quality and errors in semantic segmentation data via any model," arXiv preprint arXiv:2307.05080, 2023
- [20] G. Zheng, A. H. Awadallah, and S. Dumais, "Meta label correction for noisy label learning," in *Proceedings of the AAAI conference on artificial* intelligence, vol. 35, pp. 11053–11061, 2021.
- [21] H. Zhang, X. Xing, and L. Liu, "Dualgraph: A graph-based method for reasoning about label noise," in *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition, pp. 9654–9663, 2021.
- [22] P. A. Gutiérrez, M. Perez-Ortiz, J. Sanchez-Monedero, F. Fernandez-Navarro, and C. Hervas-Martinez, "Ordinal regression methods: survey and experimental study," *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 1, pp. 127–146, 2015.
- [23] Z. Niu, M. Zhou, L. Wang, X. Gao, and G. Hua, "Ordinal regression with multiple output cnn for age estimation," in *Proceedings of the IEEE* conference on computer vision and pattern recognition, pp. 4920–4928, 2016
- [24] W. Cao, V. Mirjalili, and S. Raschka, "Rank consistent ordinal regression for neural networks with application to age estimation," *Pattern Recognition Letters*, vol. 140, pp. 325–331, 2020.
- [25] X. Shi, W. Cao, and S. Raschka, "Deep neural networks for rank-consistent ordinal regression based on conditional probabilities," *Pattern Analysis and Applications*, vol. 26, no. 3, pp. 941–955, 2023.
- [26] V. M. Vargas, P. A. Gutiérrez, and C. Hervas-Martinez, "Cumulative link models for deep ordinal classification," *Neurocomputing*, vol. 401, pp. 48–58, 2020.
- [27] R. Herbrich, T. Graepel, and K. Obermayer, "Support vector learning for ordinal regression," in 1999 Ninth International Conference on Artificial Neural Networks ICANN 99.(Conf. Publ. No. 470), vol. 1, pp. 97–102, IET, 1999.
- [28] J. de La Torre, D. Puig, and A. Valls, "Weighted kappa loss function for multi-class classification of ordinal data in deep learning," *Pattern Recognition Letters*, vol. 105, pp. 144–154, 2018.
- [29] Q. Li, J. Wang, Z. Yao, Y. Li, P. Yang, J. Yan, C. Wang, and S. Pu, "Unimodal-concentrated loss: Fully adaptive label distribution learning for ordinal regression," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20513–20522, 2022.
- [30] B.-B. Gao, H.-Y. Zhou, J. Wu, and X. Geng, "Age estimation using expectation of label distribution learning," in *IJCAI*, vol. 1, p. 3, 2018.
- [31] C. Wen, X. Zhang, X. Yao, and J. Yang, "Ordinal label distribution learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 23481–23491, 2023.
- [32] J. Paplhám, V. Franc, et al., "A call to reflect on evaluation practices for age estimation: comparative analysis of the state-of-the-art and a unified benchmark," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1196–1205, 2024.
- [33] E. Eidinger, R. Enbar, and T. Hassner, "Age and gender estimation of unfiltered faces," *IEEE Transactions on information forensics and security*, vol. 9, no. 12, pp. 2170–2179, 2014.
- [34] E. Dugas, Jared, Jorge, and W. Cukierski, "Diabetic retinopathy detection." https://kaggle.com/competitions/diabetic-retinopathy-detection, 2015. Kaggle.
- [35] H. Liu, J. Tu, A. Gao, and C. Li, "Distributed robust support vector ordinal regression under label noise," *Neurocomputing*, vol. 598, p. 128057, 2024