IMPROVING ORDINAL CONFORMAL PREDICTION BY STEPWISE ADAPTIVE POSTERIOR ALIGNMENT

Anonymous authors

Paper under double-blind review

ABSTRACT

Ordinal classification (OC) is widely used in real-world applications to categorize instances into ordered discrete classes. In risk-sensitive scenarios, ordinal conformal prediction (OCP) is used to obtain a small contiguous prediction set containing ground-truth labels with a desired coverage guarantee. However, OC models often fail to accurately model the posterior distribution, which harms the prediction set obtained by OCP. Therefore, we introduce a new method called *Adaptive Posterior Alignment Step-by-Step* (APASS), which reduces the distribution discrepancy to improve the downstream OCP performance. It is designed as a versatile, plug-and-play solution that is easily integrated into any OC model before OCP. APASS first employs an attention-based estimator to adaptively estimate the variance of the posterior distribution using the information in the calibration set, then utilizes a stepwise temperature scaling algorithm to align the posterior variance predicted by OC models to the better variance estimation. Extensive evaluations on 10 real-world datasets demonstrate that APASS consistently boosts the OCP performance of 5 popular OC models.

028

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

029 Ordinal classification (OC) (Diaz & Marathe, 2019; Gao et al., 2017; Geng, 2016; Guo et al., 2008; Can Malli et al., 2016; Huo et al., 2016; Wen et al., 2020) plays a crucial role in high-stakes domains like healthcare (Liu et al., 2019) and finance (Manthoulis et al., 2020) by categorizing instances 031 into ordered discrete classes. Robust uncertainty quantification is critical beyond accurate point predictions to avoid costly or dangerous outcomes caused by prediction errors. To this end, various 033 methods have been developed for estimating predictive uncertainty in deep neural networks, such as 034 confidence calibration (Guo et al., 2017), MC-Dropout (Gal & Ghahramani, 2016), and Bayesian neural networks (Smith, 2013), but they lack formal guarantees. Conformal Prediction (CP) (Vovk et al., 1999; 2005; Lei et al., 2018; Wen et al., 2020; Romano et al., 2020; Angelopoulos & Bates, 037 2021; Angelopoulos et al., 2021) addresses this gap by providing a distribution-free, post-processing 038 approach that generates prediction sets (PS) guaranteed to contain the true label with a specified coverage probability, which generally design non-conformity scores to quantify the deviation the degree between the model's predictive outcomes and the data distribution. 040

041 Recent works on OC demonstrate substantial benefits of assuming the underlying conditional dis-042 tribution to be unimodal for OC tasks (Diaz & Marathe, 2019; Gao et al., 2017; Guha et al., 2024; 043 Belharbi et al., 2019; Cardoso et al., 2023). Some rely on label smoothing methods, which convert 044 one-hot target labels into unimodal prior distributions to be used as the reference for the training loss. Some works learn a non-parametric unimodal distribution as a constraint optimization problem in the loss function. In the unimodal context of ordinal classification, Ordinal Conformal Prediction 046 (OCP) (Lu et al., 2022; Xu et al., 2023) is designed to generate contiguous prediction sets using 047 the posterior distribution predicted by OC models. In contrast to Adaptive Prediction Sets (APS), 048 which calculate the scores by accumulating the sorted softmax probabilities in descending order, Ordinal-APS calculates the score by accumulating softmax probabilities of the contiguous prediction set with the minimum set size. However, the existing OCP methods neglect the possible variance 051 misalignment of the OC models, which leads to inefficient PS. 052

1053 In this work, we empirically observe the variance misalignment between the predicted posterior distribution and the oracle posterior distribution in a synthetic dataset. Specifically, a noticeable

reduction in the size of PS is observed when we align the predicted posterior to the oracle posterior.
 Further, our theoretical analysis supports the empirical findings by demonstrating a decrease in the upper bound of the prediction set size as the predicted posterior approaches the oracle posterior.

Inspired by our analytical findings, we introduce the *Adaptive Posterior Alignment Step-by-Step* (APASS), which serves as a *plug-and-play* component that can be integrated into any OCP framework to enhance its performance. The method consists of two key parts: 1) We introduce an *attention*-based estimator that adaptively estimates the variance misalignment of an input sample by examining similar samples in the calibration set; 2) a stepwise alignment algorithm that optimizes the calibrate the variance misalignment. The stepwise method can gradually amend the variance misalignment and produce more compact prediction sets.

To evaluate the effectiveness of APASS, we conduct extensive empirical assessments on real-world benchmarks, showing that APASS consistently improves the performance of OCP on 10 real-world datasets by 14.2% on average with 5 typical ordinal classification methods. The unstable performance of non-stepwise alignment baselines highlights the superiority of consistent improvement.

069 070

071

073

075

076 077

078

079

The contributions of this paper are summarized as follows:

- We identify the variance misalignment issue in current OC models that the existing OCP method neglects and theoretically prove that ignoring the misalignment will harm the efficiency of PSs in the context of OCP.
- We introduce the *Adaptive Posterior Alignment Step-by-Step* (APASS) method, a stepwise approach designed to reduce the PS size by reducing the distribution discrepancy using posterior variance alignment.
- We conduct extensive evaluations to show that APASS consistently improves the existing OCP method on various OC models. Specifically, the empirical results show the superiority of stepwise design to one-step baselines.
- 080 081

2 BACKGROUND

082 083 084

085

087

088

089

090

091

092

093

094

095

097

2.1 ORDINAL CLASSIFICATION.

In this study, we explore ordinal classification, which assigns labels to input instances based on a naturally ordered set of classes. We define the input space as $\mathcal{X} \subset \mathbb{R}^d$ and the ordered set of classes as $\mathcal{Y} = \{1, 2, \ldots, K\}$. The primary objective is to accurately predict the class label of input data using an OC model, denoted as $\hat{f} : \mathcal{X} \to \mathbb{R}^K$. Consider a scenario where the random variables X and Y are drawn from the combined space $\mathcal{X} \times \mathcal{Y}$ under a joint distribution $P_{X,Y}$. It is assumed that the true conditional distribution $P_{Y|X}$ is *unimodal*. This implies that for any given input instance $x \in \mathcal{X}$, the probability distribution P(Y = y | X = x) peaks at a certain class y. The prediction of our model, therefore, hinges on $\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{arg\,max}} \hat{p}_y(x)$, where $\hat{p}(y|x) = \operatorname{softmax}(\hat{f}_{\theta}(y|x))$ represents

the estimated probability that the input x corresponds to class y.

096 2.2 ORDINAL CONFORMAL PREDICTION.

098 Ordinal Conformal Prediction (OCP) leverages the output of ordinal classifiers, symbolized by $\hat{p}(x)$, 099 to construct a function $\mathcal{C}: \mathcal{X} \to 2^{\mathcal{Y}}$. This function maps input instances to a set of potential classes, 100 ensuring a specific, user-defined confidence level. As a distribution-free methodology, OCP generates 101 reliable prediction sets without making assumptions about the underlying data distribution. Formally, consider the following setup: 1) A calibration set comprising n i.i.d. data points $\{(X_i, Y_i)\}_{i=1}^n$. 102 These data points differ from the training data used to develop the ordinal classifier. 2) A new 103 test instance $X_{n+1} \in \mathcal{X}$ and a target variable $Y_{n+1} \in \mathcal{Y}$. The primary objective is to construct a 104 prediction set $C_{n,1-\alpha}(X_{n+1})$ that remains minimal yet while ensuring that it satisfies marginal 105 **coverage** at the confidence level $1 - \alpha$: 106

$$\mathbb{P}\Big(Y_{n+1} \in \mathcal{C}_{n,1-\alpha}(X_{n+1})\Big) \ge 1-\alpha,\tag{1}$$

Furthermore, the prediction set should provide **conditional coverage** at the same confidence level:

112

113

114 115 116

117 118 119

120 121

128

129

133 134

135

 $\mathbb{P}\Big(Y_{n+1} \in \mathcal{C}_{n,1-\alpha}(X_{n+1}) | X_{n+1} = x\Big) \ge 1 - \alpha, \quad \forall x \in \mathcal{X}.$ (2)

Oracle Prediction Set. In ideal settings, accessing the oracle conditional distribution, denoted as p(y|x), we construct an optimal prediction set satisfying Eq (2). It is formalized by:

$$\mathcal{C}_{1-\alpha}^{\text{oracle}}\left(x\right) = \left[l_{1-\alpha}^{\text{oracle}}\left(x\right), u_{1-\alpha}^{\text{oracle}}\left(x\right)\right],\tag{3}$$

where for any confidence level $\tau \in (0, 1]$, the boundaries of the prediction set are calculated as:

$$\left(l_{\tau}^{\text{oracle}}\left(x\right), u_{\tau}^{\text{oracle}}\left(x\right)\right) := \arg\min_{(l,u)\in\mathbb{R}^{2}:l\leq u} \left\{u-l: \sum_{j=l}^{u} p(y^{j}|x) \geq \tau\right\}.$$
(4)

Practical Solution. In practice, direct access to the Oracle conditional distribution, p(y|x), is often infeasible. Instead, the trained ordinal classifier is utilized to approximate this function, which we denote as $\hat{p}(y|x)$. Among various OCP methods that utilize $\hat{p}(y|x)$ to determine prediction sets, our research adopts the Ordinal-APS method (Lu et al., 2022). This approach integrates CP techniques to generate contiguous prediction sets using a calibrated threshold $\hat{\tau}$ to meet the desired coverage at the level of $1 - \alpha$:

$$\hat{\mathcal{C}}_{n,1-\alpha}(x) = \left[\hat{l}_{\hat{\tau}}(x), \hat{u}_{\hat{\tau}}(x) \right].$$
(5)

However, the estimated distribution $\hat{p}(y|x)$ often exhibits discrepancies, such as variance misalignment, compared to the oracle. Current OCP methods neglect these discrepancies, which can affect the efficiency of the PSs generated by the OCP method.

2.3 PROBABILITY CALIBRATION

Probability calibration, also known as confidence calibration (Guo et al., 2017), aims to ensure that
 the softmax probabilities predicted by neural networks accurately reflect the actual probabilities
 of correctness. To measure the degree of miscalibration, the Expected Calibration Error (ECE) is
 commonly used, quantifying the discrepancy between accuracy and confidence. ECE partitions the
 predictions into *M* equally spaced bins and calculates a weighted average of the difference between
 the accuracy and confidence within each bin. Formally, ECE is defined as:

142 143

$$ECE = \sum_{m=1}^{M} \frac{|B_m|}{n} \left| \operatorname{acc} \left(B_m \right) - \operatorname{conf} \left(B_m \right) \right|, \tag{6}$$

144

145 where $\operatorname{acc}(\cdot)$ and $\operatorname{conf}(\cdot)$ denotes the average accuracy and confidence in bin B_m .

In prior works, it has been generally assumed that probability calibration improves the quality of conformal prediction sets (Angelopoulos et al., 2021; Gibbs et al., 2023). However, the specific impact of calibration methods on conformal prediction remains uncertain (Xi et al., 2024). Furthermore, existing calibration techniques are not explicitly tailored to OC models, which often assume unimodal distributions. The influence of these methods on OCP is thus still an open question.

151 152

153

3 RELATED WORK

154 Ordinal Classification. Recent works on ordinal classification demonstrate substantial benefits of 155 assuming the underlying conditional distribution to be unimodal. Label smoothing methods convert 156 one-hot target labels into unimodal prior distributions to be used as the reference for the training 157 loss. SORD (Diaz & Marathe, 2019) constructs the ground-truth probability distribution using linear 158 exponentially decaying distributions based on a metric loss function $\ell(y^t, y^i)$ that penalizes the 159 distance between the actual value y^t and the *i*-th prediction value y^i . This method uses the crossentropy loss to train a neural network model. DLDL (Gao et al., 2017), on the other hand, constructs 160 the probability of the *i*-th prediction using a normal probability density function and minimizes the 161 Kullback-Leibler divergence between the predicted probability distribution and the ground-truth



Figure 1: **Overview of APASS.** APASS is a plug-and-play module that adjusts the conditional distribution as predicted by ordinal classifiers before applying the OCP approach. This versatility allows it to be integrated seamlessly with various ordinal classifiers and OCP methods. Details of **APASS-Calibration** and **APASS-Prediction** are presented in Algorithm 1 & 2.

178
179labels. R2CCP (Guha et al., 2024) proposes a loss function similar to label smoothing losses, which
penalizes the probability based on the distance between the actual value y^t and the *i*-th prediction
value y^i and uses a Shannon entropy regularizer to prevent the density estimator from collapsing to a
Dirac distribution. But these methods are often sub-optimal since the assumed priors might not reflect
the true distribution, classes might not be equispaced categories, and additionally, test predictions
might not necessarily be unimodal. Some other methods (Belharbi et al., 2019; Cardoso et al., 2023)
learn a non-parametric unimodal distribution as a constraint optimization problem in the loss function,
which is not only difficult to optimize but also does not guarantee unimodality on testing data.

186

172 173

174

175

176

177

Conformal Prediction Conformal prediction is a statistical framework characterized by a finite-sample coverage guarantee. One of the main goals of CP methods is to generate a compact prediction set. APS (Romano et al., 2020) introduces techniques aimed at achieving coverage that is similar across regions of feature space. RAPS (Angelopoulos et al., 2020) presents a regularized version of APS for Imagenet. The first work proposing the CP method for ordinal classification is Ordinal-APS (Lu et al., 2022), and another work proposes a similar approach in the context of ordinal conformal risk control (Xu et al., 2023).

The primary focal points of CP are reducing prediction set size and enhancing coverage rate. COPOC (Dey et al., 2023) proposes a special neural network model to ensure the ordinal classifier outputs an unimodal distribution and thus reduces the size of the ordinal prediction set size. In contrast, our model-agnostic method can be applied to different ordinal classifiers. Closely related to our insight, some works also utilize the information in the calibration set to generate compact prediction sets NCP (Ghosh et al., 2023) proposes to use non-parametric nearest neighbors for calibration, and in the context of sequential data, HopCPT (Auer et al., 2024) uses Modern Hopfield Networks to model the data similarity, and reweight the non-conformity score based on the similarity.

Probability Calibration Guo et al. (2017) investigates the problem of confidence calibration
in modern neural networks and finds post-processing methods like TS can effectively calibrate
predictions. Standard TS improves average calibration but reduces confidence for all predictions.
AdaTS (Joy et al., 2023) predicts a different temperature value for each input, allowing it to selectively
increase or decrease confidence as needed. Esaki et al. (2024) proposes an accuracy-preserving
calibration method using the Concrete distribution as a probabilistic model on the probability simplex,
which outperforms previous TS methods in accuracy-preserving calibration tasks.

209 210

4 Method

211

In this section, we *empirically* and *theoretically* analyze the influence of distribution discrepancies on
the efficiency of PSs generated by the OCP method: PSs will be smaller if distribution discrepancies
are smaller. Therefore, we propose a *plug-and-play* method to optimize the PS efficiency, named
Adaptive Posterior Alignment Step-by-Step (APASS). APASS first employs a variance estimator
(attention-based or kNN-based, see Section 4.2) to estimate the distribution discrepancies using data



Figure 2: Empirical Evidence for Variance Alignment. (a) Variance discrepancy between oracle and predicted distribution. (b) Distributions after variance alignment using temperature scaling

in the calibration set. Then, to make use of the estimated variance misalignment, we proposed a stepwise method that gradually adjusts the predicted posterior using temperature scaling (TS) with a small step size to ensure a monotonic decrease in PS size in Section 4.3.

233 4.1 MOTIVATION

226

227 228 229

230

231

232

266

267 268

Recent works on OC demonstrate substantial benefits of assuming the underlying conditional distribution to be unimodal for OC tasks. To encourage unimodality, label-smoothing methods convert one-hot target labels into unimodal prior distributions to be used as the reference for the training loss, and some other methods use non-parametric unimodal distribution as a constraint in the loss function. However, these methods are optimized for accuracy, not posterior distribution discrepancy.

240 To investigate the impact of distribution discrepancies between the prediction posterior distribution, $\hat{p}(y|x)$, and the oracle distribution, p(y|x). We first generated a *heteroscedastic* synthetic dataset, 241 then trained an OC model in this dataset. This allows us to access the Oracle distribution p(y|x) and 242 better mirror real-world data scenarios. Figure 2 (a) highlights the existing discrepancies: the OC 243 models tend to overestimate the variance when the actual variance is low and underestimate it when 244 the actual variance is high. Moreover, the oracle prediction set has an average size of 8.54, contrasting 245 with 9.62 for the set derived from $\hat{p}(y|x)$. This considerable difference highlights a critical variance 246 misalignment problem, which compromises the efficiency of the prediction set size. 247

Inspired by the probability calibration method, we then employ temperature scaling to align the 248 prediction distribution to the oracle distribution. Specifically, we use a grid-search strategy to 249 minimize the discrepancy between the variance of the predicted posterior distribution and the variance 250 of the oracle posterior distribution. In Figure 2 (b), we illustrate the aligned distribution. The results 251 show a more closely aligned distribution, with the average size of the prediction set reduced to 8.84. 252 These findings demonstrate that temperature scaling not only facilitates variance alignment but also 253 enhances the efficiency of the prediction set, effectively resolving issues of variance misalignment. 254 Then, we establish a theoretical explanation for our empirical finding, which explains how distribution 255 discrepancy affects the relationship between the estimated prediction set, $C_{n,1-\alpha}(X_{n+1})$, and the 256 oracle prediction set, $C_{1-\alpha}^{\text{oracle}}(X_{n+1})$. We begin by defining some necessary assumptions: 257

Assumption 1 (i.i.d. data). The data $\{(X_i, Y_i)\}_{i=1}^{n+1}$ are i.i.d. from some unknown joint distribution. Assumption 2 (Unimodality). For any $x \in \mathbb{R}^m$, the conditional distribution of Y|X = x is unimodal; i.e. there exists $y^0 \in \mathbb{R}$ (depending on x), such that $p(y^0 + y''|x) \le p(y^0 + y')$ if $y'' \ge y' \ge 0$, and $p(y^0 + y''|x) \le p(y^0 + y')$ if $y'' \le y' \le 0$.

Assumption 3 (η -inconsistency). Let F(y|x) denote the cumulative distribution function of Y|X = x, and define $\hat{F}(y|x)$ as the estimate of cumulative distribution function, i.e., $\hat{F}(y^j|x) := \sum_{i=1}^{j} \hat{p}_{\theta}(y^i|x)$. Then, we assume for all $j \in 1, ..., K$,

$$\mathbb{P}\left[\sup_{j\in\{1,\dots,K\}}\left[|\hat{F}(y^j|X) - F(y^j|X)|\right] \le \eta\right] \ge 1 - \eta.$$
(7)

Assumption 4 (Regularity). For any $x \in \mathbb{R}^m$ and $j \in \{1, 2, \dots, K\}$, $1/H < p(y^j|x) < 2/H$, for some H > 0.

270 Algorithm 1 APASS Calibration 271 **Input:** logits of calibration data $\hat{f}_{1:n} = \hat{f}_{\theta}(\boldsymbol{x}_{1:n})$ and variance estimator $\operatorname{Var}_{\psi}(Y|X, \mathcal{D}_{calib})$ 272 **Output:** aligned posterior $\hat{p}_{1:n}^{\text{APASS}}$ and the TS steps \hat{s} 273 1: Calculate $t(x_{1:n})$ by Eq 9, 10, 12 ▷ distribution discrepancy on the calibration set 274 2: $\hat{p}_{1:n} \leftarrow \operatorname{softmax}(\hat{f}_{1:n})$ 3: $|\hat{\mathcal{C}}_{n,1-\alpha}(\boldsymbol{x}_{1:n})| \leftarrow \text{OCP-Calibration}(\hat{f}_{1:n})$ ▷ run OCP-Calibration to get average PS size 276 4: $best \leftarrow |\hat{\mathcal{C}}_{n,1-\alpha}(\boldsymbol{x}_{1:n})|$ ▷ initialize best PS size 5: for $s \in \{1, ..., s_{max}\}$ do 278 $\hat{f}_{1:n} \leftarrow \hat{f}_{1:n}/t(\boldsymbol{x}_{1:n})$ \triangleright perform $TS(t(\boldsymbol{x}_{1:n}))$ 6: 279 $|\hat{\mathcal{C}}_{n,1-\alpha}(\boldsymbol{x}_{1:n})| \leftarrow \text{OCP-Calibration}(\hat{f}_{1:n})$ 7: ⊳ get new PS size 281 8: if $|\hat{\mathcal{C}}_{n,1-\alpha}(\boldsymbol{x}_{1:n})| \leq best$ then \triangleright stop until PS size does not reduce $best \leftarrow |\hat{\mathcal{C}}_{n,1-\alpha}(\boldsymbol{x}_{1:n})|$ 9: 10: else 284 11: Break 12: end if 13: end for 14: $\hat{p}_{1:n}^{\text{APASS}} \leftarrow \text{softmax}(\hat{f}_{1:n} * t(\boldsymbol{x}_{1:n})), \hat{s} \leftarrow s - 1$ 15: return $\hat{p}_{1:n}^{\text{APASS}}, \hat{s}$ 287 289

Assumption 4 allows us to quantify the distribution discrepancy using η , where a higher η value signifies a greater discrepancy. We leverage this quantification to theoretically analyze and establish an upper bound for $|\hat{C}_{n,1-\alpha}(X_{n+1})|$:

Theorem 1. For any $\alpha \in (0, 1]$, let $\hat{C}_{n,1-\alpha}(X_{n+1})$ denote the prediction set at level $1 - \alpha$ for Y_{n+1} obtained by applying OCP. The size prediction set $\hat{C}_{n,1-\alpha}(X_{n+1})$ is bounded by the set of oracle prediction set $C_{n,1-\alpha}^{oracle}(X_{n+1})$ as

$$\mathbb{P}\left[\left|\hat{\mathcal{C}}_{n,1-\alpha}(X_{n+1})\right| \le \left|\mathcal{C}_{n,1-\alpha}^{oracle}(X_{n+1})\right| + \gamma_n\right] \ge 1 - \xi_n,\tag{8}$$

300 301 302

298 299

290

where $\gamma_n = 2 + H(3/n + 2\sqrt{(\log n)/n} + 5\eta)$ and $\xi_n = \eta + 2n^{-2}$.

This theorem demonstrates that decreasing η , the measure of the discrepancy between estimated and actual distributions, results in a tighter upper bound for $|\hat{C}_{n,1-\alpha}(X_{n+1})|$. Consequently, this reduction can effectively decrease the size of the prediction set $\hat{C}_{n,1-\alpha}(X_{n+1})$. This provides theoretical guidance to find a practical way to enhance the efficiency of the PS by using probability calibration methods to reduce the distribution discrepancy. However, the straightforward application of probability calibration poses a challenge: we empirically find that probability calibration may harm the PS efficiency in the OCP setting (see Section 5.1). In the following sections, we address this challenge by introducing a novel stepwise distribution alignment method.

311

312 313

4.2 MEASURE DISTRIBUTION DISCREPANCY WITH POSTERIOR VARIANCE ESTIMATOR

To optimize the PS efficiency by reducing distribution discrepancy, we need to measure it by comparing the posterior variance predicted by the OC model $\widehat{Var}_{\theta}(Y|X)$ and the posterior variance estimated using data from the calibration set $\widehat{Var}(Y|X, \mathcal{D}_{calib})$. It is straightforward to calculate the posterior variance predicted by the OC model:

$$\widehat{\operatorname{Var}}_{\theta}(Y|X=x) = \sum_{j=1}^{K} \widehat{p}(y^j|x) \cdot \left(y^j - \sum_{j=1}^{N} \widehat{p}(y^j|x) \cdot y^j\right)^2.$$
(9)

322 Inspired from Auer et al. (2024), we estimate the posterior variance considering the calibration set 323 using a weighted sum of prediction errors derived from calibration data. This sum prioritizes points similar to the test sample and aggregates their deviations. Specifically, we define the squared residual

324 Algorithm 2 APASS Prediction 325

334 335 336

337

338 339 340

341

342

343

344

345

346

351

352

353 354 355

356 357

361

364

365

Input: logits of testing data $\hat{f}_{n+1} = \hat{f}_{\theta}(x_{n+1})$, variance estimator $\widehat{\text{Var}}_{\psi}(Y|X, \mathcal{D}_{calib})$, TS steps \hat{s} 326 **Output:** logits after APASS $\hat{f}_{n+1}^{\text{APASS}}$ 327 1: Calculate $t(x_{n+1})$ by Eq 9, 10, 12 ▷ distribution discrepancy of the testing sample 328 2: for $s \in \{1, ..., \hat{s}\}$ do $\hat{f}_{n+1} \leftarrow \hat{f}_{n+1}/t(x_{n+1}))$ 3: 330 4: end for 331 5: $\hat{f}_{n+1}^{\text{APASS}} \leftarrow \operatorname{softmax}(\hat{f}_{n+1})$ 332 6: return $\hat{f}_{n+1}^{\text{APASS}}$ 333

error in the calibration set as $\epsilon_i^2 = (y_i - \hat{y}_i)^2$. For a new sample point X = x, the posterior variance considering the calibration set is estimated by:

$$\widehat{\operatorname{Var}}_{\psi}(Y|X=x, \mathcal{D}_{calib}) = \operatorname{softmax}\left(\beta\phi^{T}(x) \boldsymbol{W}_{q}^{T} \boldsymbol{W}_{k}\phi(\boldsymbol{x}_{1:n})\right) \boldsymbol{\epsilon}_{1:n}^{2},$$
(10)

where $x_{1:n}$ are samples' features and $\epsilon_{1:n}^2$ are also squared residual errors in the calibration set, W_q and W_k are learned transformations applied before associating the query with the calibration data's keys, ϕ is an encoder, transforming raw features into appropriate representation vectors, and ψ represents all model parameters. To effectively train the variance estimator, we utilize a leave-one-out (LOO) training strategy. The primary objective in the training phase is to minimize the mean squared error between the predicted squared errors $\hat{\epsilon}_{1:n}^2$ and the actual squared errors $\epsilon_{1:n}^2$. The training loss, therefore, is formulated as follows:

$$\mathcal{L}_{\psi} = \mathsf{MSE}\left(\widehat{\mathsf{Var}}_{\psi}(Y|\boldsymbol{x}_{1:n}, \mathcal{D}_{calib}), \boldsymbol{\epsilon}_{1:n}^{2}\right)$$
(11)

During the computation of $\operatorname{Var}_{\psi}(Y|x_i, \mathcal{D}_{calib})$, the weight of ϵ_i is masked as 0 to prevent leakage of actual error values into the model training process following LOO. With the two posterior variance estimators, we can finally define the distribution discrepancy of a sample x as:

$$t(x) = \left(\frac{\widehat{\operatorname{Var}}_{\psi}(Y|X=x, \mathcal{D}_{calib})}{\widehat{\operatorname{Var}}_{\theta}(Y|X=x)}\right)^{q}$$
(12)

We then use t(x) as the temperature in TS, which we refer to as TS(t). There is no variance 358 discrepancy if t(x) = 1, which means no need to scale the distribution. If t(x) > 1, indicating the 359 OC model underestimates the posterior variance, TS will make the distribution more even; and if 360 t(x) < 1, indicating the OC model overestimates the posterior variance, TS will make the distribution more narrow. q > 0 is a hyper-parameter to determine the step size of TS: larger q means a larger TS 362 step given to same variance discrepancy.

STEPWISE POSTERIOR ALIGNMENT WITH TEMPERATURE SCALING 4.3

366 The next question is how to determine the step size q. A straightforward solution is to find the optimal 367 q^* by minimizing the ECE as other probability calibration methods do, named Adaptive Posterior 368 Alignment by Confidence Calibration (APACC). However, we find this approach may harm the 369 efficiency of PSs in practice. Therefore, inspired by gradient descent, we propose a sstepwisemethod that uses a small constant step size q but looks for the optimal steps reducing the PSs' size. 370

371 Corresponding to CP methods, APASS comprises 2 components: APASS-Calibration (Algorithm 1) 372 and APASS-Prediction (Algorithm 2), as illustrated in Figure 1. APASS-Calibration calibrates the 373 posterior distribution of calibration data predicted by the OC model using distribution discrepancy measure (Eq. 12) step by step and outputs the steps of TS \hat{s} . The aligned distributions $\hat{p}_{1:n}^{\text{APASS}}$ 374 375 are then fed to OCP-Calibration to calculate the non-conformity score and the conformal threshold $\hat{\tau}$. During test time, **APASS-Prediction** will first estimate the distribution discrepancy of a testing 376 sample x_{n+1} as $t(x_{n+1})$, then run $TS(t(x_{n+1}))$ for \hat{s} steps to get aligned posterior $\hat{p}_{n+1}^{\text{APASS}}$, which 377 OCP-Prediction will take to generate PS for the testing sample x_{n+1} .

1	Datasat	Dataset Original		o method	Stepwise method (Ours)		
2	Dalasel	Originai	AdaTS	APACC-Att	APASS-kNN	APASS-Att	
	Breastcancer	$0.031_{\pm 0.054}$	$0.028_{\pm 0.055}$	$0.03_{\pm 0.051}$	$0.033_{\pm 0.05}$	$0.033_{\pm 0.051}$	
	Community	$0.005_{\pm 0.029}$	$0.005_{\pm 0.034}$	$0.005_{\pm 0.029}$	$0.005_{\pm 0.032}$	$0.005_{\pm 0.031}$	
	Concrete	$0.012_{\pm 0.031}$	$0.012_{\pm 0.034}$	$0.013_{\pm 0.033}$	$0.011_{\pm 0.033}$	$0.013_{\pm 0.032}$	
	Diabetes	$0.025_{\pm 0.046}$	$0.025_{\pm 0.055}$	$0.023_{\pm 0.053}$	$0.023_{\pm 0.052}$	$0.026_{\pm 0.054}$	
	Energy	$0.014_{\pm 0.036}$	$0.013_{\pm 0.037}$	$0.015_{\pm 0.036}$	$0.014_{\pm 0.04}$	$0.013_{\pm 0.042}$	
	Forest	$0.011_{\pm 0.038}$	$0.011_{\pm 0.036}$	$0.011_{\pm 0.035}$	$0.012_{\pm 0.037}$	$0.012_{\pm 0.036}$	
	Parkinsons	$0.002_{\pm 0.015}$	$0.002_{\pm 0.015}$	$0.002_{\pm 0.014}$	$0.002_{\pm 0.015}$	$0.002_{\pm 0.014}$	
	Pendulum	$0.035_{\pm 0.045}$	$0.037_{\pm 0.04}$	$0.033_{\pm 0.039}$	$0.037_{\pm 0.044}$	$0.034_{\pm 0.042}$	
	Solar	$0.017_{\pm 0.055}$	$0.017_{\pm 0.054}$	$0.018_{\pm 0.048}$	$0.019_{\pm 0.047}$	$0.018_{\pm 0.04}$	
	Stock	$0.015_{\pm 0.065}$	$0.014_{\pm 0.057}$	$0.015_{\pm 0.06}$	$0.014_{\pm 0.059}$	$0.016_{\pm 0.054}$	

Table 1: Averaged results of Δ Cov comparing our method with baselines on 10 datasets and across 379 5 distinct OC m

408

425

426

427

428

378

5 EXPERIMENTS

396 This section presents the evaluation of APASS, designed to efficiently generate prediction sets with 397 specified coverage for OC tasks. We tested APASS across diverse real-world datasets and established OC models, demonstrating consistent performance improvements over all baselines. An ablation study confirmed the essential contribution of each component, enhancing APASS's robustness and 399 reducing its sensitivity to hyperparameter changes. The results affirm APASS's effectiveness and 400 practicality in real-world applications. 401

402 **OC models.** To validate the effectiveness of APASS across different base OC models, we tested 5 403 popular OC models, including 3 label-smoothing methods and 2 methods that promote unimodality 404 non-parametrically. Specifically, SORD, DLDL, R2CCP build different smoothed labels to encour-405 age unimodal. ELB (Belharbi et al., 2019) enforces unimodality and label-order consistency via a set 406 of non-parametric inequality constraints over all pairs of adjacent labels. UN (Cardoso et al., 2023) 407 proposed a new neural network architecture that directly constrains the output to be unimodal.

Baselines. We compared our models against four baselines: 1) Original: The original Ordinal-APS 409 model without adaptive calibration. 2) APACC-Attn: A one-step variant of our APASS method, 410 which uses grid search to find the optimal step size \hat{q} that minimizes ECE. This \hat{q} is then applied 411 to adjust the posterior during testing. 3) AdaTS: A state-of-the-art adaptive probability calibration 412 method (Joy et al., 2023), which is incorporated into OC models before applying CP. 4) APASS-kNN: 413 A variation that substitutes our attention-based variance measure with a distance-based alternative. 414 Specifically, we use a kNN estimator to assess posterior variance Formally, 415

$$\widehat{\operatorname{Var}}_{kNN}(Y|X=x_{n+1},\mathcal{D}_{calib}) = \frac{\sum\limits_{i\in\operatorname{NN}(x_{n+1})} w_i^k |\hat{y}_i - y_i|}{\sum\limits_{i\in\operatorname{NN}(x_{n+1})} w_i^k} + \frac{\min\limits_{i\in\operatorname{NN}(x_{n+1})} d(x_i,x_{n+1})}{\max\limits_{i,j\in\mathcal{D}_{calib}} d(x_i,x_j)} \widehat{\sigma}$$
(13)

where NN (x_{n+1}) is the set of the nearest k samples, $w_i = 1 - \frac{d(x_i, x_{n+1})}{\sum\limits_{i \in \text{NN}(x_{n+1})} d(x_i, x)}$ and $\hat{\sigma} = 0$

$$\sqrt{\operatorname{Var}\left[\{y_i\}_{i\in \operatorname{NN}(x_{n+1})} \cup \{\hat{y}_{n+1}\}\right]}$$
. In our experiment, $k = 5$ and $d(\cdot, \cdot)$ is the Mahalanobis distance.

Metrics. The metrics used in our evaluation include ΔCov , defined as the absolute difference between the actual and target coverage given a specific target coverage level $1 - \alpha$, where smaller values indicate better validity and a value of zero implies perfect alignment. Additionally, we assess the average size of the prediction set, denoted as **PS**, to evaluate the efficiency of different methods.

429 **Datasets.** We evaluate our method using 10 popular real-world datasets to ensure the robustness and 430 applicability of our approach. Specifically, these datasets include several from the 10 UCI Machine 431 Learning Repository (Kolby et al., 2024). These datasets encompass a diverse range of domains and data characteristics.

433	Table 2: Averaged results of PS comparing our method with baselines on 10 datasets and across 5
434	distinct OC models. The last raw is the average percentage reduction. See Table 4 in the Appendix
435	for complete results.

36	Dataset	Original	One-ster	p method	Stepwise me	ethod (Ours)
57			Add15	AIACC-Au	AIA55-KININ	лі доб-ди
38	Breastcancer	$ 43.5_{\pm 3.5} $	$40.6_{\pm 4.6}$ (-6.8%)	$44.0_{\pm 3.8}$ (+1.2%)	$41.3_{\pm 4.1}$ (-5.0%)	$40.0_{\pm 4.9}$ (-8.2%)
39	Community	$21.2_{\pm 1.6}$	$21.7_{\pm 1.8}$ (+2.1%)	$23.9_{\pm 2.2}$ (+12.7%)	$19.9_{\pm 1.8}$ (-6.2%)	19.1 _{±1.7} (-9.9%)
40	Concrete	$17.1_{\pm 1.7}$	$17.1_{\pm 2.0}^{-}$ (+0.2%)	$17.6_{\pm 1.7}^{-}$ (+2.7%)	$16.0_{\pm 2.0}$ (-6.6%)	$15.0_{\pm 1.6}^{-}$ (-12.4%)
41	Diabetes	$35.7_{\pm 3.6}^{-}$	$35.8_{\pm 3.6}^{-}$ (+0.0%)	$33.0_{\pm 3.7}$ (-7.8%)	$34.1_{\pm 3.3}^{-1.1}$	$33.5_{\pm 3.6}$ (-6.4%)
42	Energy	$11.0_{\pm 0.3}$	$11.1_{\pm 0.5}$ (+0.7%)	$10.8_{\pm 0.2}$ (-1.2%)	$9.1_{\pm 0.4}$ (-17.1%)	8.1 _{±0.3} (-26.5%)
43	Forest	$29.6_{\pm 2.8}$	$29.2_{\pm 3.1}$ (-1.3%)	$32.3_{\pm 3.1}$ (+9.2%)	$27.6_{\pm 2.6}$ (-6.7%)	26.0 ±3.0(-12.1%)
44	Parkinsons	$9.4_{\pm 0.1}$	$9.6_{\pm 0.6}$ (+2.7%)	$9.0_{\pm -0.6}$ (-3.9%)	$8.0_{\pm 0.4}$ (-14.0%)	6.8 _{±0.1} (-27.9%)
15	Pendulum	$10.4_{\pm 1.2}$	$10.2_{\pm 0.6}$ (-2.4%)	$11.0_{\pm 0.7}$ (+5.8%)	$9.8_{\pm 0.7}$ (-6.3%)	9.5 _{±1.0} (-9.0%)
16	Solar	$17.4_{\pm 3.1}^{-}$	$17.9_{\pm 3.2}$ (+2.8%)	$24.6_{\pm 3.6}$ (+41.2%)	$14.9_{\pm 3.6}$ (-14.6%)	$13.4_{\pm 3.9}$ (-23.2%)
+0 17	Stock	$13.2_{\pm 1.0}$	$13.6_{\pm 0.9}^{-1}$ (+3.2%)	$12.9_{\pm 0.6}$ (-1.8%)	$12.7_{\pm 0.4}$ (-4.0%)	$12.3_{\pm 1.0} \text{(-6.4\%)}$
18	Averaged Red	duction \downarrow	0.12%	5.81%	-8.52%	-14.20%
49						

Evaluation Setup. For each OCP method in one dataset, we randomly split the data into folds with 70% / 30% as $\mathcal{D}_{train}/\mathcal{D}_{calib} \cup \mathcal{D}_{test}$ for 3 times to train 3 different models. We conduct 10 random splits of calibration/testing sets for each OC model to estimate the empirical coverage and PS size. For all APASS experiments, we use the same step size q = 0.05

5.1 EMPIRICAL RESULTS

457 Effectiveness of APASS across Multiple OC Models on Real-World Datasets. We begin by 458 comparing our method with several baselines on a diverse set of real-world datasets across five 459 distinct OC models. As shown in Table 1, APASS-Att maintains a consistently low Δ Cov value, 460 indicating strong alignment between the predicted labels and the target labels. In terms of |PS|, as 461 reported in Table 2, APASS-Att achieves significant reductions. On average, it reduces the size of 462 the prediction sets by 14.20%, with reductions ranging from 6.4% to 27.9%. APASS-kNN using a 463 less competitive non-parametric variance estimator brings an 8.52% reduction in average. Moreover, APASS-Att and APASS-kNN never increase the prediction set size compared to the original models, 464 demonstrating their robustness and adaptability across various datasets and models. 465

Superiority of Stepwise vs. One-Step Approaches. As reported in Table 2, both stepwise methods, APASS-Attn and APASS-kNN, consistently outperform one-step baselines, including the state-of-the-art probability calibration method AdaTS and our one-step variant APACC-Att, which fail to reduce |PS| across all datasets. These one-step approaches fail to reduce the |PS| in 60% of cases. In contrast, the stepwise approach consistently reduces prediction set sizes, showcasing its superior ability to optimize prediction efficiency while maintaining strong model alignment. This advantage reinforces the value of the stepwise framework in real-world applications.

473 Empirical Evidence of Synchronous Changes of PS Size. One of the critical contributions 474 of this work is the stepwise alignment approach, which uses the calibration set to determine the 475 optimal TS steps. We validate this empirically by analyzing how prediction set sizes evolve during 476 stepwise temperature scaling. Figure 3 demonstrates that prediction set sizes on both calibration 477 and testing sets change synchronously across three OC models on the Community and Stock dataset. 478 This validates that the stepwise method can accurately determine the TS steps based solely on the 479 calibration set. The same synchronous behavior is consistently observed across other datasets, with 480 comprehensive results included in Appendix B. 481

Robust Performance of APASS and Computation Cost. Hyperparameter sensitivity is crucial in real-world deployment, as hyperparameter selection typically requires human expertise and can be resource-intensive. Our empirical results demonstrate that APASS is largely insensitive to hyperparameter variation. Specifically, Table 3 shows the percentage reduction in |PS| when applying APASS with various step sizes *q*. While larger step sizes (e.g., 0.5 and 1.0) lead to

432

450 451

452

453

454 455



Figure 3: The sizes of the prediction set change synchronously on calibration and testing sets.

suboptimal results, using a step size below 0.1 consistently yields near-optimal outcomes. APASS is also computationally efficient. The posterior variance estimator and calibration training should be complete before deployment, so this part of the computation cost is not crucial. The following table illustrates that our method only costs about 9.2% computation overhead with q = 0.05 (our setting). However, if the step size q is set to 0.01, the computation overhead will be 62% but only bring 0.1% extra reduction.

Table 3: Average size reduction after APASS-Att and computation cost in testing time with different step sizes q. The last raw is the result of the original Ordinal-APS.

			-					
\overline{q}	0.01	0.02	0.05	0.1	0.2	0.5	1.0	/
Averaged Reduction	-14.3%	-14.2%	-14.2%	-13.8%	-8.7%	-5.3%	-1.7%	/
Running time (s)	42.3	34.2	28.5	27.8	27.2	26.5	26.2	26.1
Overhead	62.1%	31.0%	9.2%	6.5%	4.2%	1.5%	0.4%	/

6 CONCLUSION

In this paper, we find the issue of variance misalignment in popular ordinal classifiers, which will harm OCP. We empirically and theoretically show the efficiency of OCP can be improved if ordinal classifiers predict a more accurate conditional distribution. Thus, we introduce the APASS technique, which employs an attention-based variance estimator and stepwise temperature scaling to align the posterior variance modeled by ordinal classifiers with better variance estimation. Empirical evaluations on benchmark datasets demonstrated that APASS significantly enhances the performance of OCP methods without the need for hyperparameter tuning, offering a robust framework for high-stakes healthcare, finance, and beyond applications. Limitation of our methods is we only use variance to align the posterior, high-order moment such as skewness and kurtosis can be considered in future works.

540 REFERENCES 541

542 543	Anastasios Angelopoulos, Stephen Bates, Jitendra Malik, and Michael I Jordan. Uncertainty sets for image classifiers using conformal prediction. <i>arXiv preprint arXiv:2009.14193</i> , 2020.
544 545 546	Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. <i>arXiv preprint arXiv:2107.07511</i> , 2021.
547 548 549	Anastasios Nikolas Angelopoulos, Stephen Bates, Michael Jordan, and Jitendra Malik. Uncertainty sets for image classifiers using conformal prediction. In <i>International Conference on Learning Representations, ICLR</i> , 2021.
550 551 552	Andreas Auer, Martin Gauch, Daniel Klotz, and Sepp Hochreiter. Conformal prediction for time series with modern hopfield networks. <i>Advances in Neural Information Processing Systems</i> , 2024.
553 554	Soufiane Belharbi, Ismail Ben Ayed, Luke McCaffrey, and Eric Granger. Non-parametric uni-modality constraints for deep ordinal classification. <i>arXiv preprint arXiv:1911.10720</i> , 2019.
556 557 558	Refik Can Malli, Mehmet Aygun, and Hazim Kemal Ekenel. Apparent age estimation using ensemble of deep learning models. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops</i> , 2016.
559 560 561	Jaime S Cardoso, Ricardo Cruz, and Tomé Albuquerque. Unimodal distributions for ordinal regression. arXiv preprint arXiv:2303.04547, 2023.
562 563	Prasenjit Dey, Srujana Merugu, and Sivaramakrishnan R Kaveri. Conformal prediction sets for ordinal classification. In Advances in Neural Information Processing Systems, 2023.
564 565 566	Raul Diaz and Amit Marathe. Soft labels for ordinal regression. In <i>Proceedings of the IEEE/CVF</i> Conference on Computer Vision and Pattern Recognition, 2019.
567 568 569	Yasushi Esaki, Akihiro Nakamura, Keisuke Kawano, Ryoko Tokuhisa, and Takuro Kutsuna. Accuracy- preserving calibration via statistical modeling on probability simplex. In <i>International Conference</i> <i>on Artificial Intelligence and Statistics</i> , pp. 1666–1674. PMLR, 2024.
570 571 572 573	Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In <i>International Conference on Machine Learning</i> , pp. 1050–1059. PMLR, 2016.
574 575 576	Bin-Bin Gao, Chao Xing, Chen-Wei Xie, Jianxin Wu, and Xin Geng. Deep label distribution learning with label ambiguity. <i>IEEE Transactions on Image Processing</i> , 2017.
577 578	Xin Geng. Label distribution learning. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 2016.
579 580 581 582	Subhankar Ghosh, Taha Belkhouja, Yan Yan, and Janardhan Rao Doppa. Improving uncertainty quantification of deep classifiers via neighborhood conformal prediction: Novel algorithm and theoretical analysis. In <i>Proceedings of the AAAI Conference on Artificial Intelligence, AAAI</i> , 2023.
583 584 585	Isaac Gibbs, John J Cherian, and Emmanuel J Candès. Conformal prediction with conditional guarantees. <i>arXiv preprint arXiv:2305.12616</i> , 2023.
586 587 588	Etash Kumar Guha, Shlok Natarajan, Thomas Möllenhoff, Mohammad Emtiyaz Khan, and Eugene Ndiaye. Conformal prediction via regression-as-classification. In <i>International Conference on Learning Representations, ICLR</i> , 2024.
589 590 591	Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In <i>International conference on machine learning</i> , pp. 1321–1330. PMLR, 2017.
592 593	Guodong Guo, Yun Fu, Charles R Dyer, and Thomas S Huang. Image-based human age estimation by manifold learning and locally adjusted robust regression. <i>IEEE Transactions on Image Processing</i> , 2008.

- 594 Zengwei Huo, Xu Yang, Chao Xing, Ying Zhou, Peng Hou, Jiaqi Lv, and Xin Geng. Deep age 595 distribution learning for apparent age estimation. In Proceedings of the IEEE Conference on 596 Computer Vision and Pattern Recognition Workshops, 2016. 597 Tom Joy, Francesco Pinto, Ser-Nam Lim, Philip HS Torr, and Puneet K Dokania. Sample-dependent 598 adaptive temperature scaling for improved calibration. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pp. 14919–14926, 2023. 600 601 Nottingham Kolby, Longjohn Rachel, and Kelly Markelle. Uci machine learning repository. 2024. 602 Jing Lei, Max G'Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-free 603 predictive inference for regression. Journal of the American Statistical Association, 2018. 604 605 Xiaofeng Liu, Xu Han, Yukai Qiao, Yi Ge, Site Li, and Jun Lu. Unimodal-uniform constrained wasser-606 stein training for medical diagnosis. In Proceedings of the IEEE/CVF International Conference on 607 Computer Vision Workshops, 2019. 608 Charles Lu, Anastasios N Angelopoulos, and Stuart Pomerantz. Improving trustworthiness of ai 609 disease severity rating in medical imaging with ordinal conformal prediction sets. In International 610 Conference on Medical Image Computing and Computer-Assisted Intervention. Springer, 2022. 611 612 Georgios Manthoulis, Michalis Doumpos, Constantin Zopounidis, and Emilios Galariotis. An ordinal 613 classification framework for bank failure prediction: Methodology and empirical evidence for us banks. European Journal of Operational Research, 2020. 614 615 Yaniv Romano, Matteo Sesia, and Emmanuel Candes. Classification with valid and adaptive coverage. 616 Advances in Neural Information Processing Systems, 2020. 617 618 Matteo Sesia and Yaniv Romano. Conformal prediction using conditional histograms. In Advances 619 in Neural Information Processing Systems, 2021. 620 Ralph C Smith. Uncertainty quantification: theory, implementation, and applications, volume 12. 621 Siam, 2013. 622 623 Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. Algorithmic learning in a random world. 624 Springer, 2005. 625 Volodya Vovk, Alexander Gammerman, and Craig Saunders. Machine-learning applications of 626 algorithmic randomness. 1999. 627 Xin Wen, Biying Li, Haiyun Guo, Zhiwei Liu, Guosheng Hu, Ming Tang, and Jingiao Wang. Adaptive 628 variance based label distribution learning for facial age estimation. In European Conference on 629 Computer Vision. Springer, 2020. 630 631 Huajun Xi, Jianguo Huang, Lei Feng, and Hongxin Wei. Does confidence calibration help conformal 632 prediction? arXiv preprint arXiv:2402.04344, 2024. 633 Yunpeng Xu, Wenge Guo, and Zhi Wei. Conformal risk control for ordinal classification. In The 634 Conference on Uncertainty in Artificial Intelligence. PMLR, 2023. 635 636 637 638 639 640 641 642 643 644 645 646
- 647

PROOF OF THEORY А

The theoretical proof strikes the idea from Sesia & Romano (2021), which focuses on regression problems with continuous variables, whereas we concentrate on ordinal classification with discrete variables.

Lemma 1. Def the event A as

$$\mathcal{A} := \left\{ x : \sup_{j \in \{1, \dots, K\}} |\hat{F}(y^j | x) - F(y^j | x)| > \eta \right\}$$
(14)

Then, under Assumptions 1-3, for any $X \perp D^{train}$,

$$\mathbb{P}\left[X \in \mathcal{A}\right] \le \eta \tag{15}$$

$$\mathcal{D}^{\text{cal},a} := \{ i \in \{1, \dots, n\} : X_i \in \mathcal{A} \}, \quad \mathcal{D}^{\text{cal},b} := \{ i \in \{1, \dots, n\} : X_i \in \mathcal{A}^c \}$$
(16)

we have that, for any constant c > 0

$$\mathbb{P}\left[|\mathcal{D}^{cal,a}| \ge n\eta + c\sqrt{n\log n}\right] \le n^{-2c^2} \tag{17}$$

Lemma 2. Under Assumptions 1–3, for any $\tau \in (0, 1)$ and $X \perp D^{train}$,

$$\mathbb{P}\left[\left|\hat{\mathcal{C}}_{n,\tau}(X)\right| \le \left|\mathcal{C}_{n,\tau+2\eta}(X)\right| + 2\right] \ge 1 - \eta,\tag{18}$$

Lemma 3. For any $\tau \in (0,1)$, let $\hat{Q}_{\tau}(E_i)$ denote the $\lceil \tau(n+1) \rceil$ smallest value among the conformity scores $\{E_i\}$ for $i \in \mathcal{D}^{cal}$, where $i \in \mathcal{D}^{cal}$ and

$$E_i := \min\left\{\tau_t \in \{0, 1/T_n, \dots, (T_n - 1)/T_n, 1\} : Y_i \in \hat{\mathcal{C}}_{n, \tau_t}(X_i)\right\}$$
(19)

Then, under Assumptions 1-3, for any c > 0,

$$\mathbb{P}\left[\hat{Q}_{\tau}(E_i) \le \tau + \epsilon_n\right] \ge 1 - 2n^{-2c^2},\tag{20}$$

where $\epsilon_n := 3/n + 3\eta + 2c\sqrt{(\log n)/n}$

> *Proof of Theorem 1.* Define $\epsilon_n := 3/n + 3\eta + 2c\sqrt{(\log n)/n}$ for any c > 0, as in Lemma 3. In the event that $\hat{Q}_{1-\alpha}(E_i) \leq 1 - \alpha + \epsilon_n$,

$$\mathbb{P}\left[|\hat{\mathcal{C}}_{n,\hat{Q}_{1-\alpha}(E_i)}(X)| \leq |\mathcal{C}_{n,1-\alpha+\epsilon_n+2\eta}(X)| + 2\right] \\
\geq \mathbb{P}\left[|\hat{\mathcal{C}}_{n,1-\alpha+\epsilon_n}(X)| \leq |\mathcal{C}_{n,1-\alpha+\epsilon_n+2\eta}(X)| + 2\right] \\
\geq 1 - \eta,$$
(21)

where the second inequality follows by applying Lemma 2 with $\tau = 1 - \alpha + \epsilon_n$. Further, as Lemma 3 tells us, the above event occurs with a high probability,

$$\mathbb{P}\left[\hat{Q}_{1-\alpha}(E_i) \le 1 - \alpha + \epsilon_n\right] \ge 1 - 2n^{-2c^2}$$
(22)

in general, we have that

$$\mathbb{P}\left[\left|\hat{\mathcal{C}}_{n,\hat{Q}_{1-\alpha}(E_i)}(X)\right| \le \left|\mathcal{C}_{n,1-\alpha+\epsilon_n+2\eta}(X)\right| + 2\right] \ge 1 - \eta - 2n^{-2c^2}$$
(23)

By Assumption 4, $p(y^j|x) > 1/H$ for all $j \in \{1, 2, ..., K\}$. This implies $\mathcal{C}_{n,\tau}(X)$ is H-Lipschitz as a function of τ . Therefore,

703
704
$$\mathbb{P}\left[|\hat{\mathcal{C}}_{n,\hat{Q}_{1-\alpha}(E_i)}(X)| \le |\mathcal{C}_{n,1-\alpha}(X)| + 2 + H(\epsilon_n + 2\eta)\right]$$

$$\geq \mathbb{P}\left[|\hat{\mathcal{C}}_{n,\hat{Q}_{1-\alpha}(E_i)}(X)| \leq |\mathcal{C}_{n,1-\alpha+\epsilon_n+2\eta}(X)|+2\right]$$

$$\geq 1-\eta-2n^{-2c^2}.$$

Hence, setting c = 1 we have proved that

$$\mathbb{P}\left[|\hat{\mathcal{C}}_{n,\hat{Q}_{1-\alpha}(E_i)}(X)| \le |\mathcal{C}_{n,1-\alpha}(X)| + \gamma_n\right] \ge 1 - \xi_n.$$

$$(25)$$

Proof of Lemma 1. Inequation 15 and easily derived from definition of A in Eq. 14 and Assumption 3. As we know from the above that $\mathbb{P}[X \in A_n] \leq \eta$, for any $\epsilon > 0$, following Hoeffding's inequality,

$$\mathbb{P}\left[\left|\mathcal{D}^{\operatorname{cal},a}\right| \ge n\eta + \epsilon\right] \le \mathbb{P}\left[\left|\mathcal{D}^{\operatorname{cal},a}\right| \ge n\mathbb{P}\left[X \in A_n\right] + \epsilon\right]$$
$$\le \mathbb{P}\left[\frac{1}{n}\sum_{i=n+1}^{2n}\mathbb{1}\left[X_i \in A_n\right] \ge \mathbb{P}\left[X_i \in A_n\right] + \frac{\epsilon}{n}\right]$$
$$\le \exp\left(-\frac{2\epsilon^2}{n}\right).$$
(26)

Therefore, setting $\epsilon = c\sqrt{n \log n}$, for some constant c > 0, yields

$$\mathbb{P}\left[\left|\mathcal{D}^{\operatorname{cal},a}\right| \ge n\eta + c\sqrt{n\log n}\right] \le n^{-2c^2}.$$
(27)

(24)

730 Proof of Lemma 2. Consider the event \mathcal{A} defined in Lemma 1. Let's consider the case where $X \in \mathcal{A}^{c}$. 731 We can write $\hat{\mathcal{C}}_{n,\tau}(X) = [\hat{j}_1, \hat{j}_2]$ for some $\hat{j}_1, \hat{j}_2 \in \{1, \dots, K\}$ such that $\hat{F}\left(y^{\hat{j}_2}\right) - \hat{F}\left(y^{\hat{j}_1-1}\right) \ge \tau$. 732 Then, the triangle inequality implies $F\left(y^{\hat{j}_2}\right) - F\left(y^{\hat{j}_1-1}\right) \ge \tau - 2\eta$. Consider the oracle set 734 $\mathcal{C}_{n,\tau+2\eta}(X)$, which we can write in short as $[l^*, u^*]$ for some $l^*, u^* \in \mathbb{R}$ such that $F(u^*) - F(l^*) \ge \tau + 2\eta$. Define $j'_1, j'_2 \in \{1, \dots, K\}$ as the indices of the label immediately below and above l^*, u^* :

$$j'_{1} := \max \left\{ j \in \{1, \dots, m_{n}\} : y^{j} < l^{*} \right\}$$

$$j'_{2} := \min \left\{ j \in \{1, \dots, m_{n}\} : y^{j} > u^{*} \right\}$$
(28)

739 This definition implies

$$y^{j'_2} - y^{j'_1} \le u^* - l^* + 2, \tag{29}$$

742 Furthermore,

$$\hat{F}\left(y^{j_{2}'}\right) - \hat{F}\left(y^{j_{1}'}\right) \ge \hat{F}\left(u^{*}\right) - \hat{F}\left(l^{*}\right)
\ge F\left(u^{*}\right) - F\left(l^{*}\right) - 2\eta
> \tau.$$
(30)

The result implies that $\hat{j}_2 - \hat{j}_1 \leq j'_2 - j'_1$ because $\hat{j}_2 - \hat{j}_1$ is the minimal of $\hat{\mathcal{C}}_{n,\tau}(X)$. Then,

$$\begin{aligned} \hat{\mathcal{C}}_{n,\tau}(X) &= y^{\hat{j}_2} - y^{\hat{j}_1} \le y^{j'_2} - y^{j'_1} \\ &\le |\mathcal{C}_{n,\tau+2\eta}(X)| + 2 \end{aligned}$$
(31)

if $X \in \mathcal{A}^{c}$. Finally, by applying Lemma 1,

753
754
755
$$\mathbb{P}\left[\left|\hat{\mathcal{C}}_{n,\tau}(X)\right| \le |\mathcal{C}_{n,\tau+2\eta}(X)| + 2\right] = \mathbb{P}[X \in \mathcal{A}^{c}] \ge 1 - \eta$$
(32)

Proof of Lemma 3. Take any $i \in \mathcal{D}^{\operatorname{cal},b}$, where $\mathcal{D}^{\operatorname{cal},b}$ is defined as in Lemma 1:

$$\mathcal{D}^{\text{cal},b} := \{ i \in \{1, \dots, n\} : X_i \in A^c \},$$
(33)

For any fixed $t \in \{0, ..., n\}$ and $\tau_t = t/n$, omitting the explicit dependence on X and \hat{p} , we can write $\hat{\mathcal{C}}_{n,\tau_t}(X) = \begin{bmatrix} \hat{j}_1, \hat{j}_2 \end{bmatrix}$, for some $\hat{j}_1, \hat{j}_2 \in \{1, ..., K\}$ such that $\hat{F}(y^{\hat{j}_2}) - \hat{F}(y^{\hat{j}_1-1}) \ge \tau_t$. Then 762

$$\mathbb{P}[E_{i} \leq \tau_{t}] = \mathbb{P}\left[Y_{i} \in \hat{\mathcal{C}}_{n,\tau_{t}}(X)\right]$$

= $F(y^{\hat{j}_{2}}) - F(y^{\hat{j}_{1}-1})$
 $\geq \hat{F}(y^{\hat{j}_{2}}) - \hat{F}(y^{\hat{j}_{1}-1}) - 2\eta$
 $\geq \tau_{t} - 2\eta.$ (34)

Above, the first inequality follows from the definition of $\mathcal{D}^{\operatorname{cal},b}$. Equivalently, we can rewrite this as

$$\mathbb{P}\left[E_i > \tau_t + 2\eta + \delta\right] \le 1 - \tau_t - \delta,\tag{35}$$

for any $\delta > 0$. Now, partition $\mathcal{D}^{\text{cal},b}$ into the following two disjoint subsets:

$$\mathcal{D}^{\operatorname{cal},b1} := \left\{ i \in \mathcal{D}^{\operatorname{cal},b} : E_i \le \tau_t + 2\eta + \delta \right\}$$
$$\mathcal{D}^{\operatorname{cal},b2} := \left\{ i \in \mathcal{D}^{\operatorname{cal},b} : E_i > \tau_t + 2\eta + \delta \right\}$$
(36)

We bound $|\mathcal{D}^{\operatorname{cal},b2}|$ with Hoeffding's inequality. For any $i \in \mathcal{D}^{\operatorname{cal}}$, define $\tilde{E}_i = E_i$ if $i \in \mathcal{D}^{\operatorname{cal},b}$ and $E_i = \tau_t$ otherwise. For any $\epsilon > 0$,

$$\mathbb{P}\left[|\mathcal{D}^{\mathrm{cal},b^{2}}| \geq n(1-\tau_{t}-\delta)+\epsilon\right] \\
\leq \mathbb{P}\left[\frac{1}{n}\sum_{i\in\mathcal{D}^{\mathrm{cal},b}}\mathbbm{1}\left[\tilde{E}_{i}>\tau_{t}+2\eta+\delta\right]\geq\mathbb{P}\left[E_{i}>\tau_{t}+2\eta+\delta\right]+\frac{\epsilon}{n}\right] \\
= \mathbb{P}\left[\frac{1}{n}\sum_{i=1}^{n}\mathbbm{1}\left[\tilde{E}_{i}>\tau_{t}+2\eta+\delta\right]\geq\mathbb{P}\left[E_{i}>\tau_{t}+2\eta+\delta\right]+\frac{\epsilon}{n}\right] \\
\leq \mathbb{P}\left[\frac{1}{n}\sum_{i=1}^{n}\mathbbm{1}\left[\tilde{E}_{i}>\tau_{t}+2\eta+\delta\right]\geq\mathbb{P}\left[\tilde{E}_{i}>\tau_{t}+2\eta+\delta\right]+\frac{\epsilon}{n}\right] \\
\leq \exp\left(-\frac{2\epsilon^{2}}{n}\right)$$
(37)

Therefore, setting $\epsilon = c\sqrt{n \log n}$, for some constant c > 0, yields

$$\mathbb{P}\left[|\mathcal{D}^{\operatorname{cal},b^2}| \ge n(1-\tau_t-\delta) + c\sqrt{n\log n}\right] \le n^{-2c^2}$$
(38)

As $|\mathcal{D}^{\operatorname{cal},b1}| = n - |\mathcal{D}^{\operatorname{cal},a}| - |\mathcal{D}^{\operatorname{cal},b2}|$, combining the above result with that of Lemma 1 yields:

$$\mathbb{P}\left[|\mathcal{D}^{\operatorname{cal},b1}| \ge n\tau_t + n\delta - n\eta - 2c\sqrt{n\log n}\right] \ge 1 - 2n^{-2c^2}$$
(39)

If we choose $\delta = \tau_t/n + \eta + 2c\sqrt{(\log n)/n}$

$$\mathbb{P}\left[|\mathcal{D}^{\operatorname{cal},b1}| \ge \tau_t(n+1)\right] \ge 1 - 2n^{-2c^2},\tag{40}$$

which means

$$\mathbb{P}\left[\hat{Q}_{\tau_t}(E_i) \le \tau_t + \tau_t/n + 3\eta + 2c\sqrt{(\log n)/n}\right] \ge 1 - 2n^{-2c^2}$$
(41)

 $> 1 - 2n^{-2c^2}$.

Now, consider any continuous $\tau \in (0, 1]$, and $t' = \min \{t \in \{0, \dots, T_n\} : \tau_t \ge \tau\}$. As $\tau_{t'} \ge \tau$, we know $\hat{Q}_{\tau_{t'}}(E_i) \ge \hat{Q}_{\tau}(E_i)$. Therefore,

$$\mathbb{P}\left[\hat{Q}_{\tau}(E_i) \leq \tau_{t'} + \tau_{t'}/n + 3\eta + 2c\sqrt{(\log n)/n}\right]$$

$$\geq \mathbb{P}\left[\hat{Q}_{\tau_{t'}}(E_i) \leq \tau_{t'} + \tau_{t'}/n + 3\eta + 2c\sqrt{(\log n)/n}\right]$$
(42)

816 817 818

819

825

826 827

828 829

830

840

842

849 850

854 855 856

813 814 815

As $\tau_{t'} = \tau + 1$. Therefore,

$$\mathbb{P}\left[\hat{Q}_{\tau}(E_i) \le \tau + 1/n + \tau/n + 1/n^2 + 3\eta + 2c\sqrt{(\log n)/n}\right] \ge 1 - 2n^{-2c^2}.$$
(43)

Finally, as $\tau \le 1$ and $n \ge 1$, replacing $1/n + \tau/n + 1/n^2$ with 3/n will preserve the inequality and Lemma 3 is proved.

B EXPERIMENT

B.1 SYNTHETIC DATASET

831 We employed a method involving multivariate normal distributions and linear combinations to 832 generate synthetic data for our experiment. Initially, we created a random mean vector and a 833 symmetric positive-definite covariance matrix to define the multivariate normal distribution. This 834 distribution is used to generate a dataset of features. We computed the output means and variances 835 by applying linear combinations of the generated features with randomly generated coefficients to 836 produce output values. Precisely, the means are calculated as a linear combination of the features, while the variances are determined by squaring another linear combination of these features, ensuring 837 non-negativity. Finally, output values are sampled from a normal distribution using the calculated 838 means and variances, resulting in a comprehensive synthetic dataset for experimental analysis. 839

841 B.2 ORDINAL CLASSIFIER

Vanilla cross-entropy loss ignores the ordinal relationship and non-uniform separation among labels. To learn better conditional mass function approximation $\hat{p}(y|x)$, many label smoothing methods convert one-hot target labels into unimodal prior distributions to be used as the reference for the training loss. Soft ordinal classification (SORD) constructs the ground-truth p.m.f. based on a metric loss function $\ell(y^t, y^i)$ that penalizes how far the true value y^t is from the i-th prediction value y^i (Diaz & Marathe, 2019):

$$p(y^i) = \frac{e^{-\ell(y^t, y^i)}}{\sum_{k=1}^{K} e^{-\ell(y^t, y^k)}} \quad \forall y^i \in \mathcal{Y},$$

$$(44)$$

and use cross-entropy loss to train the neural network model. Deep Label Distribution Learning
 (DLDL) minimizes the KL divergence between the predicted probability and the ground-truth labels
 in a similar way:

$$p(y^{i}) = \frac{\mathcal{K}(y^{i}|\mu,\sigma)}{\sum_{k=1}^{K} \mathcal{K}(y^{k}|\mu,\sigma)} \quad \forall y^{i} \in \mathcal{Y}.$$
(45)

where $\mathcal{K}(y^i|\mu,\sigma)$ is a normal p.d.f. The mean μ is set to the actual value, i.e., $\mu = y^t$, and σ is usually determined by the data distribution. For the continuous label in the regression setting, Regressionto-Classification Conformal Prediction (R2CCP) converts regression into ordinal classification by discretizing the label space into K bins. They proposed a loss similar to label smoothing losses with a Shannon entropy regularizer, which prevents the density estimator from collapsing to one-hot output (Guha et al., 2024):

$$\mathcal{L}(\theta) = \sum_{k=1}^{K} \ell(y^t, y^k) \hat{p}_{\theta}(y^k | x) - \tau \mathcal{H}(\hat{p}_{\theta}(\cdot | x)).$$
(46)

Dataset	OC	Original	One-step	v method	Stepwise method (Ours)		
Duiusei	Model	Originai	AdaTS	APACC-Att	APASS-kNN	APASS-Att	
	DLDL	$ 42.53_{\pm 4.36}$	$45.3_{\pm 4.62}$	$43.13_{\pm 2.5}$	$39.35_{\pm 2.97}$	$36.35_{\pm 3.99}$	
	ELB	$41.9_{\pm 3.74}$	$30.17_{\pm 6.02}$	$42.48_{\pm 4.95}$	$41.83_{\pm 4.12}$	41.7 $_{\pm 5.46}$	
Breastcancer	R2CCP	$43.0_{\pm 2.05}$	$49.61_{\pm 4.76}$	$43.95_{\pm 3.95}$	$41.03_{\pm 4.68}$	$39.96_{\pm 5.49}$	
	SORD	$43.77_{\pm 3.72}$	$28.93_{\pm 2.73}$	$43.38_{\pm 2.47}$	$40.08_{\pm 4.31}$	$39.05_{\pm 5.09}$	
	UN	$46.4_{\pm 3.85}$	$48.82_{\pm 4.69}$	$47.21_{\pm 5.18}$	$44.41_{\pm 4.31}$	$42.8_{\pm 4.59}$	
	DLDL	$18.77_{\pm 1.45}$	$19.54_{\pm 2.11}$	$21.79_{\pm 2.73}$	$18.51_{\pm 1.23}$	$\textbf{18.4}_{\pm 1.54}$	
	ELB	$16.4_{\pm 1.56}$	$17.94_{\pm 1.14}$	$17.83_{\pm 2.1}$	$16.22_{\pm 2.06}$	$16.0_{\pm 1.93}$	
Community	R2CCP	$17.98_{\pm 1.26}$	$18.73_{\pm 1.97}$	$22.81_{\pm 2.41}$	$17.55_{\pm 1.62}$	$17.35_{\pm 1.28}$	
	SORD	$20.9_{\pm 1.85}$	$18.33_{\pm 1.72}$	$20.98_{\pm 1.46}$	$19.6_{\pm 2.48}$	$19.27_{\pm 2.59}$	
	UN	$32.1_{\pm 1.83}$	$33.85_{\pm 2.25}$	$36.21_{\pm 2.22}$	$27.7_{\pm 1.57}$	$24.6_{\pm 1.38}$	
	DLDL	$10.31_{\pm 1.96}$	$9.26_{\pm 1.5}$	$10.66_{\pm -0.07}$	$10.19_{\pm 1.14}$	$10.14_{\pm 1.92}$	
	ELB	$15.4_{\pm 1.75}$	$10.0_{\pm 2.38}$	$15.74_{\pm 3.55}$	$15.35_{\pm 3.41}$	$15.2_{\pm 1.73}$	
Concrete	R2CCP	$9.86_{\pm 1.3}$	$9.72_{\pm 0.49}$	$10.12_{\pm 1.18}$	$9.7_{\pm 1.79}$	$9.65_{\pm 1.17}$	
	SORD	$11.35_{\pm 1.65}$	$13.38_{\pm 2.44}$	$11.42_{\pm 2.61}$	$11.31_{\pm 1.35}$	$11.27_{\pm 1.7}$	
	UN	$38.6_{\pm 1.89}$	$43.29_{\pm 3.38}$	$39.89_{\pm 1.22}$	$33.36_{\pm 2.12}$	$28.7_{\pm 1.39}$	
	DLDL	$35.05_{\pm 3.28}$	$40.13_{\pm 2.18}$	$34.19_{\pm 3.02}$	$33.21_{\pm 3.43}$	$32.81_{\pm 3.58}$	
	ELB	$33.3_{\pm 3.89}$	$37.04_{\pm 1.71}$	$31.22_{\pm 2.99}$	$33.18_{\pm 2.68}$	$33.1_{\pm 3.56}$	
Diabetes	R2CCP	$38.16_{\pm 4.42}$	$39.42_{\pm 4.97}$	$32.5_{\pm 3.84}$	$33.76_{\pm 3.28}$	$32.21_{\pm 3.65}$	
	SORD	$34.89_{\pm 2.99}$	$22.7_{\pm 2.39}$	$35.03_{\pm 4.04}$	$34.03_{\pm 4.32}$	$33.48_{\pm 3.5}$	
	UN	$37.3_{\pm 3.68}$	$39.45_{\pm 6.63}$	$31.91_{\pm 4.7}$	$36.18_{\pm 2.92}$	$35.7_{\pm 3.52}$	
	DLDL	$2.88_{\pm 0.23}$	$2.5_{\pm -0.57}$	$2.9_{\pm -1.36}$	$2.85_{\pm -1.23}$	$2.83_{\pm 0.27}$	
	ELB	$8.38_{\pm 0.26}$	$8.95_{\pm 0.98}$	$8.17_{\pm 0.79}$	$8.01_{\pm 1.62}$	7.76 $_{\pm 0.28}$	
Energy	R2CCP	$3.07_{\pm 0.26}$	$2.14_{\pm 2.37}$	$2.96_{\pm 1.93}$	$3.03_{\pm -0.11}$	$3.02_{\pm 0.29}$	
	SORD	$2.95_{\pm 0.31}$	$3.33_{\pm 0.15}$	$3.06_{\pm 0.31}$	$2.94_{\pm 1.04}$	$2.93_{\pm 0.3}$	
	UN	$37.6_{\pm 0.24}$	$38.36_{\pm -0.48}$	$37.16_{\pm -0.68}$	$28.66_{\pm 0.51}$	$23.8_{\pm 0.28}$	
	DLDL	$31.18_{\pm 2.45}$	$23.25_{\pm 1.67}$	$38.13_{\pm 2.89}$	$30.29_{\pm 2.35}$	$29.55_{\pm 2.51}$	
	ELB	$27.0_{\pm 3.16}$	$27.31_{\pm 3.75}$	$26.87_{\pm 3.08}$	$26.45_{\pm 2.38}$	$26.2_{\pm 2.57}$	
Forest	R2CCP	$30.91_{\pm 1.45}$	$35.76_{\pm 1.48}$	$35.83_{\pm 1.88}$	$28.85_{\pm 1.85}$	$27.79_{\pm 2.5}$	
	SORD	$33.57_{\pm 3.72}$	$36.56_{\pm 3.61}$	$33.22_{\pm 4.81}$	$31.02_{\pm 3.02}$	$29.23_{\pm 4.45}$	
	UN	$ 25.4_{\pm 3.0} $	$23.32_{\pm 4.86}$	$27.64_{\pm 2.88}$	$21.47_{\pm 3.59}$	$17.4_{\pm 2.99}$	
	DLDL	$2.1_{\pm 0.08}$	$1.56_{\pm 0.92}$	$2.0_{\pm -0.85}$	$2.09_{\pm 0.18}$	$2.09_{\pm 0.08}$	
	ELB	$6.91_{\pm 0.06}$	$7.59_{\pm 0.37}$	$6.78_{\pm -1.11}$	$6.69_{\pm -0.21}$	$6.8_{\pm 0.07}$	
Parkinsons	R2CCP	$2.12_{\pm 0.04}$	$2.22_{\pm 1.29}$	$2.05_{\pm -0.5}$	$2.05_{\pm 0.21}$	$2.01_{\pm 0.04}$	
	SORD	$2.04_{\pm 0.04}$	$1.83_{\pm 0.82}$	$2.03_{\pm 0.91}$	$2.04_{\pm 1.69}$	$2.04_{\pm 0.04}$	
	UN	33.0 ± 0.06	34.85 ± -0.61	32.09 ± -1.23	27.37 ± -0.03	$20.8_{\pm 0.07}$	
	DLDL	$7.15_{\pm 1.12}$	$6.04_{\pm 0.39}$	$6.68_{\pm 1.02}$	$6.77_{\pm 1.82}$	$6.38_{\pm 1.06}$	
~	ELB	$12.9_{\pm 1.27}$	$14.51_{\pm 0.71}$	$12.39_{\pm 0.63}$	$12.77_{\pm 1.39}$	$12.7_{\pm 0.8}$	
Pendulum	R2CCP	$6.67_{\pm 0.95}$	$7.15_{\pm 0.76}$	$9.06_{\pm 0.45}$	$6.21_{\pm -1.02}$	$5.86_{\pm 0.77}$	
	SORD	$9.23_{\pm 1.35}$	$10.27_{\pm 1.43}$	$9.43_{\pm -0.05}$	$8.37_{\pm -0.75}$	$8.08_{\pm 1.22}$	
	UN	10.3 ± 1.3	13.04 ± -0.19	17.71 ± 1.23	$14.83_{\pm 2.11}$	14.5 $_{\pm 1.22}$	
	DLDL	$22.22_{\pm 2.81}$	$24.47_{\pm 2.42}$	$19.4_{\pm 2.36}$	$20.97_{\pm 2.68}$	$20.03_{\pm 3.17}$	
	ELB	$25.1_{\pm 3.91}$	$21.17_{\pm 2.86}$	$26.02_{\pm 3.65}$	$21.54_{\pm 3.42}$	$18.3_{\pm 3.55}$	
Solar	R2CCP	$6.22_{\pm 0.78}$	$6.51_{\pm -0.95}$	$26.99_{\pm 1.19}$	$3.66_{\pm 1.09}$	$1.8_{\pm 0.68}$	
	SORD	$21.98_{\pm 6.23}$	$25.38_{\pm 5.69}$	$21.24_{\pm 4.59}$	$18.22_{\pm 6.3}$	$17.53_{\pm 6.08}$	
	UN	$ 11.6_{\pm 1.86}$	$12.03_{\pm 5.88}$	$29.34_{\pm 6.03}$	$9.98_{\pm 4.75}$	9.29 ±5.96	
	DLDL	$9.14_{\pm 1.04}$	$9.83_{\pm -0.36}$	$9.25_{\pm 1.25}$	$8.65_{\pm -0.21}$	$\textbf{8.37}_{\pm 0.95}$	
	ELB	$11.4_{\pm 0.91}$	$10.57_{\pm 0.79}$	$11.07_{\pm 0.11}$	$11.34_{\pm -0.42}$	$11.3_{\pm 1.03}$	
Stock	R2CCP	$9.74_{\pm 1.15}$	$7.88_{\pm 1.14}$	$9.42_{\pm 1.68}$	$9.01_{\pm 2.42}$	$8.86_{\pm 1.25}$	
	SORD	$9.73_{\pm 0.83}$	$9.32_{\pm 1.03}$	$9.83_{\pm -1.55}$	$9.34_{\pm -0.45}$	$9.25_{\pm 0.98}$	
	UN	$ 25.9 \pm 0.87$	$30.41_{\pm 1.77}$	$25.14_{\pm 1}$ 25	24.94 ± 0.55	$23.9_{\pm 0.98}$	

B.3 ADDITIONAL EXPERIMENT RESULTS

 In this section, we provide all experiment results of how the sizes of the prediction set change in the calibration and testing sets as we do stepwise posterior alignment.



Figure 4: The prediction set size change on Breastcancer, Community, Concrete, and Diabetes



Figure 5: The prediction set size change on Energy, Forest, Parkinsons, and Pendulum



Figure 6: The prediction set size change on Solar and Stock