# MBCT: Tree-Based Feature-Aware Binning for Individual Uncertainty Calibration

Siguang Huang[1], Yunli Wang[1], Lili Mou[2], Huayue Zhang[1], Han Zhu[1], Chuan Yu[1], Bo Zheng[1]

[1]Alibaba Group, China

[2]Dept. Computing Science, Alberta Machine Intelligence Institute (Amii), University of Alberta, Canada

{siguang.hsg, ruoyu.wyl, huayue.zhy, zhuhan.zh, yuchuan.yc, bozheng}@alibaba-inc.com

doublepower.mou@gmail.com

## ABSTRACT

Most machine learning classifiers only concern classification accuracy, while certain applications (such as medical diagnosis, meteorological forecasting, and computation advertising) require the model to predict the true probability, known as a calibrated estimate. In previous work, researchers have developed several calibration methods to post-process the outputs of a predictor to obtain calibrated values, such as binning and scaling methods. Compared with scaling, binning methods are shown to have distribution-free theoretical guarantees, which motivates us to prefer binning methods for calibration. However, we notice that existing binning methods have several drawbacks: (a) the binning scheme only considers the original prediction values, thus limiting the calibration performance; and (b) the binning approach is non-individual, mapping multiple samples in a bin to the same value, and thus is not suitable for order-sensitive applications. In this paper, we propose a feature-aware binning framework, called Multiple Boosting Calibration Trees (MBCT), along with a multi-view calibration loss to tackle the above issues. Our MBCT optimizes the binning scheme by the tree structures of features, and adopts a linear function in a tree node to achieve individual calibration. Our MBCT is non-monotonic, and has the potential to improve order accuracy, due to its learnable binning scheme and the individual calibration. We conduct comprehensive experiments on three datasets in different fields. Results show that our method outperforms all competing models in terms of both calibration error and order accuracy. We also conduct simulation experiments, justifying that the proposed multi-view calibration loss is a better metric in modeling calibration error. In addition, our approach is deployed in a real-world online advertising platform; an A/B test over two weeks further demonstrates the effectiveness and great business value of our approach.

## CCS CONCEPTS

• **Computing methodologies → Machine learning**.

## KEYWORDS

Uncertainty calibration, Feature-aware binning, Multiple boosting calibration trees, Multi-view calibration

## 1 INTRODUCTION

In recent years, machine learning is widely used for classification in various fields. A large number of tasks not only concern the classification accuracy, but also expect that the output of the machine learning model reflects the true frequency of an event, such as medical diagnosis [10], meteorological forecasting [3, 5, 7], and computation advertising [21, 22]. Take computation advertising as an example. If the click-through rate on 1000 page views is 0.05, we expect approximately 50 clicks. Unfortunately, the prediction of most machine learning models (e.g., logistic regression and deep neural networks) usually deviates from the true probability, either overestimated or underestimated. This in turn leads to the lack of robustness and interpretability for these models, limiting their application in safe-critical tasks.

To mitigate the estimation deviation, researchers have developed various methods to better estimate the true probability. Several studies directly learn a calibrated predictor [26], whereas others adopt a post-hoc calibrator for the predictor [13, 14, 25]. The post-hoc methods, which are commonly known as *calibration*, are widely used, as they can be applied to off-the-shelf predictors. Existing post-hoc calibration can be further divided into three groups: parametric (e.g., scaling), non-parametric (e.g., binning), and their hybrid. Scaling methods adopt a parametric function with a certain characteristic assumption, which is usually data efficient, but lacks theoretical guarantees [9, 13]. Binning methods divide the data samples into finite buckets (bins), and then calibrate the samples in the bucket by its average prediction, which provides distribution-free theoretical guarantees but requires more data samples. Kumar et al. [13] propose the scaling-binning method that combines Platt scaling [17] and histogram binning [23] to balance the data efficiency and distribution-free theoretical guarantees. In our work, we also consider the hybrid of binning and scaling.

Although the widely used binning methods do not have distribution assumptions, we find the binning scheme of previous methods is too simple, limiting the calibration performance. Here, the binning scheme is a way to divide data samples into a finite number

of bins. For example, histogram binning with the uniform-mass scheme [23] sorts the samples by the predictions and then splits the samples into multiple bins, in which the number of samples are approximately the same. Histogram binning further assumes samples with close predictions have the same true probability, which is a widely used binning approach. However, even though the average prediction probability is close to the ground-truth, some samples in a bin may be overestimated, whereas others may be underestimated, because the samples in the same bin may still exhibit different bias patterns. We observe that the overestimation and underestimation are related to both the data sample itself and the predictor outputs. Thus, we should consider both the input features and predicted output into consideration for calibration, and put samples sharing a similar bias pattern into the same bin.

To this end, we propose a feature-aware binning framework that learns an improved binning scheme for calibration, which takes both input features and the output prediction of a predictor as the calibrator's inputs. We further employ multiple boosting trees as the backbone model for feature-aware binning due to the following considerations: 1) a path from the tree root to the leaf is a natural bin, which largely improves model interpretation; 2) tree-based methods enjoy the flexibility of handling various loss functions for node split; and 3) boosting combines multiple weak tree learners to a strong learner for better performance. We name the proposed framework as Multiple Boosting Calibration Trees (MBCT).

Moreover, we notice that traditional evaluation metrics usually adopt the uniform-mass scheme with sorted samples, which only reflects the calibration error on a certain group of partitions. Ideally, the loss of a well-performing calibrator should be near 0 under any partition. Therefore, we propose a multi-view calibration error (MVCE) as both the node-splitting criteria in training and the performance evaluation measure in experiments.

Finally, we apply a linear function in each node of MBCT to achieve individual calibration. We observe that traditional binning is typically non-individual, i.e., predicting the same value for all the samples in a bin. This restricts model capacity, making it impossible to improve order accuracy [21, 22]. By contrast, our linear calibrator function predicts a calibrated value for an individual sample.

To verify the effectiveness and generalization of our method, we conduct comprehensive experiments on three datasets in two application fields. One dataset is constructed by ourselves from real industrial data, and the other two are publicly available benchmark datasets. Experiments show that our method significantly outperforms competing methods in calibration error, and even improves the order accuracy and classification accuracy on certain datasets. Our MBCT is deployed in an online advertising system, achieving significant improvements on business metrics.

In general, our main contributions are four-fold: 1) We propose MBCT, a novel calibration method that combines binning and scaling methods. It improves the binning scheme by adopting multiple boosting trees and learning the bias pattern of the predictor with data. It also adopts elements of scaling methods, as we fit a linear function to achieve individual calibration, breaking the bottleneck of traditional calibration methods on order accuracy. 2) We propose a MVCE metric for calibration evaluation. It evaluates the calibration error from multiple perspectives, and is closer to the theoretical value than previous metrics. 3) We construct a dataset

of click-through rate prediction with real industrial logs. It fertilizes the research of uncertainty calibration and click-through rate prediction. 4) We conduct comprehensive experiments to verify the effectiveness of our method and metric. We further conduct online experiments in a real-world application to verify the impact on uncertainty calibration. Results show great commercial value of uncertainty calibration and our proposed method.

## 2 RELATED WORK

With the increasing demand for reliable and transparent machine learning in various applications, a large amount of calibration research work has recently emerged [1, 9, 13, 18, 20]. Most of the calibration methods serve as a post-processing step, which learns a mapping function to calibrate the output of an off-the-shelf classifier. These methods can be mainly divided into three categories: non-parametric, parametric, and hybrid methods.

**Parametric methods** often assume a parametric data distribution, and use the calibration data to fit the function for training. The Platt scaling [17] uses logistic transformation to modify the model output into a calibrated value. This method is further extended to the temperature scaling [8, 14] and Dirichlet scaling [11] that can be applied to multi-classification problems. Beta calibration [12] focuses on refining the performance on non-logistic predictors. Scaling methods are data efficient, but largely depend on the correctness of the distribution assumption. According to Gupta et al. [9], scaling methods cannot achieve distribution-free approximate calibration.

**Non-parametric methods** do not have distribution assumptions; examples include binning methods [15, 16] and isotonic regression [24]. The most popular binning method is histogram binning [23], which groups samples into multiple buckets (bins). Then for each sample, the average prediction results of its bucket is set to be the calibration result. Most of the binning methods (such as histogram binning) is non-strictly mon otonic. Although binning methods have theoretical guarantees in calibration accuracy [9, 13], the non-strict monotonicity may hurt applications that concern order accuracy. Isotonic regression [24] also divides the samples into bins and use an isotonic function for calibration with strictly monotonic constraints.

To take both advantages of scaling and binning methods, Kumar et al. [13] propose a **hybrid** method, called scaling-binning, to balance the sample efficiency and the verification guarantees. Zhang et al. [25] develop an ensemble method for both parametric and non-parametric calibrations. In our work, we also consider making full use of the parametric and non-parametric methods. Unlike Kumar et al. [13] and Zhang et al. [25], we adopt the binning framework because of its theoretical guarantee, but introduce parametric functions to our approach in consideration of order accuracy.

The **evaluation method** is another key issue in the research of calibration. Most studies use expected calibration error (ECE) and maximum calibration error (MCE) [8, 15, 16], but they are shown to be biased and usually underestimated [13]. Roelofs et al. [19] point out that the bucketing scheme, the number of buckets, and the amount of data will all affect ECE's stability, and claim that a good metric for calibration should not be so sensitive as ECE. So they propose a strategy to determine the partition method and bin

number for obtaining a more stable metric (named $\text{ECE}_{\text{sweep}}$). Experiments show that $\text{ECE}_{\text{sweep}}$ is less biased than existing metrics, such as ECE and MCE. In our paper, We further explore the calibration evaluation methods and propose the multi-view calibration error (MVCE).

## 3 PROBLEM FORMULATION

In Sections 3 and 4, we formulate the problem in the binary classification setting for simplicity; the extension to multi-class classification is straightforward.

Let $\mathcal{X}$ and $\mathcal{Y} = \{0, 1\}$ denote the feature and label spaces, respectively. Let $f : \mathcal{X} \to \mathcal{Y}$ be a predictor trained on the training set $D_{\text{train}} = \{(X_i, Y_i)\}_{i=1}^n$, where $X_i$ and $Y_i$ are drawn i.i.d from data distribution. Note that $Y$ is a random variable even when $X$ is given, and we would like $f$ to predict the expectation $\mathbb{E}[Y|X]$. The calibration error of $f$ with respect to the ground-truth under $\ell_p$-norm is defined as

$$\text{TCE}_p(f) = (\mathbb{E}_X[|\mathbb{E}[Y|X] - f(X)|^p])^{\frac{1}{p}} \qquad (1)$$

Typically, $p$ is set to 2 (see [19]). The prediction $f$ is *perfectly calibrated* if the calibration error is zero. In practice, most predictors $f$ (such as logistic regression and deep networks) are not perfectly calibrated and one usually adopts a post-hoc calibrator $h$, which is trained on a calibration set $D = \{(X_i, Y_i)\}_{i=1}^m$.

According to Gupta et al. [9], perfect calibration is impossible in practice, but approximate and asymptotic grouped calibration are possible for finite partitions of $\mathcal{X}$. Here, we restate them as follows for our subsequent analysis:

DEFINITION 1 (APPROXIMATE GROUPED CALIBRATION). *Suppose $\mathcal{X}$ is divided into $B$ partitions, denoted as $\{\mathcal{X}_b\}_{b=1}^B$. A calibrator $h$ is $(\varepsilon_b, \alpha)$-approximately grouped calibrated for the $b$th group for some $\alpha \in (0, 1)$ and $\varepsilon_b \in [0, 1]$, if with probability at least $1 - \alpha$,*

$$|\mathbb{E}[Y_b|X_b] - \mathbb{E}[h(f(X_b))]| \leq \varepsilon_b \qquad (2)$$

*for every $X_b \in \mathcal{X}_b$.*

DEFINITION 2 (ASYMPTOTIC GROUPED CALIBRATION). *A calibrator $h$ is asymptotically calibrated at level $\alpha \in (0, 1)$ if, for every $b \in \{1, \cdots, B\}$, $h$ is $(\varepsilon_b, \alpha)$-approximately calibrated for some $\varepsilon_b \in [0, 1]$ such that $\varepsilon_b = o_P(1)$, i.e., $\varepsilon_b$ converges in probability to 0 with the size of $D$ increasing.*

## 4 APPROACH

Our method generally follows the binning framework, but explores a better way of partitioning and further develops individual and non-monotonic calibration to break the bottleneck of traditional calibration methods. Concretely, we adopt Multiple Boosting Trees to learn a feature-aware binning scheme that optimizes a multi-view calibration loss. Our approach is individual but non-monotonic calibration, which is important to order-sensitive applications. We call our method **M**ultiple **B**oosting **C**alibration **T**rees (MBCT).

In the following subsections, we first introduce the weakness of traditional binning methods and present feature-aware binning to mitigate this problem. Then, we discuss the necessity of individual and monotonic calibration, and finally, we explain our MBCT method in detail.

## 4.1 Feature-Aware Binning

Previous work [9] has proved that binning-based methods can achieve distribution-free approximate grouped calibration and asymptotic grouped calibration. It also provides a theoretical bound for each partition (also referred to as a *bin*):

THEOREM 1. *Let $D_b$ denote the set of $D$ samples that fall into the partition $b$. For any $\alpha \in (0, 1)$, partition $b$, predictor $f$ and binning-based calibrator $h$, we have with probability at least $1 - \alpha$,*

$$|\mathbb{E}[Y_b] - \hat{y}_b| \leq \sqrt{\frac{2\hat{V}_b \ln(3B/\alpha)}{c_b}} + \frac{3\ln(3B/\alpha)}{c_b} \qquad (3)$$

*where $\hat{y}_b = \frac{1}{c_b} \sum_{i:X_i \in D_b} Y_i$, $\hat{V}_b = \frac{1}{c_b} \sum_{i:X_i \in D_b} (Y_i - \hat{y}_b)^2$, $c_b$ is the number of samples in the $b$th bin, and $B$ is the number of bins. For binning methods, we have $\hat{y}_b = \frac{1}{c_b} \sum_{X_i \in D_b} h(f(X_i))$.*

Note that the guarantee provided by Theorem 1 is distribution-free. On the contrary, scaling-based methods can only provide theoretical guarantee with distribution assumptions; thus, scaling-based methods are only suitable for specific tasks.

Although binning-based methods appear to be more general and have some theoretical guarantee, such a guarantee is only for specific partitions determined by the calibration method. **In the following, we explain the shortcomings of traditional binning methods in detail and propose feature-aware binning to address the limitations.**

Suppose $h$ is a binning-based calibrator, $h$ will have a division scheme $\text{div}_h$ that divides a dataset into finite subsets (bins). A function $g_b$ is learned to calibrate each bin by minimizing the expected partition calibration error.

DEFINITION 3. *For binning-based methods, a dataset $D$ will be divided into $\mathcal{D}^{(h)} = \{D_b\}_{b=1}^B$ by the division scheme $\text{div}_h$ of $h$. For a sample set $D_b = \{(X_i, Y_i)\}_{i=1}^{c_b}$ in $\mathcal{D}^{(h)}$, the **partition calibration error** is*

$$\text{PCE}(D_b) = |\hat{y}_b^{\text{pred}} - \hat{y}_b| \qquad (4)$$

*We may also define **expected partition calibration error** (EPCE) as:*

$$\text{EPCE}(D_b) = |\hat{y}_b^{\text{pred}} - \mathbb{E}[Y_b]| \qquad (5)$$

*where $\hat{y}_b^{\text{pred}} = \frac{1}{c_b} \sum_{X_i \in D_b} g_b(f(X_i))$.*

Note that the parameters for each bin $b$ may be different, so we use $g_b(f(\cdot))$ to denote the calibration with respect to $b$. We may abbreviate it as $g(f(\cdot))$ if there is no ambiguity. Binning methods ensure that $\text{PCE}(D_b) = 0$ for every $D_b \in \mathcal{D}^{(h)}$.

According to Definition 3 and Theorem 1, we have EPCE $\leq$ PCE $+ |\mathbb{E}[Y_b] - \hat{y}_b| \leq \sqrt{\frac{2\hat{V}_b \ln(3B/\alpha)}{c_b}} + \frac{3\ln(3B/\alpha)}{c_b}$. If the variance of $D_b$ is finite and its size $c_b$ is large enough, then $\text{EPCE}(D_b) \approx 0$. This ensures that the binning-based methods can still optimize the EPCE of the division even if the ground-truth cannot be obtained.

DEFINITION 4. *For a bin $D_b$, the **partition underestimated degree** (PUD) with respect to $f$ and $h$ is:*

$$\text{PUD}(D_b) = \hat{y}_b^{\text{pred}} / \hat{y}_b \qquad (6)$$

Similarly, for one sample $X_i$, we define the sample underestimated degree (SUD) as:

$$\mathrm{SUD}(X_i) = g(f(X_i))/\mathbb{E}[Y_i|X_i] \tag{7}$$

Apparently, the bin $D_b$ is underestimated when $\mathrm{PUD}(D_b) < 1$, or overestimated when $\mathrm{PUD}(D_b) > 1$. For SUD, it is not possible to calculate its value in practice, it will facilitate our theoretical analysis.

We say a bin $D_b$ is $(\beta, h)$-**well-calibrated** when $\mathrm{EPCE}(D_b) \le \beta$. If $\beta$ is small enough in real scenarios, we say $D_b$ is a **well calibrated partition**.

For a bin $D_b$, we may further split it by a division scheme div into $k$ equally sized subsets $\{D_{b_i}\}_{i=1}^{k}$. For simplicity, we define the $(k, \mathrm{div})$-**balanced finer-grained partition calibration error** (BFGPCE) of $D_b$ as:

$$\mathrm{BFGPCE}_{k,\mathrm{div}}(D_b) = \frac{1}{k} \sum_{i=1}^{k} |\hat{y}_{b_i}^{\mathrm{pred}} - \hat{y}_{b_i}| \tag{8}$$

Further, $(k, \mathrm{div})$-**balanced finer-grained expected partition calibration error** (BFGEPCE) is defined as

$$\mathrm{BFGEPCE}_{k,\mathrm{div}}(D_b) = \frac{1}{k} \sum_{i=1}^{k} |\hat{y}_{b_i}^{\mathrm{pred}} - \mathbb{E}[y_{b_i}]| \tag{9}$$

For a certain well-calibrated subset $D_b$, $\mathrm{BFGPCE}_{k,\mathrm{div}}(x_b)$ and $\mathrm{BFGEPCE}_{k,\mathrm{div}}(x_b)$ may still be large if a different division is used, in which case the calibration performance may not be satisfactory. Then, a natural question follows: Is it possible to achieve near-perfect calibration based on the binning method? Perfect calibration under binning methods can be restated as: For each bin $D_i \in \mathcal{D}^{(h)}$, (1) $D_b$ is well calibrated under $h$, and (2) each sample $(X_i, Y_i) \in D_b$ satisfies $\mathbb{E}[Y_i|X_i] = g_b(f(X_i))$. Here, we call $g_b$ the bias pattern of the bin $D_i$. Notice that condition (1) can be achieved by all the binning-based methods with sufficient samples. Condition (2) is a sufficient condition for (1). Thus, if we can put the samples that have the same bias pattern into the same bin, we will achieve a near-perfect calibration.

Traditional binning methods aim to find a constant bias pattern and usually assume that the samples' groundtruth probabilities are approximately the same if the predicted probabilities by $f$ are close; hence, the samples are sorted by the predicted probabilities and divided into partitions in an uniform-mass [23] or equal-width [8] manner. However, the assumption of traditional binning-based methods is too strong. In fact, the data may deviate from this assumption; especially, the difference is large when the estimation accuracy is low. Kumar et al. [13] also prove that there exist better partition schemes than uniform-mass and equal-width methods with sorted samples, although they do not provide a method to find them.

In this paper, we aim to find a better partition scheme for binning-based calibration. We notice that, if $h$ can achieve approximately perfect calibration, we will have the following property: for any subset $D_{sub}$ (with sufficient samples) of $D$, we get $\mathrm{EPCE}(D_{\mathrm{sub}}) \approx 0$. In other words, the more partitions with $\mathrm{EPCE} \approx 0$, the closer to perfect calibration $h$ would be. Thus, we propose the multi-view calibration error for evaluating the calibration with respect to a set

of partitions $\{\mathrm{div}_i\}_{i=1}^{r}$:

$$\mathrm{MVCE}_{f,h,\{\mathrm{div}_i\}_{i=1}^{r}}(D_b) = \left( \frac{1}{r} \sum_{i=1}^{r} \left( \frac{1}{t_i} \sum_{j=1}^{t_i} \mathrm{PCE}(D_{i,j}) \right)^p \right)^{\frac{1}{p}} \tag{10}$$

where $t_i$ is the number of bins under $\mathrm{div}_i$ and $D_{i,j} \in \mathrm{div}_i(D)$. Intuitively, MVCE is suitable for measuring how close the calibration is to the optimum. Another intuition is the bias between $f(X)$ and $\mathbb{E}[Y]$ would have different widespread patterns related to $X$ that can be captured by machine learning methods. Therefore, **we propose a feature-aware binning framework to learn which samples belong to a bias pattern**. It employs a model $M$ which takes the features $X$ as input and optimizes the MVCE. Optimizing MVCE can force the model to put the samples with a similar bias pattern into the same bin.

Specifically, we employ a linear function $g(f(\cdot), b) = k_b f(\cdot)$ as the bias pattern to achieve individual calibration (Section 4.2). In other words, we would like to find samples that have the same SUD. Details of our methods are presented in Section 4.3.

## 4.2 Individual and Monotonic Calibration

Besides the true probability, order accuracy, often measured by the area under the curve (AUC), is also important to various applications, such as computation advertising [27, 28] and web search [2]. Traditional binning-based methods are non-strictly monotonic because they set the same calibrated for samples in the same bin, which is also considered as non-individual. In this part, we first give the definitions of individual and monotonic calibration and then discuss their necessity in uncertainty calibration.

**Individual** calibration means that we predict a calibrated value for an individual sample. By contrast, a *non-individual* calibrator predicts a value for a set of pooled samples.

**Monotonic** calibration means that the order is preserved. Concretely, for a dataset $D = \{(X_i, Y_i)\}_{i=1}^{m}$, we call $h$ a **non-strictly monotonic** calibrator if for every predictor $f$ and every two samples $(X_i, Y_i), (X_j, Y_j) \in D$, we have:

$$(f(X_i) - f(X_j)) \cdot (h(f(X_i)) - h(f(X_j))) \ge 0 \tag{11}$$

If the equality holds only when $f(X_i) = f(X_j)$, then $h$ is a **strictly monotonic** calibrator.

If there exists two samples violating Equation (11), we say $h$ is a **non-monotonic** calibrator.

In practice, non-individual calibration is harmful to order-sensitive applications, such as computation advertising. For example, we would like to select samples with top-3 probabilities from a sample set; non-individual calibrators may give 10 samples in a bin with the same predicted probability. Thus, individual calibration is crucial to order-sensitive uncertainty calibration, and our MBCT makes samples distinguishable by learning a linear function $g$ for a bin.

On the other hand, traditional methods often apply non-strictly or strictly monotonic calibrators, assuming that the base classifier $f$ can predict the correct order. However, this assumption appears too strong, and restricts the power of the calibrator. In our work, we lift such constraint and our MBCT may be non-monotonic.
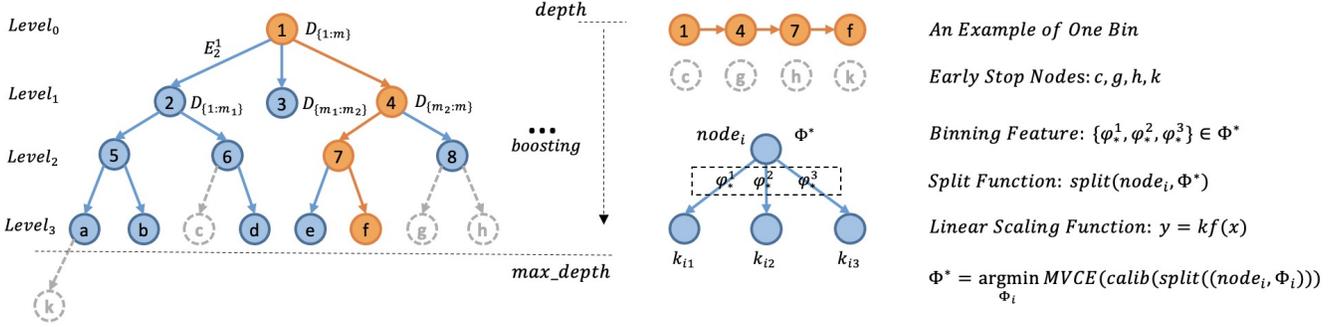
**Figure 1: An illustration of our Multiple Boosting Calibration Trees (MBCT). The left-hand side shows an example of a binning tree. MBCT may have multiple boosting trees to minimize our MVCE loss. The right-hand side shows several key elements and concepts in our MBCT. Consider a node $i$ with corresponding data $D_i$. It selects a feature $\Phi_*$ by MVCE, and splits the node by the feature values $\varphi_*^1, \varphi_*^2, \varphi_*^3$. For each child node, a linear function $g_j(x) = k_j x$ is fitted for calibration.**

## 4.3 Multiple Boosting Calibration Trees

As mentioned, our paper proposes a feature-aware binning framework, which adopts a machine learning model $M$ to optimize the MVCE loss; our calibrator should be individual but non-monotonic.

Concretely, we adopt multiple boosting trees as the model $M$ mainly for the following considerations: (1) Tree-based models are more friendly to interpretation, and it is natural to assign a bin to each node in the tree. (2) It is easy to optimize customized losses (such as MVCE) with boosting trees. (3) Boosting methods are shown to be effective in various applications [4, 6].

Figure 1 shows an overview of our approach, Multiple Boosting Calibration Trees (MBCT). It takes discretized features as input. Let $\mathscr{F} = \{\Phi_i\}_{i=1}^d$ be the feature space, each feature $\Phi_i$ taking values $\varphi_i^1, \cdots, \varphi_i^a$. We adopt the random shuffle uniform-mass partition scheme as division of MVCE, i.e., we first randomly shuffle the samples (instead of sorting by the predictions), and then perform a uniform-mass division. This enables us to obtain multiple uniform-mass divisions.

There are three key operations during the growth of a binning tree: (1) For each leaf node, select the best feature to split the node. For a node with its data $D'$, the best feature is obtained by

$$\Phi_* = \operatorname{argmin}_{\Phi_i \in \mathscr{F}} \operatorname{MVCE}(\operatorname{calib}(\operatorname{split}(D', \Phi_i)))) \qquad (12)$$

where split($\cdot$) is an operation that splits a node in the calibration tree by the values of $\Phi_i$, and calib($\cdot$) calculates the calibration parameters (in our model, it is a linear function). (2) Split the leaf node by its $\Phi^*$ and calibrate its children, if the early stop conditions are not met and $\Phi^*$ can split the data of the node at least into 2 subsets. (3) If there is no leaf node can be extended, finish the growth of the binning tree.

The stopping conditions for growing one tree include: (1) The node reaches the maximum depth set as a hyper-parameter. (2) The samples in the node is less than $\beta$, which is heuristically estimated by Equation (3). Concretely, we first calculate $\hat{V}_D$ and $\hat{y}_D$ by Equation (3) on the whole dataset $D$. Then, we assume MBCT is close to a uniform-mass binning which satisfies $B = \frac{|D|}{c}$, where $c$ is the sample number in each bin. Finally, given a permissible absolute percentage error limit $e$ and confidence $\alpha$, we obtain the

maximum $c$ that satisfies:

$$\hat{y}_D \leq \frac{1}{e}\left(\sqrt{\frac{2\hat{V}_D \ln(3B|D|/c\alpha)}{c}} + \frac{3\ln(3B|D|/c\alpha)}{c}\right) \qquad (13)$$

We set $\beta$ as $c$, and we simply set the bin size of the loss as $\frac{c}{2}$ because we need at least 2 bins to calculate the calibration loss. $\alpha$ and $e$ are hyper-parameters for our method. (3) The local loss rises after splitting and calibrating the node. Here, the local loss refers to the loss on the subset corresponding to the node to be divided and calibrated. We do not adopt global loss because it will greatly increase the complexity for MBCT.

MBCT allows to generate multiple binning trees. We decide whether to generate a new binning tree according to the change of global loss. A path from the root node to the leaf node represents a specific bin. Multiple binning trees are actually doing the ensemble of different binning schemes for further optimizing the calibration loss. With these operations, MBCT optimizes calibration loss in a greedy way. Although greedy optimization may not be optimal, it is considered as a very effective optimization algorithm in various models, including boosting trees.

## 5 EXPERIMENTS

### 5.1 Experimental Setup

We first conducted simulation experiments to verify the rationale for using MVCE to evaluate calibration errors, which is also applied to our online advertising system for calibrating the predicted Click-Through Rate (pCTR). Then, we conducted both offline and online experiments to verify the effectiveness of our method and the impact on the real-world application.

For the offline experiment, we dumped the impression and click logs from our commercially deployed online display advertising system, and we call our industrial dataset CACTRDC[1], which is the acronym for Computation Advertising Click-Through Prediction Dataset for Calibration.

---

[1]https://github.com/huangsg1/Tree-Based-Feature-Aware-Binning-for-Individual-Uncertainty-Calibration

Table 1: Dataset statistics. "Predictor Train" refers to the
training set of the base predictor $f$. "Calibration Train" and
"Calibration Test" are the training and test sets for the cali-
brator $h$, respectively.

| Dataset | Predictor Train | Calibration Train | Calibration Test |
|---|---|---|---|
| CACTRDC | 48M | 11M | 1M |
| Porto Seguro | 357K | 208K | 30K |
| Avazu | 24M | 12M | 4M |

In addition, we tested our method on two public datasets, Porto
Seguro[2] and Avazu[3], showing the generality of our method. Table 1
presents the statistics of these datasets.

Regarding the base predictor, we used a deep neural network
(DNN) for CACTRDC and our online A/B testing, as already de-
ployed in the system. For the Porto Seguro and Avazu datasets, we
use the Gradient Boosting Decision Tree (GBDT), because we would
like to experiment with non-DNN models. In principle, our calibra-
tion is agnostic to the base predictor; this further demonstrates the
generality of our calibration approach.

In our experiments, we conducted feature engineering and tuned
hyperparameters on the CACTRDC dataset, which is related to
our deployed commercial system. For the Porto Seguro and Avazu
benchmark datasets, we mostly adopted the same settings, as the
main focus of our approach is not hyperparameter tuning; this,
again, shows the generality of our MBCT. For discrete features,
we directly take their values when splitting a node by the feature;
continuous features are discretized by quantiles. We also convert
the outputs of the predictor into discrete features for MBCT. More
feature engineering details are shown in appendix A.1.

For hyperparameters, we set $\alpha$ to 0.05 and $e$ to 0.1 (Section 4.3),
so the bin sizes were 800, 2000, and 1000 for Porto Seguro, Avazu,
and CACTRDC, respectively. The max tree depth was 5 and the
max tree number was 8.

For evaluation of the calibration methods, we adopt MVCE for
calibration accuracy and area under the curve (AUC) for order
accuracy. We set the $r$ of MVCE to 100.

## 5.2 Competing Methods

We compare our MBCT model with the following state-of-the-art
methods in previous studies.

- **Platt Scaling.** Platt et al. [17] propose a scaling method that
  uses logistic function to transform the model predictions into
  calibrated probabilities.
- **Beta Calibration.** Kull et al. [12] refine the assumption of iso-
  tonic calibration, which is more suitable for non-logistics predic-
  tors, such as Naïve Bayes and AdaBoost.
- **Histogram Binning.** Zadrozny and Elkan [23] propose a method
  that divides the samples into bins and calibrates each sample by
  the average probability of its bin as the calibrated value.
- **Isotonic Regression.** Zadrozny and Elkan [24] propose a method
  that learns an isotonic function for calibration with the strictly
  monotonic constraint.

- **Scaling-Binning.** Kumar et al. [13] propose a scaling-binning
  calibrator that combines the Platt scaling and Histogram binning.
  The approach first fits a parametric function to reduce variance
  and then puts the function values into bins for calibration. Note
  that it is a non-individual and non-strictly monotonic calibrator,
  which is different from our method.

## 5.3 Analysis of the Metrics for Calibration Error

As stated in Section 4.1, MVCE is intuitively an appropriate metric
for evaluating calibration error. Also, we notice that previous work
often uses ECE or $\text{ECE}_{\text{sweep}}$ to estimate the calibration error. ECE
with $n$ bins in $\ell_p$-norm can be formulated as:

$$\text{ECE}_n = \left(\frac{1}{N} \sum_{i=1}^{N} \text{PCE}(D_i)^p\right)^{\frac{1}{p}} \tag{14}$$

where $D_i$ represents the samples in the bin $i$, and $n$ is a hyper-
parameter referring to the number of the bins for evaluation. Note
that in ECE samples are sorted by calibration results, based on
which uniform-mass binning is performed. $\text{ECE}_{\text{sweep}}$ is proposed
by Roelofs et al. [19] and has a strategy to determine the partition
method and bin number, which is shown to be less biased than ECE.

Intuitively, our MVCE evaluates the calibration error from a more
comprehensive perspective; ECE and $\text{ECE}_{\text{sweep}}$ are special cases in
MVCE. In this part, we will evaluate these metrics with a simulation
experiment. In fact, the metric for calibration ultimately concerns
the theoretical calibration error (TCE), which unfortunately cannot
be computed on real-world data, because the true probability of a
sample in unknown. Using simulated data with certain assumptions
allows us to compute TCE and evaluates a calibration metric.

Therefore, we design a simulation experiment, where we assume
$f(X)$ follows a certain known distribution and bias pattern, so
that the TCE can be computed. Concretely, we assume $h(X) \sim$
$\text{Beta}(0.2, 0.7)$ and $\mathbb{E}[Y|h(X) = c] = c^2$, where we hide the role of
$X$ as it is irrelevant to our simulation. Then, we can compute the
TCE according to Equation (1) with the $\ell_2$-norm. The TCE in our
simulation experiment is 0.0868.

We evaluate a metric $\mu$ (e.g., ECE, $\text{ECE}_{\text{sweep}}$, and MVCE; also
with $\ell_2$-norm) by comparing it with TCE. The expected difference
for a sample $X$ is formulated as:

$$E_{bias} = |\mathbb{E}[\mu(X, h)] - \text{TCE}(X, h)| \tag{15}$$

We used the Monte Carlo method to estimate $E_{bias}$, where we
conducted $m = 200$ experiments, each with random $n$ samples. The
empirically estimated difference is given by

$$\hat{E}_{bias}(n) = \frac{1}{m} \sum_{i=1}^{m} |\mu(X_i^{(n)}, h) - \text{TCE}(X, h)| \tag{16}$$

where $X_i^{(n)}$ means a set of $n$ samples.

Figure 2 shows the estimated difference of ECE, $\text{ECE}_{\text{sweep}}$, and
MVCE under various sample sizes $n$. We see that $\text{ECE}_{\text{sweep}}$ is better
than ECE, which is consistent with the experimental conclusion
of Roelofs et al. [19]. MVCE achieves the best result, indicating
that MVCE is a better metric for calibration error than ECE and
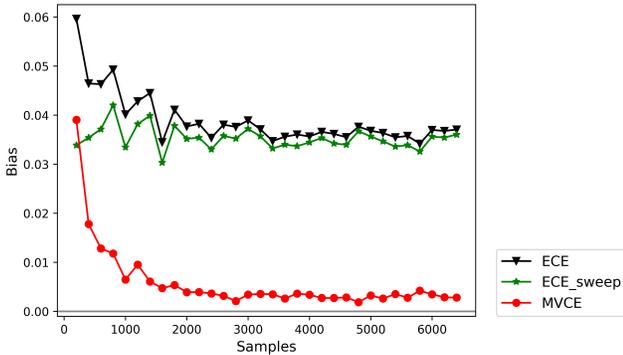$\text{ECE}_{\text{sweep}}$.

Figure 2: Main simulation results of ECE, ECE$_{sweep}$ and MVCE. The bin numbers of ECE and MVCE are set as 32 in this experiment. The bin number of ECE$_{sweep}$ is fixed and determined by itself.
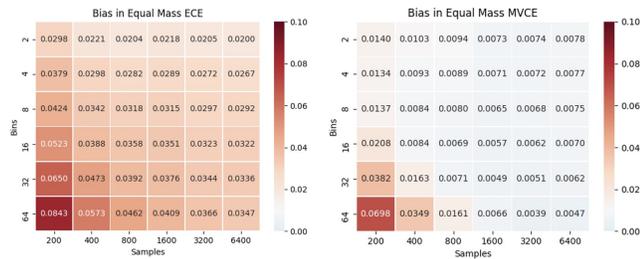


Figure 3: $\hat{E}_{bias}$ of MVCE and ECE under different numbers of bins and samples.

We further study the impact of different combinations of bin and sample numbers on ECE and MVCE in terms of $\hat{E}_{bias}$ in Figure 3. We see that (1) MVCE is consistently better than ECE in different cases; that (2) more data lead to less bias for ECE under a certain bin number, but the larger the amount of data, the smaller the effect is; and that (3) the bin number should be moderate for MVCE with respect to the sample number, and neither too large nor too small is desired. We provide more supplement results with different settings in appendix B.2.

## 5.4 Main Results

Table 2 presents the main results of our approach on three datasets. Numbers in bold indicate the best performance. A higher score of AUC indicates a better model; for MVCE, a lower score is better.

We see that MBCT consistently outperforms the original base predictor and all the baseline calibration methods in terms of both MVCE and AUC metrics, which indicates that our method can bring improvement on both calibration and order accuracy, and is generalizable to different datasets. The AUC of other methods fluctuates after calibration, which is because these methods may turn different estimates into exactly the same.

Moreover, we conduct an ablation study of the boosting strategy in the last two rows of Table 2. As seen, boosting consistently brings in improvement of AUC and MVCE on CACTRDC, Porto Seguro, and Avazu. Even without boosting, our method achieves

Table 2: Results for uncertainty calibration on the benchmark datasets.

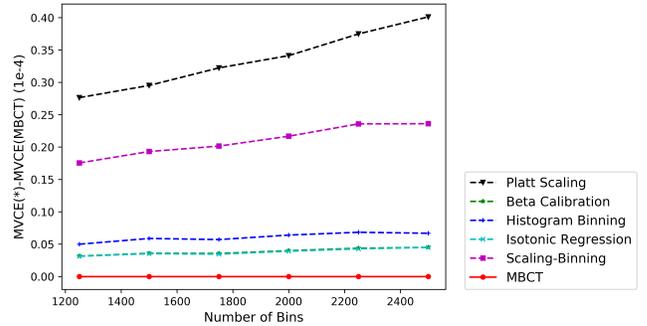| Method | CACTRDC | | Porto Seguro | | Avazu | |
|---|---|---|---|---|---|---|
| | MVCE ↓ | AUC ↑ | MVCE ↓ | AUC ↑ | MVCE ↓ | AUC ↑ |
| Original Predictions | 0.00394 | 0.77902 | 0.00619 | 0.62869 | 0.00976 | 0.71880 |
| Platt Scaling | 0.00374 | 0.77902 | 0.00604 | 0.62869 | 0.00792 | 0.71880 |
| Beta Calibration | 0.00371 | 0.77902 | 0.00601 | 0.62869 | 0.00789 | 0.71880 |
| Histogram Bining | 0.00372 | 0.77895 | 0.00597 | 0.62998 | 0.00787 | 0.72381 |
| Isotonic Regression | 0.00371 | 0.77915 | 0.00598 | 0.62936 | 0.00787 | 0.72030 |
| Scaling-Binning | 0.00373 | 0.77892 | 0.00605 | 0.62880 | 0.00792 | 0.71871 |
| **Our full MBCT** | **0.00368** | **0.78693** | **0.00586** | **0.63097** | **0.00780** | **0.74177** |
| **Our model w/o boosting** | 0.00374 | 0.78373 | 0.00595 | 0.62999 | 0.00784 | 0.73797 |



Figure 4: The difference of MVCE scores in comparison with our MBCT. Natually, our approach yields a value of 0. All competing models have a positive value, showing that our approach is the best.

the best AUC on all the benchmarks and achieves the best MVCE on Porto Seguro and Avazu. We notice that boosting brings more improvement of AUC on large datasets; this is probably becasue larger datasets can benefit more from the enlarged model capacity brought by boosting.

## 5.5 In-Depth Analysis

In this part, we conduct in-depth analysis for our model. Due to the limit of time and space, we take the CACTRDC dataset as the testbed, unless otherwise stated.

In Table 2, we simply set the bin size of MVCE and ECE to the minimum bin size ($\beta$ in Section 4.3) during training. The MVCE score under different bin sizes is also a curious question, and we show the trend in Figure 4. We see that MBCT consistently outperforms the baselines on MVCE under all bin numbers, which further illustrates the effectiveness of our method.

Then, we analyze the effect of training losses. We compare our MVCE loss with traditional ECE and ECE$_{sweep}$ losses in Figure 5. We see that using MVCE as the training loss consistently outperforms ECE and ECE$_{sweep}$ on the three datasets, which shows that it is beneficial to consider the average PCEs from different perspectives during training.

We further conduct an analysis to show whether MBCT puts samples with similar bias patterns for $f$ into the same bin. We randomly select 5 bins of MBCT and Histogram methods, and divide each bin into 4 sub-groups. Each bin of MBCT is corresponding to a leaf node. Figure 6 shows the PUD of these sub-groups. It is
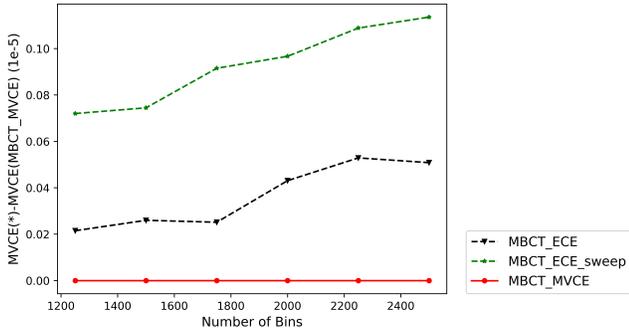
**Figure 5: Ablation study of the training loss (in comparison with MVCE training).**
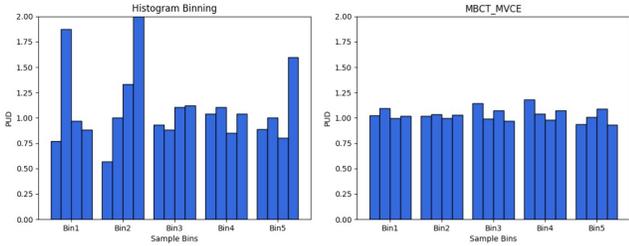


**Figure 6: PUD of finer-grained partitions in MBCT and Histogram Binning.**

obvious that the sub-groups of MBCT is closer to 1, which implies that MBCT learns better relationships between the inputs and the bias pattern of $f$ than naive binning methods.

To better understand our MBCT model, we present a case study in Figure 7, showing the growth of the calibration tree. It shows a few nodes in the top 3 layers of the first tree of MBCT. As seen, our approach splits a node by optimizing the local MVCE, and this greedy way can effectively reduce the global loss.



**Figure 7: Visualization of the greedy optimization in MBCT. Local MVCE of one node is in blue, and the Global MVCE of the whole dataset is in orange.**

## 5.6 The Impact of Calibration on Real-world Applications

To the best of our knowledge, previous work has not applied their calibration methods to real-world applications. It is important to test whether the calibration can improve performance in the real world. We deploy MBCT to calibrate the click-through rate for an online advertising system and conduct an online A/B test for 15 days. We use click-through rate (CTR) and the effective cost per mille (eCPM) to evaluate the models, which are the key performance metrics:

$$\text{CTR} = \frac{\text{\# of clicks}}{\text{\# of impressions}}, \ \text{eCPM} = \frac{\text{Ad revenue}}{\text{\# of impressions}} \quad (17)$$

We build three buckets on the online traffic platform to conduct fair comparison. We reserve bucket A with no calibration. Isotonic regression and MCBT are applied to buckets B and C, respectively. Each bucket has 3% of our online traffic, which is large enough considering the overall page view requests. Table 3 lists the promotion of the two main online metrics. 5.7% growth on AUC exhibits that MBCT could improve the order accuracy of the online predictor. As for CTR, it has a 22.1% improvement, indicating that MBCT can bring more clicks because order accuracy is highly related to the allocation efficiency. Regarding eCPM, a 8.4% lift demonstrates more revenue for the platform. These improvements are significant for an online advertising system, demonstrating great business value. We provide more details and discuss the relationship between calibration and revenue in Appendix A.2 and C.

**Table 3: Online results of 15 days in comparison with no calibration and Isotonic Regression.**

| Metric | No Calibration | Isotonic Regression | MBCT |
|--------|----------------|---------------------|--------|
| AUC    | 0.0%           | +0.3%               | +5.7%  |
| CTR    | 0.0%           | +8.8%               | +22.1% |
| eCPM   | 0.0%           | +4.5%               | +8.4%  |

## 6 CONCLUSION

In this paper, we propose a novel calibration method, called MBCT, which extends the binning framework in several ways. We adopt multiple boosting trees in the feature-aware binning framework, and propose the multi-view calibration error (MVCE) to learn a better binning scheme. We also learn a linear function to calibrate the samples in a bin. Thus, our approach is individual but non-monotonic, which is important to order-sensitive applications.

We conduct experiments on both public datasets and an industrial dataset, covering different application fields. Experimental results show that MBCT significantly outperforms the competing methods in terms of calibration error and order accuracy; profoundly, it outperforms the original prediction in order accuracy, which cannot be achieved by monotonic calibrators. We also conduct simulation experiments to verify the effectiveness of our proposed MVCE metric. Our approach is deployed in an online advertising system. Results show high performance and great commercial value of our approach in real-world applications.

## REFERENCES

[1] Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Abbas Khosravi, U Rajendra Acharya, Vladimir Makarenkov, et al. 2020. A review of uncertainty quantification in deep learning: techniques, applications and challenges. *arXiv preprint arXiv:2011.06225* (2020).

[2] Alexey Borisov, Julia Kiseleva, Ilya Markov, and Maarten de Rijke. 2018. Calibration: A simple way to improve click models. In *ACM International Conference on Information and Knowledge Management*. 1503–1506.

[3] Jochen Bröcker. 2009. Reliability, sufficiency, and the decomposition of proper scores. *Journal of the Atmospheric Sciences and Physical Oceanography* 135, 643 (2009), 1512–1519.

[4] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: a scalable tree boosting system. In *ACM SIGKDD international conference on Knowledge discovery and data mining*. 785–794.

[5] Morris H DeGroot and Stephen E Fienberg. 1983. The comparison and evaluation of forecasters. *Journal of the Royal Statistical Society: Series D (The Statistician)* 32, 1-2 (1983), 12–22.

[6] Yoav Freund and Robert E Schapire. 1997. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences* 55, 1 (1997), 119–139.

[7] Tilmann Gneiting and Adrian E Raftery. 2005. Weather forecasting with ensemble methods. *Science* 310, 5746 (2005), 248–249.

[8] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017. On calibration of modern neural networks. In *International Conference on Machine Learning*. 1321–1330.

[9] Chirag Gupta, Aleksandr Podkopaev, and Aaditya Ramdas. 2020. Distribution-free binary classification: prediction sets, confidence intervals and calibration. In *Advances in Neural Information Processing Systems*.

[10] Xiaoqian Jiang, Melanie Osl, Jihoon Kim, and Lucila Ohno-Machado. 2012. Calibrating predictive model estimates to support personalized medicine. *Journal of the American Medical Informatics Association* 19, 2 (2012), 263–274.

[11] Meelis Kull, Miquel Perello-Nieto, Markus Kängsepp, Hao Song, Peter Flach, et al. 2019. Beyond temperature scaling: obtaining well-calibrated multiclass probabilities with dirichlet calibration. *arXiv preprint arXiv:1910.12656* (2019).

[12] Meelis Kull, Telmo Silva Filho, and Peter Flach. 2017. Beta calibration: a well-founded and easily implemented improvement on logistic calibration for binary classifiers. In *Artificial Intelligence and Statistics*. 623–631.

[13] Ananya Kumar, Percy Liang, and Tengyu Ma. 2019. Verified uncertainty calibration. In *Advances in Neural Information Processing Systems*. 3787–3798.

[14] Azadeh Sadat Mozafari, Hugo Siqueira Gomes, Wilson Leão, Steeven Janny, and Christian Gagné. 2018. Attended temperature scaling: a practical approach for calibrating deep neural networks. *arXiv preprint arXiv:1810.11586* (2018).

[15] Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. 2015. Obtaining well calibrated probabilities using bayesian binning. In *AAAI Conference on Artificial Intelligence*, Vol. 29.

[16] Mahdi Pakdaman Naeini, Gregory F Cooper, and Milos Hauskrecht. 2014. Binary classifier calibration: non-parametric approach. *arXiv preprint arXiv:1401.3390* (2014).

[17] John Platt et al. 1999. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in Large Margin Classifiers* 10, 3 (1999), 61–74.

[18] Amir Rahimi, Kartik Gupta, Thalaiyasingam Ajanthan, Thomas Mensink, Cristian Sminchisescu, and Richard Hartley. 2020. Post-hoc calibration of neural networks. *arXiv preprint arXiv:2006.12807* (2020).

[19] Rebecca Roelofs, Nicholas Cain, Jonathon Shlens, and Michael C Mozer. 2020. Mitigating bias in calibration error estimation. *arXiv preprint arXiv:2012.08668* (2020).

[20] Seonguk Seo, Paul Hongsuck Seo, and Bohyung Han. 2019. Learning for single-shot confidence calibration in deep neural networks through stochastic inferences. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9030–9038.

[21] Di Wu, Xiujun Chen, Xun Yang, Hao Wang, Qing Tan, Xiaoxun Zhang, Jian Xu, and Kun Gai. 2018. Budget constrained bidding by model-free reinforcement learning in display advertising. In *ACM International Conference on Information and Knowledge Management*. 1443–1451.

[22] Xun Yang, Yasong Li, Hao Wang, Di Wu, Qing Tan, Jian Xu, and Kun Gai. 2019. Bid optimization by multivariable control in display advertising. In *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1966–1974.

[23] Bianca Zadrozny and Charles Elkan. 2001. Obtaining calibrated probability estimates from decision trees and naive Bayesian classifiers. In *International Conference on Machine Learning*. 609–616.

[24] Bianca Zadrozny and Charles Elkan. 2002. Transforming classifier scores into accurate multiclass probability estimates. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 694–699.

[25] Jize Zhang, Bhavya Kailkhura, and T Yong-Jin Han. 2020. Mix-n-match: ensemble and compositional methods for uncertainty calibration in deep learning. In *International Conference on Machine Learning*. 11117–11128.

[26] Shengjia Zhao, Tengyu Ma, and Stefano Ermon. 2020. Individual calibration with randomized forecasting. In *International Conference on Machine Learning*. 11387–11397.

[27] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. 2019. Deep interest evolution network for click-through rate prediction. In *AAAI Conference on Artificial Intelligence*, Vol. 33. 5941–5948.

[28] Han Zhu, Daqing Chang, Ziru Xu, Pengye Zhang, Xiang Li, Jie He, Han Li, Jian Xu, and Kun Gai. 2019. Joint optimization of tree-based index and deep model for recommender systems. In *Advances in Neural Information Processing Systems*. 3973–3982.

## A IMPLEMENTATION DETAILS

We describe the implementation details for both offline and online experiments to support the reproducibility of our work.

### A.1 Feature Engineering of Offline Experiments

- **CACTRDC Dataset**: We extracted eight features from the raw features. Seven were discrete features, and the other was 100-equal-width discretized predicted click-through rate (pCTR). The features were selected according to our online feature importance evaluation system. We only adopted discretized features to facilitate not only the application of the tree structure, but also the engineering optimization of online performance (run-time and the consumption of computing resources), detailed in Section A.2.
- **Porto Seguro and Avazu Datasets**: To make the experimental settings as consistent with CACTRDC, we only considered discretized features and selected seven discrete features based on GBDT's feature importance. We also took the 100-equal-width discretized predicted value as a calibration feature. For Porto Seguro, the selected seven features are ps_ind_05_cat, ps_ind_17_bin, and ps_car_01_cat, ps_car_06_cat, ps_car_07_cat, and ps_car_09_cat. For Avazu, the selected seven features are banner_pos, site_id, site_category, device_type, C16, C19, and C20.

### A.2 Online Deployment

Our online model for CTR prediction is a DNN model and its AUC reaches about 0.77. The model uses several hundred features, mainly involving user profiles (e.g., age and area), user behaviors (e.g., recently clicked ads), position, current time, and properties of ad (e.g., text, image, and ID). Our offline CACTRDC experiment considered same set of features during feature selection.

In order to meet the stability and latency requirements of online service, we could transform the trained MBCT model into a systematic "if...then..." rules, which benefits a lot from the interpretability of tree model. Then, we used a high-speed cache to store these "rules" for online deployment. When new user-generated data were accumulated, we frequently re-trained and updated our MBCT model to ensure real-time calibration; this improved business profits for the online advertising platforms.

We also deployed the following engineering techniques for MBCT to accelerate the running time and reduce the memory overhead:

**Data aggregation**: The MBCT algorithm takes discrete binning features as input. If a feature is continuous, we converted it to discrete one, so the combinations of the features are finite and we can map the samples into finite groups. This means we can aggregate the samples with the same feature group and use the average statistics of one group as a new sample to replace the samples in this group. The data aggregation can largely reduce the use of computational resources.

**Parallelization**: We follow the parallelization method proposed by XGBoost [4] for accelerating the running time of MBCT. In the process of constructing the calibration tree, we parallelize $s$ processes at any level of the tree for selecting the optimal binning feature, where $s$ is the number of features. This parallelization method
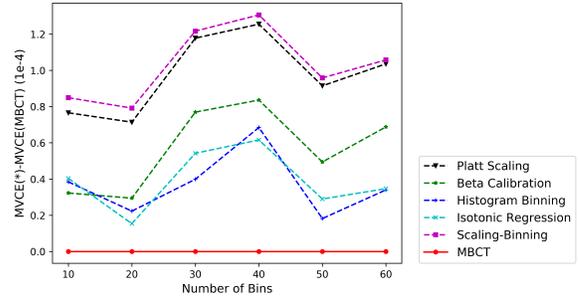


**Figure 8: The difference of MVCE scores in comparison with our MBCT on Porto Seguro. Natually, our approach yields a value of 0. All competing models have a positive value, showing that our approach is the best.**
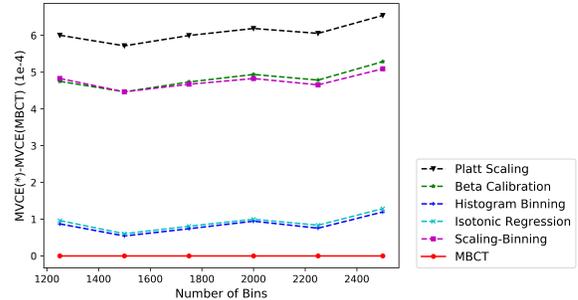


**Figure 9: The difference of MVCE scores in comparison with our MBCT on Avazu.**

guarantees load balancing to the largest extent, resulting in roughly $s$ times run-time acceleration. Concretely, the time complexity is optimized from $O(k * \min(s, d) * s * r * n)$ to $O(k * \min(s, d) * r * n)$, where $k$ is the tree number, $s$ is the number of features, $d$ is the max depth of the tree, $r$ is the hyper-parameter for MVCE, and $n$ is the number of train samples.

## B ADDITIONAL EXPERIMENTAL RESULTS

### B.1 In-Depth Analysis

We show additional analysis on the other two datasets, which are not presented in the main text due to space limit (cf., Section 5.5).

Figures 8 and 9 show the MVCE score under different bin sizes of the baselines and our method on Porto Seguro and Avazu, respectively. We see that MBCT also achieves the best MVCE score on these two public datasets, which is consistent with the results on real-world data of our online system. Figures 10 and 11 show the visualization of the PUD of finer-grained partitions on Porto Seguro and Avazu, verifying the effectiveness of feature-aware binning. All of the results in this section further illustrates that our method has a high generalization ability.

### B.2 Analysis of Calibration Metrics

In this part, we provide additional simulation experiments with different data distributions.
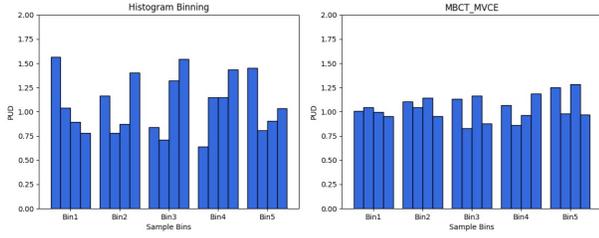
Figure 10: PUD of the finer-grained partitions in MBCT and Histogram Bining on Porto Seguro.
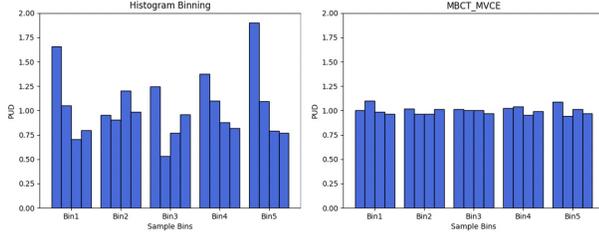


Figure 11: PUD of the finer-grained partitions in MBCT and Histogram Bining on Avazu.
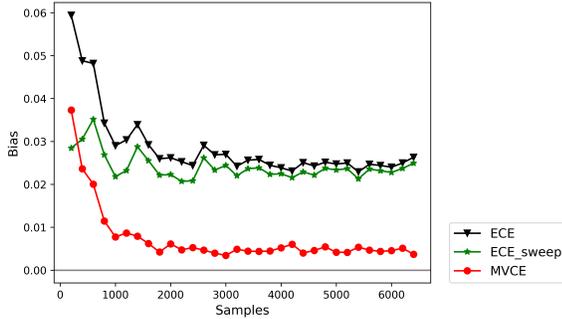


Figure 12: Main simulation results of ECE, ECE$_{sweep}$ and MVCE, in which $h(X) \sim \text{Beta}(0.4, 0.7)$ and $\mathbb{E}[Y|h(X) = c] = c^2$.

Figures 12 and 13 show the simulation results under $h(X) \sim \text{Beta}(0.4, 0.7)$ and $\mathbb{E}[Y|h(X) = c] = c^2$. Figure 14 and 15 show the simulation results under $h(X) \sim \text{Beta}(0.6, 0.7)$ and $\mathbb{E}[Y|h(X) = c] = c^3$. We see that the experimental results under these settings are consistent with the simulation experiment in the main text (in Section 5.3), which further indicates that MVCE is a more appropriate metric for calibration.

## C DISCUSSIONS OF CALIBRATION'S IMPACT ON ONLINE ADVERTISING SYSTEMS

For online real-time bidding advertising platforms, the revenue of the platform is closely related to pCTR (predicted click-through rate). Take the simplest click-to-pay advertisement as an example, advertisers give their *bid* on per click, the system allocate each chance of page-view to one advertiser and charge based on his *bid*
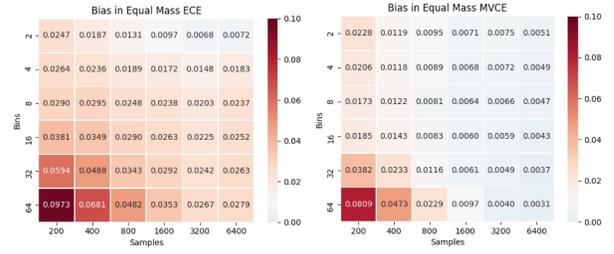


Figure 13: $\hat{E}_{bias}$ of MVCE and ECE under different numbers of bins and samples, in which $h(X) \sim \text{Beta}(0.4, 0.7)$ and $\mathbb{E}[Y|h(X) = c] = c^2$.
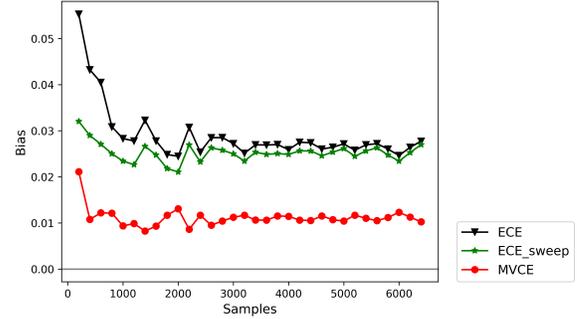


Figure 14: Main simulation results of ECE, ECE$_{sweep}$ and MVCE, in which $h(X) \sim \text{Beta}(0.6, 0.7)$ and $\mathbb{E}[Y|h(X) = c] = c^3$.
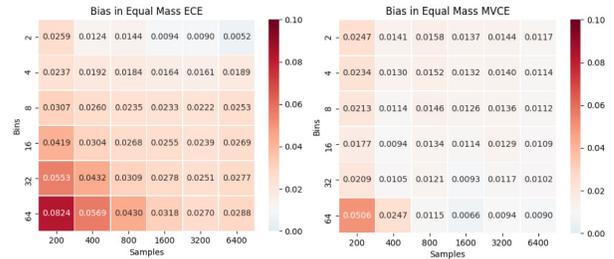


Figure 15: $\hat{E}_{bias}$ of MVCE and ECE under different numbers of bins and samples, in which $h(X) \sim \text{Beta}(0.6, 0.7)$ and $\mathbb{E}[Y|h(X) = c] = c^3$.

when the ad is clicked. The system should allocate the chance of a page-view to the advertiser whose ad achieves the highest score of $bid * pCTR$ to maximize its revenue. Thus, the order accuracy of $bid * pCTR$ score of the ads is crucial for the revenue of the platform. The *bid* is given by advertisers, so the issue becomes to predict the true probability of the click-through rate. Therefore, calibration is quite important for real-time bidding advertising systems. The results of online experiments in the main text (Section 5.6) provides empirical evidence that MBCT brings great business value to our online advertising system, even if the improvement of offline MVCE is seeming small (about 0.8% relative improvement).