

T-Cal: An Optimal Test for the Calibration of Predictive Models

Donghwan Lee*

DH7401@SAS.UPENN.EDU

*Graduate Group in Applied Mathematics and Computational Science
University of Pennsylvania
Philadelphia, PA 19104-6340, USA*

Xinmeng Huang*

XINMENGH@SAS.UPENN.EDU

*Graduate Group in Applied Mathematics and Computational Science
University of Pennsylvania
Philadelphia, PA 19104-6340, USA*

Hamed Hassani

HASSANI@SEAS.UPENN.EDU

*Department of Electrical and Systems Engineering
University of Pennsylvania
Philadelphia, PA 19104-6340, USA*

Edgar Dobriban

DOBRIBAN@WHARTON.UPENN.EDU

*Department of Statistics and Data Science
University of Pennsylvania
Philadelphia, PA 19104-6340, USA*

Editor: Ambuj Tewari

Abstract

The prediction accuracy of machine learning methods is steadily increasing, but the calibration of their uncertainty predictions poses a significant challenge. Numerous works focus on obtaining well-calibrated predictive models, but less is known about reliably assessing model calibration. This limits our ability to know when algorithms for improving calibration have a real effect, and when their improvements are merely artifacts due to random noise in finite datasets. In this work, we consider detecting mis-calibration of predictive models using a finite validation dataset as a hypothesis testing problem. The null hypothesis is that the predictive model is calibrated, while the alternative hypothesis is that the deviation from calibration is sufficiently large.

We find that detecting mis-calibration is only possible when the conditional probabilities of the classes are sufficiently smooth functions of the predictions. When the conditional class probabilities are Hölder continuous, we propose *T-Cal*, a minimax optimal test for calibration based on a debiased plug-in estimator of the ℓ_2 -Expected Calibration Error (ECE). We further propose *adaptive T-Cal*, a version that is adaptive to unknown smoothness. We verify our theoretical findings with a broad range of experiments, including with several popular deep neural net architectures and several standard post-hoc calibration methods. T-Cal is a practical general-purpose tool, which—combined with classical tests for discrete-valued predictors—can be used to test the calibration of virtually any probabilistic classification method. T-Cal is available at <https://github.com/dh7401/T-Cal>.

*. Equal Contribution.

Keywords: uncertainty quantification, calibration, nonparametric statistics, hypothesis testing, minimax optimality

1 Introduction

The prediction accuracy of contemporary machine learning methods such as deep neural networks is steadily increasing, leading to adoption in more and more safety-critical fields such as medical diagnosis (Esteva et al., 2017), self-driving vehicles (Bojarski et al., 2016), and recidivism forecasting (Berk, 2017). In these applications and beyond, machine learning models are required not only to be accurate but also to be well-calibrated: giving precise probability estimates for the correctness of their predictions.

To be concrete, consider a classification problem where the goal is to classify features \mathbf{x} (such as images) into one of several classes \mathbf{y} (such as a building, vehicle, etc.). A probabilistic classifier (or, probability predictor) f assigns to each input \mathbf{x} a probability distribution $f(\mathbf{x})$ over the classes. For a given input \mathbf{x} , the entries of $f(\mathbf{x})$ represent the probabilities assigned by the classifier to the event that the outcome belongs to the k -th class, for any $k = 1, \dots, K$. This classifier is *calibrated* if for any value \mathbf{z} taken by $f(\mathbf{x})$, and for all classes k , the probability that the outcome belongs to the k -th class, i.e., $[\mathbf{y}]_k = 1$, equals the predicted probability, i.e., the k -th coordinate $[\mathbf{z}]_k$ of \mathbf{z} :

$$P([\mathbf{y}]_k = 1 | f(\mathbf{x}) = \mathbf{z}) = [\mathbf{z}]_k.$$

This form of calibration is an important part of uncertainty quantification, decision science, analytics, and forecasting (see e.g., Hilden et al., 1978; Miller et al., 1991, 1993; Steyerberg et al., 2010; Hand, 1997; Jolliffe and Stephenson, 2012; Van Calster and Vickers, 2015; Harrell, 2015; Tetlock and Gardner, 2016; Shah et al., 2018; Steyerberg, 2019, etc). Unfortunately, however, recent works starting from at least Guo et al. (2017) have reported that modern machine learning methods are often poorly calibrated despite their high accuracy; which can lead to harmful consequences (e.g., Van Calster and Vickers, 2015; Steyerberg, 2019).

To address this problem, there has been a surge of works aimed at improving the calibration of machine learning models. These methods seek to achieve calibration either by modifying the training procedure (Harrell, 2015; Lakshminarayanan et al., 2017; Kumar et al., 2018; Thulasidasan et al., 2019; Zhang et al., 2020; Mukhoti et al., 2020) or by learning a re-calibration function that transforms, in a post-hoc way, the predictions to well-calibrated ones (Cox, 1958; Mincer and Zarnowitz, 1969; Steyerberg et al., 2010; Platt, 1999; Zadrozny and Elkan, 2001, 2002; Guo et al., 2017; Kumar et al., 2019; Kisamori et al., 2020).

In this regard, a key challenge is to rigorously assess and compare the performance of calibration methods. Without such assessments, we have limited ability to know when algorithms for improving calibration have a real effect, and when their improvements are merely artifacts due to random noise in finite-size datasets. As it turns out, existing works do not offer a satisfactory solution to this challenge.

In more detail, in this work, we consider the problem of detecting mis-calibration of predictive models using a finite validation dataset. We focus on models whose probability predictions are continuously distributed—which is generally reasonable for many modern

machine learning methods, including deep neural nets. We develop efficient and provably optimal algorithms to test their calibration.

Detecting mis-calibration has been studied from the perspective of statistical hypothesis testing. The seminal work of Cox (1958) formulated a test of calibration for a collection of binary (yes-no) predictions, and proposed using a score test for a logistic regression model. This has been widely used and further developed, leading to various tests for the so-called calibration slope and calibration intercept, which can validate various qualitative versions of model calibration, see e.g., Hosmer and Lemeshow (1980); Miller et al. (1991); Steyerberg (2019) and references therein. In pioneering work, Miller (1962) suggested a chi-squared test for testing calibration of multiple series of binary predictions. To deal with the challenging problem of setting critical values (i.e., how large of an empirical mis-calibration is statistically significant?) for testing calibration, bootstrap methods have become common, see e.g., Harrell (2015). We refer to Section 1.1 for more details and for a discussion of other related works.

In contrast to the above works that aim to test calibration slopes and intercepts, we aim to develop a *nonparametric hypothesis test* for calibration, which does not assume a specific functional form (such as a logistic regression model), for the deviations to be detected from perfect calibration. A nonparametric approach has the advantage that it can detect subtle forms of mis-calibration even after re-calibration by parametric methods. However, existing approaches for nonparametric testing often rely on ad hoc techniques for binning the probability predictions, which is a limitation because the results can depend on the way that the binning has been performed (Harrell, 2015; Steyerberg, 2019). In contrast, our *adaptive tests* automatically select an optimal binning scheme. Finally, as a new development in the area of testing calibration, T-Cal has theoretically guaranteed *minimax optimality properties* for detecting certain reasonable types of smooth mis-calibration. These properties make T-Cal both practically and theoretically appealing.

We consider a given multi-class probabilistic classifier, and are interested in testing if it is calibrated. We make the following contributions:

- As a candidate test statistic, we consider the plug-in estimator of ℓ_2 -expected calibration error (ECE), which is the expectation of the squared distance between the probability predictions and class probabilities given these predictions. This is also known as the mean calibration error (e.g., Harrell, 2015, p. 105). While the plug-in estimator is biased (i.e., its expectation is not zero even under perfect calibration), we show how to construct a *debiased plug-in estimator* (DPE).

We consider detecting mis-calibration when the deviation between predicted class probabilities and their true values—the “mis-calibration curve”—satisfies a classical smoothness condition known as Hölder continuity. We later show that a smoothness condition is essentially unavoidable. Under this condition, we show that T-Cal can detect mis-calibration if the ECE is sufficiently large and the number of bins is chosen appropriately, depending on the smoothness (Theorem 3).

- To make T-Cal practical, we present a version that is adaptive to the unknown smoothness parameter (Theorem 5). This makes T-Cal fully tuning-free and practically useful. From a theoretical perspective, adaptivity only requires a minor additional increase in the level of mis-calibration that can be detected; by a log n factor.

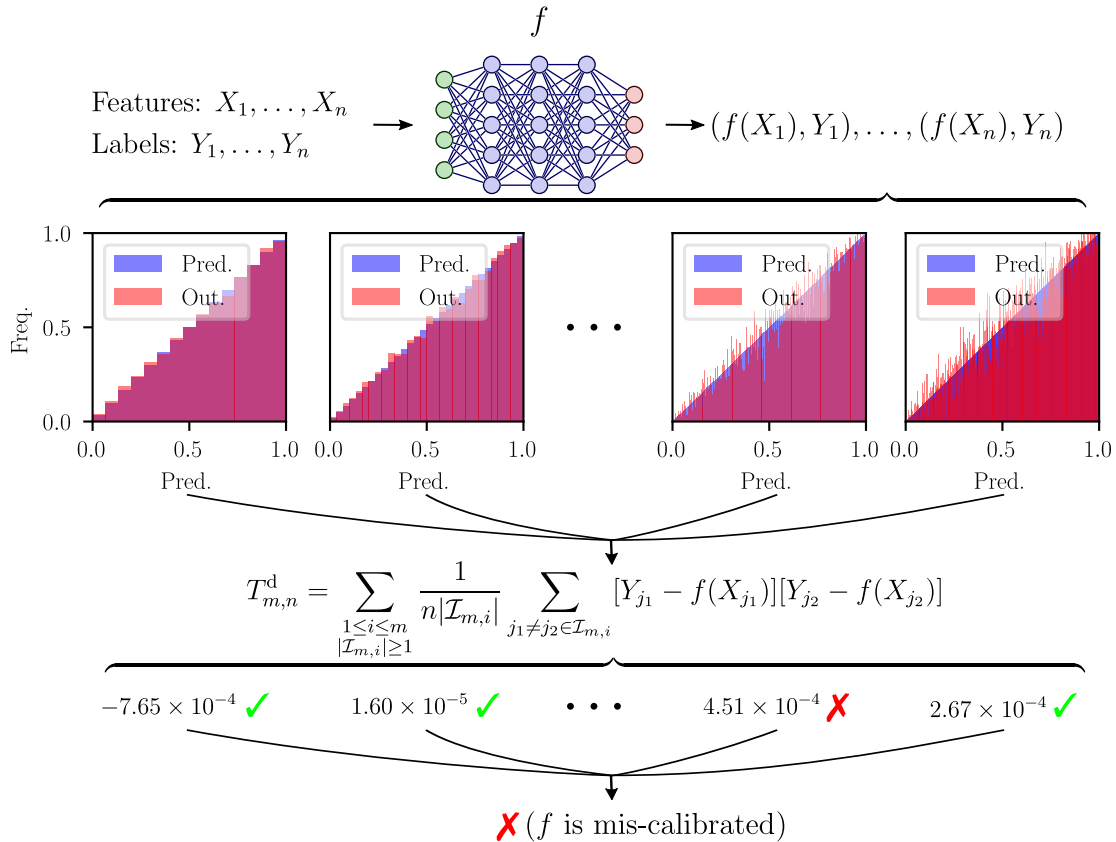


Figure 1: An overview of adaptive T-Cal. For a given probability predictor f , we compute $T_{m;n}^d$, the debiased plug-in estimator (DPE), binned over several scales (See (6) for the definition). We then compare each value with the hypothetical distribution of DPE that we would get if the model were perfectly calibrated. The hypothesis of perfect calibration is rejected if at least one of the scales is detected to be mis-calibrated. This multi-scale approach ensures that T-Cal adaptively detects mis-calibration.

- We support our theoretical results with a broad range of experiments. We provide simulations, which support our theoretical optimality results. We also provide experiments with several popular deep neural net architectures (ResNet-50, VGG-19, DenseNet-121, etc), on benchmark datasets (CIFAR 10 and 100, ImageNet) and several standard post-hoc calibration methods (Platt scaling, histogram binning, isotonic regression, etc).
- To complement these results, we argue that T-Cal is optimal, by providing a number of fundamental lower bounds. We prove that detecting mis-calibration from a finite dataset is only possible when the mis-calibration curve is sufficiently smooth, and it is not possible when the curve is just continuous (Proposition 9).

When the mis-calibration curves are Hölder smooth, we show that the calibration error required for reliable detection of mis-calibration has to be appropriately large (Theorem 10). This minimax result relies on Ingster’s (or the chi-squared) method. Combined with our previous results, this shows that T-Cal is minimax optimal.

- To further put our problem in context, we show that testing calibration can be reduced to a well-known problem in statistical inference—the two-sample goodness-of-fit problem—by a novel randomization technique. Based on this insight, and building on the results of Arias-Castro et al. (2018); Kim et al. (2022) on goodness-of-fit testing, we present another asymptotically minimax optimal test for mis-calibration that matches the lower bound (Theorem 12). While this method is theoretically optimal, it relies on sample splitting and is not as sample-efficient as our previous method in experiments.
- In the proofs, we have the following innovations:
 1. We introduce an *equal-volume binning scheme* for the probability simplex Δ_{K-1} (Appendix B.3). We decompose the probability simplex into hypersimplices by taking intersections with smaller hypercubes composing the unit hypercube $[0, 1]^K$. Then we further decompose the hypersimplices into equal-volume simplices using results in polyhedral combinatorics. This construction enables us to extend proof techniques from the nonparametric hypothesis testing literature to our setting.
 2. To analyze our plug-in estimator, we need to deal with terms involving probability scores of inputs, which are continuous random variables. This is different from the structure of chi-squared statistics such as that of Ingster (1987). Thus, computing the mean and variance of the DPE requires a different analysis.
 3. While densities on the probability simplex can take arbitrary positive values, the conditional expectation of probability predictions has to lie in the probability simplex. This requires a careful construction of alternative distributions to use Ingster’s method.

Our numerical results can be reproduced with code available at <https://github.com/dh7401/T-Cal>.

We now summarize some key takeaways:

- **The need for statistical significance to claim calibration.** It is crucial to perform rigorous statistical tests to assess the calibration of machine learning methods. While models with smaller empirical ECE generally tend to be better calibrated, these values can be highly influenced by noise and randomness inherent in finite datasets. Hence it is crucial to develop and use tools to assess statistical significance—such as the hypothesis tests of calibration that we develop—when claiming improved calibration.
- **Potential suboptimality of popular approaches.** The currently prevalent usage of popular metrics, such as the empirical ECE, may be suboptimal. The current standard is to evaluate mis-calibration metrics using a fixed number of bins (such as 15) of the probability scores, for all prediction models (ResNet, VGG, etc), and all datasets. Our results show theoretically that the optimal number of bins increases with the level of oscillations and non-smoothness expected in the probability predictor. Modern machine learning methods are becoming more and more over-parametrized and data-adaptive. This suggests that it is ever more important to use a careful model- and data-adaptive test (and number of bins) when testing calibration.

1.1 Related Works

There is a great body of related work on evaluating the calibration of prediction methods, on improving calibration accuracy, and on nonparametric hypothesis testing techniques. We review the most closely related works.

Broader Context. Broadly speaking, the study of calibration is an important part of the study of classification, prediction, analytics, and forecasting (e.g., Hilden et al., 1978; Miller et al., 1991, 1993; Steyerberg et al., 2010; Hand, 1997; Jolliffe and Stephenson, 2012; Gneiting and Katzfuss, 2014; Van Calster and Vickers, 2015; Harrell, 2015; Tetlock and Gardner, 2016; Shah et al., 2018; Steyerberg, 2019, etc).

Calibration. As recounted in Lichtenstein et al. (1977), research on calibration dates back at least to the early 1900s, when meteorologists suggested expressing predictions as probabilities and comparing them to observed empirical frequencies. Calibration has since been studied in a variety of areas, including meteorology, statistics, medicine, computer science, and social science; and under a variety of names, such as realism or realism of confidence, appropriateness of confidence, validity, external validity, secondary validity, and reliability (Lichtenstein et al., 1977). A general finding in this area is that human forecasters are often overconfident and thus mis-calibrated (e.g., Keren, 1991, etc), as codified for instance in Tversky and Kahneman’s celebrated work on prospect theory (Kahneman and Tversky, 2013).

Beyond our hypothesis testing perspective, approaches to study calibration include Bayesian perspectives (e.g., Dawid, 1982; Kadane and Lichtenstein, 1982, etc.) and on-line settings (e.g., Foster and Vohra, 1998; Vovk and Shafer, 2005, etc.). See also Section 10.9 of Harrell (2015), Section 15.3 of Steyerberg (2019), and Hastie and Tibshirani (1998); Ivanov et al. (1999); Garczarek (2002); Buja et al. (2005); Toll et al. (2008); Gebel (2009); Serrano (2012); Van Calster et al. (2019); Huang et al. (2020), among others.

Calibration Measures. Proper scoring rules (Good, 1952; De Finetti, 1962; Savage, 1971; Winkler et al., 1996; DeGroot and Fienberg, 1983; Gneiting et al., 2007) such as

the Brier score (Brier, 1950) and negative log-likelihood (e.g., Winkler et al., 1996, etc) are objective functions of two probability distributions (a true distribution and a predicted distribution). They are minimized when the predicted distribution equals the true distribution; see also Bickel (2007).

As discussed in Sections 4.5 and 10.9 of Harrell (2015), some of the standard techniques in the area include plotting calibration curves, also known as reliability diagrams (estimated probabilities against predicted ones); which can be bias-corrected using the bootstrap; and re-calibration by fitting statistical models to these curves (e.g., Austin and Steyerberg, 2014, etc). More recently, the notion of ECE, also known as mean absolute calibration error (e.g., Harrell, 2015, p. 105) is popularized in Naeini et al. (2015) and later generalized to multi-class settings in Vaicenavicius et al. (2019). Gupta et al. (2021) develop a binning-free calibration measure based on the Kolmogorov-Smirnov test. Arrieta-Ibarra et al. (2022) introduce calibration metrics building on Komogorov-Smirnov and Kuiper statistics.

Calibration in Modern Machine Learning. Guo et al. (2017) draw attention to the mis-calibration of modern neural networks and compare different recalibration methods based on ECE. Many other works (Milius et al., 2018; Kull et al., 2019; Zhang et al., 2020) also evaluate their methods using ECE or its variants. In the following works of Vaicenavicius et al. (2019); Kumar et al. (2019), it has been recognized that ECE evaluated on a fixed binning scheme can underestimate the calibration error. The limitation of fixed binning has been known for the analogous problems of testing probability distributions and densities, see e.g., Mann and Wald (1942), or page 19 of Ingster (2012). Kumar et al. (2019) proposes a debiased ECE, but only for probability predictors with a finite number of outputs. Nixon et al. (2019) empirically study various versions of ECE obtained by adjusting hyperparameters involved in the estimator of ECE (such as norm, binning scheme, and class conditionality), and find that the choice of calibration measure is crucial when comparing different calibration methods. ? propose a heuristic for choosing an optimal number of bins when computing ECE. Zhang et al. (2020) use kernel density estimation to estimate ECE without relying on a binning scheme. Zhao et al. (2020) show that individual calibration is possible by randomized predictions and propose a training objective to enforce individual calibration. See also Niculescu-Mizil and Caruana (2005); Kull et al. (2017); Bai et al. (2021), among others.

There has been interest in a variety of forms of calibration. We study the strongest form, multi-class calibration, which is stronger than other definitions such as marginal calibration and confidence (top) calibration (Vaicenavicius et al., 2019; Widmann et al., 2019).

Nonparametric Hypothesis Testing. Ingster (1986) derives the minimax testing rate for two-sample testing where s -Hölder continuous densities on $[0, 1]$ are separated in an L^2 sense, and shows that the chi-squared test achieves the minimax optimal rate $n^{-2s/(4s+1)}$. Ingster (1987) extends this result to L^p metrics and derives the minimax optimal rate $n^{-s/(2s+1-\max\{2,p\})}$ for $1 < p < \infty$ and $(n/\log n)^{s/(2s+1)}$ for $p = 1$. Ingster (2000) proposes an adaptive version of the test at the cost of $(\log \log n)^{s/(4s+1)}$ factor in the minimax rate. Arias-Castro et al. (2018) extend these results to densities on $[0, 1]^d$ and show the minimax rate $n^{-2s/(4s+d)}$. Kim et al. (2022) prove that a permutation test can also achieve the same optimal rate. Butucea and Tribouley (2006) study two-sample testing for one-dimensional densities in Besov spaces; they also prove adaptivity. See also Balakrishnan

and Wasserman (2018); Donoho and Jin (2015); Jin and Ke (2016); Chhor and Carpentier (2021); Dubois et al. (2021); Berrett et al. (2021) for reviews and further related works.

Nonparametric Functional Estimation. Bickel and Ritov (1988) study the problem of estimating the quadratic integral functional of the k -th derivative of s -Hölder probability densities on \mathbb{R} , and prove that the optimal convergence rate is $n^{-[4(s-k)=(4s+1)\wedge 1=2]}$. Donoho and Nussbaum (1990); Brown and Low (1996) show an analogous result for regression functions in the Gaussian white noise model. Birgé and Massart (1995) generalize these results to smooth integral functionals of Hölder densities and their derivatives, and prove the same convergence rate. Kerkycharian and Picard (1996) provide optimal Haar wavelet-based estimators of cubic functionals of densities over the broader class of Besov spaces; and also discuss estimating integrals of other powers of the density. Laurent (1996) studies estimation of functionals of the form $\int \phi(f(x), x) d\mu(x)$ of densities f , where ϕ is a sufficiently smooth function and μ is a measure. This work constructs estimators attaining the optimal parametric rate using orthogonal projections, including showing semiparametric efficiency, when the smoothness $s > d/4$ in dimension d . Robins et al. (2008); Giné and Nickl (2008); Tchetgen et al. (2008) introduce estimation method using higher-order U-statistics. Efro-movich and Low (1996); Cai and Low (2006); Giné and Nickl (2008); Mukherjee et al. (2015) propose estimation methods adaptive to unknown smoothness s based on Lepski’s method (Lepski, 1991; Lepski and Spokoiny, 1997). See Giné and Nickl (2021) for a more thorough review of related literature.

Hypothesis Testing for Calibration. Cox (1958) formulates a test of calibration for a collection of Bernoulli random variables, as a test that their success probabilities are equal to some given values; and proposed using a score test for a logistic regression model. These tests are referred to as testing the calibration slope and intercept, and they are part of a broader hierarchy of calibration (Van Calster et al., 2016). See also Miller et al. (1991); Steyerberg (2019) and references therein. Miller (1962), Section 5, suggests a chi-squared test for testing calibration of a collection of sequences of Bernoulli random variables. Spiegelhalter (1986) proposes a test of calibration based on the Brier score, for discrete-valued probability predictors. The Hosmer-Lemeshow test (Hosmer and Lemeshow, 1980) is a goodness-of-fit test for logistic regression models. The test is based on a chi-squared statistic that measures differences between expected and observed numbers of events in subgroups, and thus has, on the surface, a similarity to the types of test statistics we consider. There are also related tests for comparing predictors (Schervish, 1989; Diebold and Mariano, 1995).

Seillier-Moiseiwitsch and Dawid (1993) study testing the calibration of sequential probability forecasts. Bröcker and Smith (2007) study the bootstrap-based procedure they call consistency resampling to produce standard error bars in reliability diagrams; without focusing on its optimality. For testing the calibration of forecasted densities, Dawid (1984); Diebold et al. (1998) propose the probability integral transform (PIT). Held et al. (2010) propose a score-based approach for testing calibration. Vaicenavicius et al. (2019) use consistency resampling to test a hypothesis of perfect calibration; again without studying its optimality. Widmann et al. (2019) propose kernel-based mis-calibration measures together with their estimators, and argue that the estimators can be viewed as calibration test statistics. Tamás and Csáji (2021) suggest distribution-free hypothesis tests for the null $H_0 : E[Y | X] = X$ based on conditional kernel mean embedding.

Note on Terminology. The term calibration sometimes has a different meaning in a variety of areas of human activity, including measurement technology, engineering, economics, and even statistics, etc., see e.g., Franklin (1999); Dawkins et al. (2001); Kodovský and Fridrich (2009); Osborne (1991); Vovk et al. (2020); Angelopoulos et al. (2021). These generally mean adjusting a measurement to agree with a desired standard, within a specified accuracy. However, in our work, we focus on the notion of probabilistic calibration described so far.

1.2 Notations

For an integer $d \geq 1$, and a vector $\mathbf{v} \in \mathbb{R}^d$, we refer to the coordinates of \mathbf{v} as both $[\mathbf{v}]_1, \dots, [\mathbf{v}]_d$ and v_1, \dots, v_d . For any $p \geq 1$, and for an integer $K \geq 2$, we denote the ℓ_p -norm of $\mathbf{x} = (x_1, \dots, x_K)^\top \in \mathbb{R}^K$ by $\|\mathbf{x}\|_p := (\sum_{i=1}^K |x_i|^p)^{1/p}$. When p is unspecified, $\|\cdot\|$ stands for $\|\cdot\|_2$. For an event A , we denote by $I(A)$ its indicator random variable, where $I(A) = 1$ if event A happens, and $I(A) = 0$ otherwise. For two real numbers a, b , we denote $a \wedge b := \min(a, b)$. For two sequences $(a_n)_{n \geq 1}$ and $(b_n)_{n \geq 1}$ with $b_n \neq 0$, we write $a_n \sim b_n$ if $0 < \liminf_n a_n/b_n \leq \limsup_n a_n/b_n < \infty$. When the index n is self-evident, we may omit it above. We use the Bachmann-Landau asymptotic notations $\Omega(\cdot), \Theta(\cdot)$ to hide constant factors in inequalities and use $\tilde{\Omega}(\cdot), \tilde{\Theta}(\cdot)$ to also hide logarithmic factors. For a Lebesgue measurable set $A \subseteq \mathbb{R}^d$, we denote by $1_A : \mathbb{R}^d \rightarrow \{0, 1\}$ its indicator function where $1_A(\mathbf{x}) = 1$ if $\mathbf{x} \in A$ and $1_A(\mathbf{x}) = 0$ otherwise. For a real number $s \in \mathbb{R}$, we denote the largest integer less than or equal to s by $\lfloor s \rfloor$. Also, the smallest integer greater than or equal to s is denoted by $\lceil s \rceil$.

2 Definitions and Setup

For $K \geq 2$, consider a K -class classification problem where $X \in \mathcal{X}$ is the input feature vector (for instance, an image) and $Y \in \mathcal{Y} := \{y\} = \{y_1, \dots, y_K\}^\top \in \{0, 1\}^K : \sum_{i=1}^K y_i = 1$ is the one-hot encoded output label (for instance, the indicator of the class of the image: building, vehicle, etc).

We consider a probabilistic classifier f mapping the feature space to probability distributions over K classes. Formally, the output space is the $(K - 1)$ -dimensional probability simplex Δ_{K-1} ,

$$\Delta_{K-1} := \{\mathbf{z} = (z_1, \dots, z_K)^\top \in [0, 1]^K : z_1 + \dots + z_K = 1\},$$

i.e., $f : \mathcal{X} \rightarrow \Delta_{K-1}$. For any $k \in \{1, \dots, K\}$, the individual component $[f(X)]_k$ denotes the predicted probability of the k -th class. Thus, f is also referred to as a probability predictor. The probability predictor f is assumed to be pre-trained on data that are independent of our calibration data at hand.

We assume that the feature-label pair (X, Y) has an unknown joint probability distribution P on $\mathcal{X} \times \mathcal{Y}$. Calibration requires that the predicted probabilities of correctness are equal to the true probabilities. Thus, given that we predicted the probabilities $f(X) = \mathbf{z}$, and thus $[f(X)]_k = \mathbf{z}_k$, the true probability that $[Y]_k = 1$ should be equal to $[f(X)]_k = \mathbf{z}_k$. Thus, for almost every \mathbf{z} , calibration requires that for all $k = 1, \dots, K$,

$$P[[Y]_k = 1 \mid f(X) = \mathbf{z}] = \mathbf{z}_k.$$

We can reformulate this in a way that is more convenient to study. The map $(\mathbf{x}, \mathbf{y}) \mapsto (f(\mathbf{x}), \mathbf{y})$ induces a probability distribution on $\Delta_{K-1} \times Y$; where we can think of (\mathbf{x}, \mathbf{y}) as a realization of (X, Y) . As will be discussed shortly, calibration only depends on the joint distribution of $(f(X), Y)$. For this reason, we also denote the joint distribution of $(f(X), Y)$ by P when there is no confusion. We write $Z := f(X)$ for the predicted probabilities corresponding to X .

We define the *regression function* $\text{reg}_f : \Delta_{K-1} \rightarrow \Delta_{K-1}$ as

$$\text{reg}_f(\mathbf{z}) := \mathbb{E}[Y \mid f(X) = \mathbf{z}] = \mathbb{E}[Y \mid Z = \mathbf{z}],$$

where the expectation is conditioned on the score Z with $(Z, Y) \sim P$. Note that each component, for $k = 1, \dots, K$, has the form $\mathbb{E}[[Y]_k \mid f(X) = \mathbf{z}] = P[[Y]_k = 1 \mid f(X) = \mathbf{z}]$. Especially for binary classification, this is also referred to as the *calibration curve* of the probabilistic classifier f (Harrell, 2015). Since we are particularly interested in continuous probability predictors, we assume that the marginal distribution P_Z of Z has a density with respect to the uniform measure on Δ_{K-1} . Then, this expectation is well-defined almost everywhere.

In this language, the probabilistic classifier f is *perfectly calibrated* if $\text{reg}_f(Z) = Z$ almost everywhere.¹ Further, it turns out that it is important to study *the deviations from calibration*. For this reason, we define the *residual function* $\text{res}_f : \Delta_{K-1} \rightarrow \mathbb{R}^K$ as

$$\text{res}_f(\mathbf{z}) := \text{reg}_f(\mathbf{z}) - \mathbf{z},$$

so that perfect calibration amounts to $\text{res}_f(Z) = 0$ almost everywhere. When (Z, Y) have a joint distribution P , we sometimes write $\text{res}_f = \text{res}_{f,P}$ to display the dependence of the mis-calibration curve on P . As we will see, the structure of the residual function crucially determines our ability to detect mis-calibration. In analogy to the notion of calibration curves mentioned above, we may also call res_f the *mis-calibration curve* of the probabilistic classifier f .

We observe calibration data $(Z_i, Y_i) \in \Delta_{K-1} \times Y, i \in \{1, \dots, n\}$, sampled i.i.d. from P , and denote their joint product distribution as P^n . Our goal is to rigorously test if f is perfectly calibrated based on this finite calibration dataset. The calibration properties of the probabilistic classifier f can be expressed equivalently in terms of the distribution P of $(f(X), Y) = (Z, Y)$. Therefore, we will sometimes refer to testing the calibration of the distribution P , and the probabilistic classifier will be implicit.

Expected Calibration Error. The ℓ_p -ECE (Expected Calibration Error) for the distribution P , also known as the mean calibration error (e.g., Harrell, 2015, p. 105), is

$$\ell_p\text{-ECE}(f) = \ell_p\text{-ECE}_P(f) = \mathbb{E}_{Z \sim P_Z} \left[\sum_{k=1}^K |k \text{reg}_f(Z) - Z| k^{\frac{1}{p}} \right] = \mathbb{E}_{Z \sim P_Z} \left[\sum_{k=1}^K |k \text{res}_f(Z)| k^{\frac{1}{p}} \right]. \quad (1)$$

1. In the binary case ($K = 2$), we identify Δ_{K-1} with $[0, 1]$ via the map $(z; 1 - z) \mapsto z$ and use $Y = \{0, 1\}$ instead of the one-hot encoded output space. We say f is perfectly calibrated if $\text{reg}_f(z) := P(Y = 1 \mid f(X) = z) = z$ almost everywhere.

In words, this quantity measures the average over all classes $k = 1, \dots, K$ and over the data distribution $X \sim P_X$ of the per-class error $[\text{res}_f(\mathbf{z})]_k = \mathbb{E}[[Y]_k \mid f(X) = \mathbf{z}] - [\mathbf{z}]_k$ between the predicted probability of class k for input X —i.e., $[\mathbf{z}]_k = [f(X)]_k$ —and the actual probability $\mathbb{E}[[Y]_k \mid f(X) = \mathbf{z}] = P[[Y]_k = 1 \mid f(X) = \mathbf{z}]$ of that class. For instance, when the number of classes is $K = 2$, and the power is $p = 1$, we have $\ell_1\text{-ECE}(f) = \mathbb{E}_{X \sim P_X} \frac{1}{2} \sum_{k=1}^2 |P[[Y]_k = 1 \mid f(X)] - [f(X)]_k| = 2 \mathbb{E}_{X \sim P_X} |P[[Y]_1 = 1 \mid f(X)] - [f(X)]_1|$.

Hölder Continuity. We describe the notion of Hölder continuity for functions defined on Δ_{K-1} . For simplicity, we only provide the definition for $K = 2$. See Appendix B.1 for the complete definition for general $K \geq 2$.

Identifying Δ_1 with $[0, 1]$ via the map $(z, 1 - z)^\top \mapsto z$, a function $g : \Delta_1 \rightarrow \mathbb{R}$ can be equivalently understood as a function $g : [0, 1] \rightarrow \mathbb{R}$. For an integer $d \geq 0$ and a function $g : [0, 1] \rightarrow \mathbb{R}$, let $g^{(d)}$ be the d -th derivative of the function g . For a real number s , we denote the smallest integer greater than or equal to s by $\lceil s \rceil$.

For a Hölder smoothness parameter $s > 0$ and a Hölder constant $L > 0$, let $H_K(s, L)$ be the class of (s, L) -Hölder continuous functions $g : [0, 1] \rightarrow \mathbb{R}$ satisfying, for all $x_1, x_2 \in [0, 1]$

$$|g^{(\lceil s \rceil - 1)}(x_1) - g^{(\lceil s \rceil - 1)}(x_2)| \leq L |x_1 - x_2|^{s - \lceil s \rceil + 1}. \quad (2)$$

In particular, $H_K(1, L)$ denotes all L -Lipschitz functions. We consider $L > 0$ as an arbitrary fixed constant, and we do not display the dependence of our results on its value. For instance, when the Lipschitz constant is $L = 1$, and the Hölder smoothness parameter is $s = 1.5$, this is the set of real-valued functions g defined on $[0, 1]$ such that for all $x_1, x_2 \in [0, 1]$, $|g'(x_1) - g'(x_2)| \leq L |x_1 - x_2|^{0.5}$.

Goal. Our goal is to test the null hypothesis of perfect calibration, i.e., $\text{res}_f = 0$, against the alternative hypothesis that the model is mis-calibrated. To quantify mis-calibration, we use the notion of the ℓ_p -ECE(f) from (1). We study the signal strength needed so that reliable mis-calibration detection is possible. Further, we assume that the mis-calibration curves are Hölder continuous because we will show that by only assuming continuity, reliable detection of mis-calibration is impossible. In Remark 11, We will also discuss what happens when the mis-calibration function is not Hölder smooth.

Let \mathcal{P} be the family of all distributions P over $(Z, Y) \in \Delta_{K-1} \times \mathcal{Y}$ such that the marginal distribution P_Z of Z has a density with respect to the uniform measure on Δ_{K-1} . Define the collection \mathcal{P}_0 of joint distributions P of (Z, Y) under which the probability predictor f is perfectly calibrated:

$$\mathcal{P}_0 := \{P \in \mathcal{P} : \text{res}_{f,P}(Z) = 0, P_Z\text{-a.s.}\}.$$

For a Hölder smoothness parameter s and a Hölder constant L , let $\mathcal{P}_{s;L;K}$ be the family of probability distributions $P \in \mathcal{P}$ over the predictions and labels $(f(X), Y) = (Z, Y) \in \Delta_{K-1} \times \mathcal{Y}$ under which the residual map $\mathbf{z} \mapsto [\text{res}_{f,P}(\mathbf{z})]_k$ (i.e., the map res_f under the distribution $(Z, Y) \sim P$) belongs to the class of (s, L) -Hölder continuous functions $H_K(s, L)$ for every $k \in \{1, \dots, K\}$. For a separation rate $\varepsilon > 0$, define the collection $\mathcal{P}_1(\varepsilon, p, s)$ of joint distributions $P \in \mathcal{P}_{s;L;K}$ under which the ℓ_p -ECE of f is at least ε :

$$\mathcal{P}_1(\varepsilon, p, s) := \{P \in \mathcal{P}_{s;L;K} : \ell_p\text{-ECE}_P(f) \geq \varepsilon\}. \quad (3)$$

We will also refer to these distributions as ε -*mis-calibrated*. Our goal is to test the *null hypothesis* of calibration against the *alternative* of an ε -calibration error:

$$H_0 : P \in \mathcal{P}_0 \quad \text{versus} \quad H_1 : P \in \mathcal{P}_1(\varepsilon, p, s). \quad (4)$$

Although we consider the null hypothesis of perfect calibration, we generally do not expect a model trained on finite data to be perfectly calibrated. In this regard, the purpose of testing (4) is to check if there is statistically significant evidence of mis-calibration, and *not* to check whether the predictor f is perfectly calibrated. As usual in hypothesis testing, not rejecting the null hypothesis does not mean that we accept that f is perfectly calibrated but means that there is no statistically significant evidence of mis-calibration. In this case, to gain more confidence that the model is calibrated, one may consider testing other hypotheses about calibration—such as top- k calibration, (Guo et al., 2017)—or collecting more data; of course, this may require dealing with multiple testing problems. Meanwhile, since the null of calibration is not rejected, one may use the classifier as if it was calibrated until evidence to the contrary is presented.

Moreover, in Remark 4, we also provide results for the null hypothesis of a small enough calibration error.

Hypothesis Testing. We recall some notions from hypothesis testing (e.g., Lehmann and Romano, 2005; Ingster, 2012, etc) that we use to formulate our problem. A *test* ξ is a function² $\xi : (\Delta_{\mathcal{K}-1} \times \mathcal{Y})^n \rightarrow \{0, 1\}$ of the data, given a dataset $S = (X_i, Y_i)_{i=1}^n \in (\Delta_{\mathcal{K}-1} \times \mathcal{Y})^n$, the decision $\xi(S)$ of rejecting the null hypothesis. In other words, for a given dataset S , $\xi(S) = 1$ means that we detect mis-calibration, and $\xi(S) = 0$ means that we do not detect mis-calibration.

Denote the set of all *level* $\alpha \in (0, 1)$ tests, which have a *false detection rate* (or, *false positive rate*; *type I error*) bounded by α , as

$$\Phi_n(\alpha) := \left\{ \xi : \sup_{P \in \mathcal{P}_0} P(\xi = 1) \leq \alpha \right\}.$$

The probability $P(\xi = 1)$ is taken with respect to the distribution of the sample. For $\varepsilon > 0$ and $P \in \mathcal{P}_1(\varepsilon, p, s)$ from (3), we want to minimize the *false negative rate* (*type II error*) $P(\xi = 0)$, the probability of not detecting mis-calibration. We consider the worst possible value (maximum or rather supremum) $\sup_{P \in \mathcal{P}_1(\varepsilon, p, s)} P(\xi = 0)$ of the type II error, over all distributions $P \in \mathcal{P}_1(\varepsilon, p, s)$. We then want to minimize this over all tests $\xi \in \Phi_n(\alpha)$ that appropriately control the level, leading to the *minimax risk* (minimax type II error)

$$R_n(\varepsilon, p, s) := \inf_{\xi \in \Phi_n(\alpha)} \sup_{P \in \mathcal{P}_1(\varepsilon, p, s)} P(\xi = 0).$$

In words, among all tests that have a false detection rate of $\alpha < 1$ using a sample of size n , we want to find the one with the best possible (smallest) mis-detection rate over all ε -mis-calibrated distributions.

We consider $\alpha \in (0, 1)$ as a fixed constant, and we do not display the dependence of our results on its value. We want to understand how large the ℓ_p -ECE (as measured by

2. To be rigorous, a Borel measurable function.

ε in $\mathcal{P}_1(\varepsilon, p, s)$) needs to be to ensure reliable detection of mis-calibration. This amounts to finding ε' such that the best possible worst-case risk $R_n(\varepsilon', p, s)$ is small. For a fixed $\beta \geq (0, 1 - \alpha)$, the minimum separation (signal strength) for s -Hölder functions, in the ℓ_p -norm, needed for a minimax type II error of at most β is defined as

$$\varepsilon_n(\beta; p, s) := \varepsilon_n(p, s) = \inf \{ \varepsilon' : R_n(\varepsilon', p, s) \leq \beta g \}.$$

Since $\beta \geq (0, 1 - \alpha)$ is fixed, we usually omit the dependence of ε_n on this value.

Remark 1 (Comparison with classical nonparametric hypothesis testing) *As we summarize in Section 1.1, prior works such as Ingster (1987, 2000, 2012); Berman et al. (2014) have studied the problem of testing that the L^p norm of a function is zero against the alternative that it is nonzero, where the function is either a probability density or a regression function in the Gaussian white noise model. Our task here is different from the classical problem since reg_f is not a probability density, and we are not provided independent observations of the function reg_f or res_f in the Gaussian white noise model. Rather, our observation model is closer to multinomial regression; which is heteroskedastic and differs from the above models. While our proposed test shares ideas with the chi-squared test of Ingster (1987, 2000), it requires a different analysis for the above-mentioned reasons.*

3 An Adaptive Debiased Calibration Test

Here we describe our main test for calibration. This relies on a debiased plug-in estimator for $\ell_2\text{-ECE}(f)^2$. We prove that the test is minimax optimal and discuss why debiasing is necessary. We also provide an adaptive plug-in test, which can adapt to an unknown Hölder smoothness parameter s .

3.1 Debiased Plug-in Estimator

The calibration error of a continuous probability predictor f is often estimated by a discretized plug-in estimator associated with a partition (or binning) of the probability simplex Δ_{K-1} (e.g., Cox, 1958; Harrell, 2015). The early work of Cox (1958) already recommended grouping together similar probability forecasts. More recently, Guo et al. (2017) divide the interval $[0, 1]$ into bins of equal width and compute the (top-1) ECE by averaging the difference between confidence and accuracy in each bin. Vaicenavicius et al. (2019) generalize this idea to K -class classification and data-dependent partitions.

In this work, we use an equal-volume partition B_m of the probability simplex Δ_{K-1} , which is parametrized by a binning scheme parameter $m \geq \mathbb{N}_+$. The partition B_m consists of m^{K-1} simplices with equal volumes and diameters proportional to m^{-1} . To construct such a partition, we first divide the simplex Δ_{K-1} into $K - 1$ hypersimplices—generalizations of the standard probability simplex that can have more vertices and edges—by taking intersections with m^{-1} -scaled and translated K -dimensional hypercubes. The hypersimplices are further divided into unit volume simplices using the result of Stanley (1977); Sturmfels (1996). The construction of B_m is elaborated in Appendix B.3. The purpose of using an equal-volume partition B_m is only for a simpler description of our results, and any partition with $\Theta(m^{-K+1})$ volumes and $\Theta(m^{-1})$ diameters can be used.

Let us denote the sets comprising the partition as $B_m = \{B_1; \dots; B_{m^k-1}\}$. For each $i \in \{1; \dots; m^k-1\}$, define the indices of data points falling into the bin B_i as $I_{m;i} := \{j : Z_j \in B_i; 1 \leq j \leq n\}$. Then, for each $i \in \{1; \dots; m^k-1\}$, the averaged difference between probability predictions $Z_j = f(X_j)$ and true labels Y_j for the probability predictions in B_i is $\frac{1}{|I_{m;i}|} \sum_{j \in I_{m;i}} (Y_j - Z_j)$. This estimates $E[Y - Z | Z \in B_i] = E[\text{res}_f(Z) | Z \in B_i]$. Now, the quantity $\int_2\text{-ECE}(f)^2 = E[\text{res}_f(Z)^2]$ can be approximated by piecewise averaging as $\frac{1}{m^k-1} \sum_{i=1}^{m^k-1} P_Z(B_i) kE[\text{res}_f(Z) | Z \in B_i]k^2$. Plugging in the estimate $\frac{1}{|I_{m;i}|} \sum_{j \in I_{m;i}} (Y_j - Z_j)$ of $kE[\text{res}_f(Z) | Z \in B_i]k^2$, we can define a plug-in estimator of $\int_2\text{-ECE}(f)^2$ as follows:

$$T_{m;n}^b := \sum_{i=1}^{m^k-1} \frac{|I_{m;i}|}{n} \frac{1}{|I_{m;i}|} \sum_{j \in I_{m;i}} (Y_j - Z_j)^2 \quad (5)$$

Above, the sum is taken over bins B_i containing at least one datapoint. As will be discussed in Section 3.3, the plug-in estimator is biased in the sense that its expectation is not zero under perfectly calibrated distributions. Moreover, it does not lead to an optimal test statistic. Informally, this happens because we are estimating both $E[Y | Z \in B_i]$ and $E[Z | Z \in B_i]$ with the same sample $\{(X_i; Y_i); i \in \{1; \dots; n\}\}$. We hence define the Debiased Plug-in Estimator (DPE):

$$T_{m;n}^d := \sum_{i=1}^{m^k-1} \frac{|I_{m;i}|}{n} \frac{1}{|I_{m;i}|} \sum_{j \in I_{m;i}} (Y_j - Z_j)^2 - \frac{1}{|I_{m;i}|^2} \sum_{j \in I_{m;i}} kY_j - Z_jk^2 \quad (6)$$

The debiasing term in (6) ensures that $T_{m;n}^d$ has mean zero under a distribution $P \in \mathcal{P}_0$ under which f is a calibrated probability predictor. Due to the discretization, the mean of $T_{m;n}^d$ is not exactly $\int_2\text{-ECE}(f)^2$ under $P \in \mathcal{P}_1(\cdot; p; s)$, but the debiasing makes it comparable to $\int_2\text{-ECE}(f)^2$. This will be a crucial step when proving the optimality of $T_{m;n}^d$.

Remark 2 (Connection to nonparametric functional estimation) The definition of $T_{m;n}^d$ is closely related to the U-statistic for estimating the quadratic integral functional of a probability density (Kerkycharian and Picard, 1996; Laurent, 1996). To see this, let $f_i(x) = P_Z(B_i) \mathbb{1}_{B_i}(x) : i = 1; \dots; m^k-1$ be the Haar scaling functions associated to the partition B_m . For each $1 \leq k \leq K$, the U-statistic

$$\frac{1}{n(n-1)} \sum_{1 \leq i < j \leq n} \sum_{1 \leq k < l \leq K} [Y_{j_1} - Z_{j_1}]_k [Y_{j_2} - Z_{j_2}]_l f_i(Z_{j_1}) f_l(Z_{j_2})$$

is an unbiased estimate of $\int_1 \sum_{i=1}^{m^k-1} P_Z(B_i) \int_1^R [\text{res}_f(z)]_k dP_Z(z)^2$. Summing over $1 \leq k \leq K$ and plugging in $P_Z(B_i) = \frac{|I_{m;i}|}{n}$, we recover (6) with the minor modification of changing $n! \rightarrow n-1$ in the scaling.

However, as noted in Remark 1, our problem differs from those studied in classical nonparametric statistic literature. Specifically, our definition of $T_{m;n}^d$ additionally requires an estimation of $P_Z(B_i)$ by $\frac{|I_{m;i}|}{n}$. Therefore, prior results on nonparametric functional

Algorithm 1 T-Cal: an optimal test for calibration (based on debiased plug-in estimation of the calibration error)

Input: Probability predictor $f : X \rightarrow \mathcal{K}$; i.i.d. sample $(X_i, Y_i)_{i=1}^n$; false detection rate $\beta \in (0, 1)$; true detection rate $\gamma \in (0, 1)$; Hölder smoothness s

Initialize: $m = \lfloor n^{2/(4s+K-1)} \rfloor$; $T_{m;n}^d = 0$; $Z_i = f(X_i)$ for $1 \leq i \leq n$; define B_1, \dots, B_m as in Appendix B.3

for $i = 1$ to m do

$l_{m;i} = \sum_{j \in B_i} Z_j \mathbb{1}_{\{j \in B_i\}}$

$T_{m;n}^d = T_{m;n}^d + \frac{\mathbb{1}_{\{l_{m;i} \neq 0\}}}{|l_{m;i}|} \sum_{j \in B_i} (Y_j - Z_j)^2 \frac{\mathbb{1}_{\{j \in B_i\}}}{|l_{m;i}|^2} \sum_{k \in \mathcal{K}} Y_j - Z_j k^2$

end for

$\alpha_{m;n} = \frac{\beta}{1 - \beta} \frac{2K}{m^{K/2} n^{1/2} \wedge m^{K/2}}$

Output: Reject H_0 if $\alpha_{m;n} = 1$

estimation (Bickel and Ritov, 1988; Donoho and Nussbaum, 1990; Birgé and Massart, 1995) cannot be directly applied to χ^2 -ECE(f)².

Wang et al. (2008); Shen et al. (2020) consider quadratic functional estimation for an unknown distribution of covariates and show that the minimax rate also depends on the Hölder smoothness of the covariate density function.

In the following theorem, we prove that $T_{m;n}^d$ leads to a minimax optimal test when the number of bins is chosen in a specific way, namely $m = \lfloor n^{2/(4s+K-1)} \rfloor$. Crucially, the number of bins required decreases with the smoothness parameter s . In this sense, our result parallels the well-known results on the optimal choice of the number of bins for testing probability distributions and densities (Mann and Wald, 1942; Ingster, 2012).

The guarantee on the power (or, Type II error control) requires the following mild condition, stated in Assumption 1. This ensures that the probability of each bin is proportional to the inverse of the number of bins up to some absolute constant. In particular, this holds if the density of the probabilities predicted is close to uniform. This assumption is necessary when extending the results of Arias-Castro et al. (2018); Kim et al. (2022) to a general base probability measure μ of the probability predictions over the probability simplex. See Appendix B.2 for more discussion.

Assumption 1 (Bounded marginal density) Let μ be the uniform probability measure on the probability simplex \mathcal{K} . There exist constants $\beta, \alpha > 0$ such that $\beta \leq \mu(Z \in B) \leq \alpha$ almost everywhere.

Theorem 3 (Calibration test via debiased plug-in estimation) Suppose $\beta \in (0, 1)$ and assume that the Hölder smoothness parameter s is known. For a binning scheme parameter $m \in \mathbb{N}_+$, let

$$\alpha_{m;n}(\beta) = \alpha_{m;n} := \frac{\beta}{1 - \beta} \frac{2K}{m^{K/2} n^{1/2} \wedge m^{K/2}};$$

Under Assumption 1 and for $m = bn^{2s+(K-1)c}$, we have

1. False detection rate control. For every P for which f is perfectly calibrated, i.e., for $P \in \mathcal{P}_0$, the probability of falsely claiming mis-calibration is at most ϵ , i.e., $P(T_{m;n}^d = 1) \leq \epsilon$.
2. True detection rate control. There exists $c > 0$ depending only on $(s; L; K; \epsilon; \epsilon_0; \epsilon_1)$ such that when

$$\epsilon \leq cn^{2s+(K-1)c},$$

then for every $P \in \mathcal{P}_1(\epsilon; p; s)$ i.e., when f is mis-calibrated with an ϵ -ECE of ϵ the power (true positive rate) is bounded as $P(T_{m;n}^d = 1) \geq 1 - \epsilon$.

The proof can be found in Appendix A.1. The proof follows the classical structure of upper bound arguments in nonparametric hypothesis testing, see e.g., Arias-Castro et al. (2018); Kim et al. (2022) for recent examples. We compute the mean and variance of $T_{m;n}^d$ under null distributions $P_0 \in \mathcal{P}_0$ and alternative distributions $P_1 \in \mathcal{P}_{s;L;K}$ with a large ECE. Using Lemma 13, we can find a lower bound on $E_{P_1}[T_{m;n}^d] - E_{P_0}[T_{m;n}^d]$. The variances $\text{Var}_{P_0}(T_{m;n}^d)$ and $\text{Var}_{P_1}(T_{m;n}^d)$ can be also upper bounded. We argue that the mean difference $E_{P_1}[T_{m;n}^d] - E_{P_0}[T_{m;n}^d]$ is significantly larger than the square root of the variances $\text{Var}_{P_0}(T_{m;n}^d)$ and $\text{Var}_{P_1}(T_{m;n}^d)$. The conclusion follows from Chebyshev's inequality.

Combined with our lower bound in Theorem 10, this result shows the desired property that our test is minimax optimal. This holds for all $p \geq 2$, so that the test is minimax optimal even when the mis-calibration is measured in the ϵ_p norm with $p < 2$. This is consistent with experimental findings such as those of Nixon et al. (2019), where the empirical ϵ_2 -ECE performs better than the empirical ϵ_1 -ECE as a measure of calibration error. Also see Section 4.1 for a comparison of the empirical ϵ_1 -ECE and ϵ_2 -ECE as a test statistic.

Although we present explicit critical values in Theorem 3, they can be conservative in practice, as in other works in nonparametric testing (Ingster, 1987; Arias-Castro et al., 2018; Kim et al., 2022). Therefore, we recommend choosing the critical values via a version of bootstrap: consistency resampling (Brockner and Smith, 2007; Vaicenavicius et al., 2019). See Appendix C.3 for further details on choosing critical values.

Remark 4 So far, we considered the null hypothesis of perfect calibration. However, since the predictor f is trained on a finite dataset, we cannot expect it to be perfectly calibrated. We can extend Theorem 3 to the null hypothesis of "small enough" mis-calibration, namely, for any given constant $c_0 > 0$, an ϵ -ECE of at most $c_0 n^{2s+(K-1)c}$. Then, the true and false positive rates of the test

$$p_{m;n}^{\text{comp}} := \mathbb{P}(T_{m;n}^d \leq \frac{K}{2} \frac{q}{m^K} \frac{1}{n^{2s+(K-1)c}} + 5c_0^2 n^{\frac{4s}{4s+K-1}} \frac{1}{n^{1+(K-1)c}})$$

can be controlled as in Theorem 3. See Appendix A.3 for the proof.

3.2 An Adaptive Test

The binning scheme used in our plug-in test requires knowing the smoothness parameter s to be minimax optimal. However, in practice, this parameter is usually unknown. Can we design an adaptive test that does not require knowing this parameter? Here we answer this question in the affirmative. As in prior works in nonparametric hypothesis testing, e.g., Ingster (2000); Arias-Castro et al. (2018); Kim et al. (2022), we propose an adaptive test that can adapt to an unknown Hölder smoothness parameters. The idea is to evaluate the plug-in test over a variety of partitions, and thus be able to detect mis-calibration at various different scales.

In more detail, we evaluate the test with a number of bins ranging over a dyadic grid $2; 2^2; \dots; 2^B$. In addition, to make sure that we control the false detection rate, we need to divide the level α by the number of tests performed. Thus, for a number $B = \lfloor \frac{2}{K-1} \log_2(n) \rfloor$ of tests performed, we let the adaptive test

$$T_n^{\text{ad}} := \max_{1 \leq b \leq B} T_{2^b; n} \quad (7)$$

detect mis-calibration if any of the debiased plug-in tests $T_{2^b; n}(\alpha/B)$, with the number of bins $2^b, b \in 1; \dots; B$, detects mis-calibration at level α/B . We summarize the procedure in Algorithm 2.

Theorem 5 (Adaptive plug-in test) Suppose $\alpha \geq 2$. Under Assumption 1, the adaptive test from (7) enjoys

1. **False detection rate control.** For every P for which f is perfectly calibrated, i.e., for $P \in \mathcal{P}_0$, the probability of falsely claiming mis-calibration is at most α , i.e., $\mathbb{P}_P(T_n^{\text{ad}} = 1) \leq \alpha$.
2. **True detection rate control.** There exists $c_{\text{ad}} > 0$ depending on $(s; L; K; \mu; \nu; \gamma)$ such that the power (true positive rate) is lower bounded $\mathbb{P}(T_n^{\text{ad}} = 1) \geq 1 - \alpha$ for every $P \in \mathcal{P}_1(\alpha; p; s)$ i.e., when f is mis-calibrated with an α -ECE of at least $\frac{c_{\text{ad}}}{2^{s+(4s+K-1) \log_2(n)}}$.

See Appendix A.4 for the proof. Compared to the non-adaptive test, this test requires a mild additional factor of $(\log n)^{s+(4s+K-1)}$ in the separation rate α to guarantee detection. It is well understood in the area of nonparametric hypothesis testing that some adaptation cost is unavoidable, see for instance Spokoiny (1996); Ingster (2000). For more discussion, see Remark 14.

Remark 6 We remark that the false detection rate control of Theorem 3 and 5 does not require a Hölder smoothness assumption.

3.3 Necessity of Debiasing

Recall from (5) that $T_{m; n}^b$ is the plug-in estimator of $\|f\|_2^2$ without the debiasing term in (6). We argue that this biased estimator is not an optimal test statistic, even for $m = m = \lfloor \frac{n}{2^{s+(4s+K-1) \log_2(n)}} \rfloor$ from Theorem 3 (which is optimal for the debiased test), by presenting a failure case in the following example.

Algorithm 2 Adaptive T-Cal: an adaptive test for calibration

Input: Probability predictor $f : X \rightarrow \mathcal{K} \setminus \{1\}$; i.i.d. sample $f(X_i; Y_i) \stackrel{P_0}{\sim} X \times Y \stackrel{P_0}{\sim} \{1, \dots, n\}$; false detection rate $\beta \in (0, 1)$; true detection rate $\gamma \in (0, 1)$
 Initialize: $B = \frac{2}{K-1} \log_2(n) \log_2(n) e$; $Z_i = f(X_i)$ for $1 \leq i \leq n$; $\frac{ad}{n} = 0$
 for $b = 1$ to B do
 Compute $T_{m;n}^{b; \gamma}(\bar{B})$ as in Algorithm 1
 if $T_{m;n}^{b; \gamma}(\bar{B}) = 1$ then
 $\frac{ad}{n} = 1$
 break
 end if
 end for
 Output: Reject H_0 if $\frac{ad}{n} = 1$

Example 1 (Failure of naive plug-in) Consider binary classification problem with $K = 2$, $m = bn^{2=(4s+1)c}$ (assumed to be divisible by four), and the partition

$$B_m = \{B_1, \dots, B_m\}, g = \left(0, \frac{1}{m}, \dots, \frac{m-1}{m}, 1\right)$$

Let P_0 be the distribution over $(Z; Y) \in [0, 1] \times \{0, 1\}$ given by $Z \stackrel{P_0}{\sim} \text{Unif}([0, 1])$ and $Y | Z = z \stackrel{P_0}{\sim} \text{Ber}(z)$ for all $z \in [0, 1]$. Under P_0 , the probability predictor f is perfectly calibrated, i.e., $P_0 \in \mathcal{P}_0$. Let $\phi : \mathcal{R} \rightarrow \mathcal{R}$ be the function defined by

$$\phi(x) := e^{-\frac{1}{x(1-x)}} \mathbf{1}_{(0,1)}(x) \tag{8}$$

Let $g : [0, 1] \rightarrow [0, 1]$ be the function (corresponding to the calibration curve of the probability predictor f)

$$g(z) := z \prod_{j=0}^{m-1} \frac{z - \frac{j}{m}}{z - \frac{j-1}{m}} + \prod_{j=\frac{m}{4}}^{m-1} \frac{z - \frac{j}{m}}{z - \frac{j-1}{m}} \tag{9}$$

for $s \in (\frac{1}{4}, \frac{1}{2})$, $\gamma > 0$, and ϕ defined in (8). Define the distribution P_1 over $(Z; Y)$ by

$$Z \stackrel{P_1}{\sim} \text{Unif}([0, 1]) \text{ and } Y | Z = z \stackrel{P_1}{\sim} \text{Ber}(g(z))$$

for all $z \in [0, 1]$. As we will show in the proof of Theorem 10, (1) the mis-calibration curve $g(z) - z = \text{res}_f(z)$ is s -Holder and (2) $\int_0^1 g(z) dz - \int_0^1 z dz = \text{ECE}_{P_1}(f) = \frac{1}{n^{2s=(4s+1)c}}$.

As can be seen in Figure 2a, the probability predictor f under P_1 is an example of a mis-calibrated predictor, as $f(X)$ is smaller than $E[Y | f(X)]$ when $f(X)$ is above 0.5; and vice versa. However, the mean of $T_{m;n}^b$ under the mis-calibrated distribution P_1 is surprisingly smaller than the mean under the calibrated distribution P_0 when n is large enough (Proposition 7).

That is, the statistic $T_{m;n}^b$ does not capture the amount of mis-calibration, and therefore the calibration test based on it will not perform well. Figure 2b confirms this finding, and Figure 2c displays that this effect can be removed by using the debiased statistic $T_{m;n}^d$.

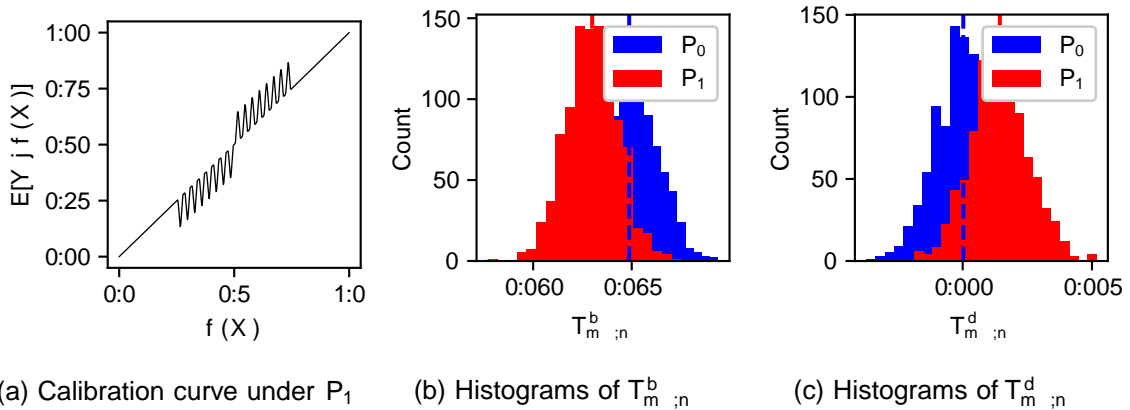


Figure 2: (a) A graph of the calibration curve $z \mapsto g(z) = E_{P_1}[Y | f(X) = z]$ defined in (9). When the true label probability is above/below 0.5, the model outputs a smaller/larger score. Hence f is a mis-calibrated probability predictor under P_1 . (b) Histograms of $T_{m;n}^b$ and $T_{m;n}^d$ under P_0 and P_1 are obtained from 1,000 independent observations. We use the parameters $m = 10,000$, $s = 0.3$, and $n = 100$. The dashed line indicates the empirical mean of each distribution. Note that the biased estimator $T_{m;n}^b$ has a smaller mean under P_1 , which aligns with Proposition 7. (c) We see this effect disappears after debiasing and that the mean of $T_{m;n}^d$ becomes zero.

Proposition 7 (Failure of naive plug-in test) Let P_0 and P_1 be the distributions defined in Example 1, and $m = bn^{2=(4s+1)}$. Then $E_{P_0}[T_{m;n}^b] < E_{P_1}[T_{m;n}^b]$ for all large enough $n \in \mathbb{N}_+$.

See Appendix A.5 for the proof. We remark that it is possible to avoid the phenomenon in Proposition 7, by choosing a different m . Proposition 7 aims only to highlight that the effect of the bias in $T_{m;n}^b$ can be extreme in certain cases, and we do not claim that $m = m$ is also the optimal choice for the biased statistic.

We finally comment on the related results of Brocker (2012); Ferro and Fricker (2012); Kumar et al. (2019). In Brocker (2012); Ferro and Fricker (2012), the plug-in estimator of the squared χ^2 -ECE is decomposed into terms related to reliability and resolution. Based on this observation, Kumar et al. (2019) propose a debiased estimator for the squared χ^2 -ECE and show an improved sample complexity for estimation. However, their analysis is restricted to the binary classification case and probability predictors with only finitely many output values. It is not clear how to adapt their method to predictors with continuous outputs, because this would require discretizing the outputs. Our debiased plug-in estimator $T_{m;n}^d$ is more general, as it can be used for multi-class problems and continuous probability predictors f . Also, our reason to introduce $T_{m;n}^d$ (testing) differs from theirs (estimation).

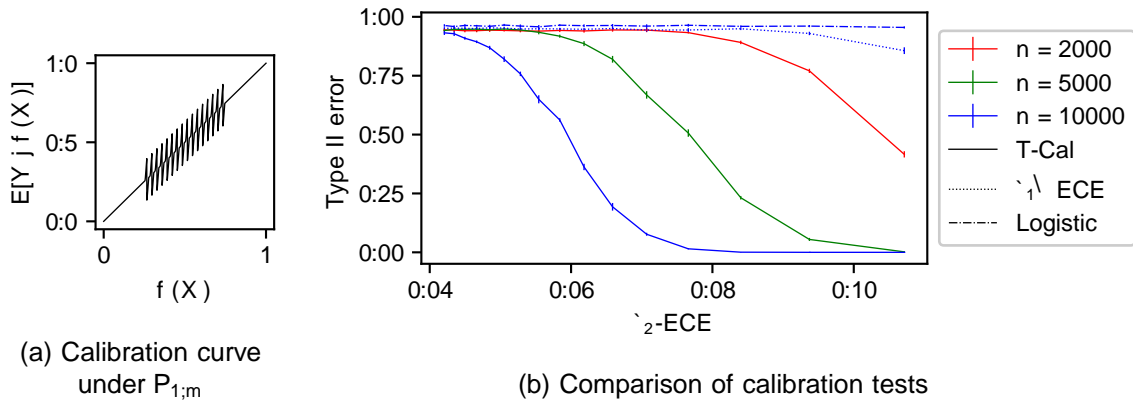


Figure 3: (a) A graph of the calibration curve $z \mapsto g_m(z) = E_{P_{1;m}}[Y | f(X) = z]$ defined in (10). The mis-calibration curve alternates between negative and positive values, making detection challenging. (b) We compare our test $\hat{\rho}_{2;m;n}$ with other commonly used calibration tests. Since our test optimally adjusts the number of bins $m = \lfloor n^{2/(4s+K-1)} \rfloor$ according to the sample size n , it can detect mis-calibration over smaller and smaller intervals as n grows. On the other hand, the plug-in test $\hat{\rho}_{1;m;n}$, with a fixed n binning scheme parameter m , fails to detect mis-calibration over intervals smaller than the bin width. This issue remains when the sample size increases. The test based on the calibration slope and intercept also suffers from the same issue. Standard error bars are plotted over 10 repetitions.

4 Experiments

We perform experiments on both synthetic and empirical datasets to support our theoretical results. These experiments suggest that T-Cal is in general superior to state-of-the-art methods.

4.1 Synthetic Data: Power Analysis

Let $P_0 \neq P_0$ be the distribution defined in Example 1|a distribution under which f is perfectly calibrated. For $m \in \mathbb{N}_+$, $s > 0$, $\alpha > 0$, and $\beta : \mathbb{R} \rightarrow \mathbb{R}$ from (8), define $g_m : [0; 1] \rightarrow [0; 1]$ by

$$g_m(z) := z + \frac{\alpha}{m^s} \sum_{j=0}^{\lfloor m^s z \rfloor} (-1)^j \frac{2mz}{m} - j \quad (10)$$

This function oscillates strongly, as shown in Figure 3a. Let $P_{1;m}$ be the distribution over $(Z; Y) \in [0; 1] \times [0; 1]$ given by $Z \sim P_{1;m} \text{Unif}([0; 1])$ and $Y | Z = z \sim \text{Ber}(g_m(z))$ for all $z \in [0; 1]$. Under $P_{1;m}$, the probability predictor f is mis-calibrated with an $\hat{\rho}_p$ -ECE of at least $\alpha = \frac{\alpha}{K_L p} m^{-s}$. However, since the mis-calibration curve g_m oscillates strongly, mis-calibration can be challenging to detect.

We study the type II error of tests against the alternative where the mis-calibration is specified as $H_1 : (Z; Y) \sim P_{1;m}$. This gives a lower bound on the worst-case type II error

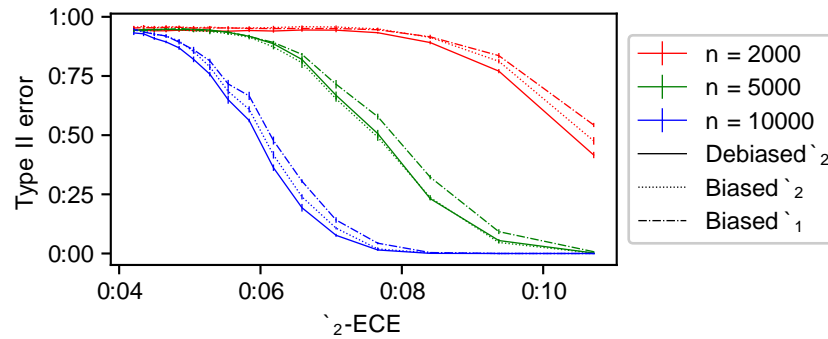


Figure 4: Type II error comparison for $T_{m;n}^d$ (T-Cal), $T_{m;n}^b$, and $T_{m;n}^1$. Using $\hat{\epsilon}_2$ is better than $\hat{\epsilon}_1$, and debiased $\hat{\epsilon}_2$ (T-Cal) is better than biased $\hat{\epsilon}_2$. Standard error bars are plotted over 10 repetitions.

over the alternative hypothesis $P_1(\cdot; p; s)$. We repeat the experiment for different values of m to obtain a plot of $\hat{\epsilon}_2$ -ECE versus type II error.

Comparison of Tests. We compare the test $T_{m;n}$ with classical calibration tests dating back to Cox (1958), and discussed in Harrell (2015); Vaicenavicius et al. (2019). Harrell (2015) fits a logistic model

$$P(Y = 1 | Z) = \frac{1}{1 + \exp(\beta_0 + \beta_1 \log \frac{Z}{1-Z})}$$

on the sample $(Z_i; Y_i) : i = 1, \dots, n$ and tests the null hypothesis of $\beta_0 = 0$ and $\beta_1 = 1$. Specifically, we perform the score test (Rao, 1948; Silvey, 1959), with the test statistic derived from the gradient of log-likelihood with respect to the tested parameters. There are several approaches to set the critical values, including by using the asymptotic distribution theory of sampling statistics under the null hypothesis, or by data reuse methods such as the bootstrap. We estimate the critical values via 1000 Monte Carlo simulations.

Vaicenavicius et al. (2019) use $\hat{\epsilon}_1$ -ECE, the plug-in estimator for ϵ_1 -ECE, as their test statistic. They approximate the distribution of $\hat{\epsilon}_1$ -ECE by a bootstrapping procedure called consistency resampling (in which both the probability predictions and the labels are resampled) and compute a p-value based on this approximation. This test also uses a plug-in estimator as the test statistic but differs from T-Cal as it is neither debiased nor adaptive. Since the data-generating distribution is known in this synthetic experiment, we set the critical value via 1000 Monte Carlo simulations.

We control the false detection rate at a level $\alpha = 0.05$ and run experiments for $n = 2,000; 5,000$, and $10,000$. We find that our proposed test achieves the lowest type II error, see Figure 3b. We also find that other tests do not leverage the growing sample size. For this reason, we only display $n = 10,000$ for the other two tests. As can be seen in Figure 3b, T-Cal outperforms other testing methods in true detection rate by a large margin.

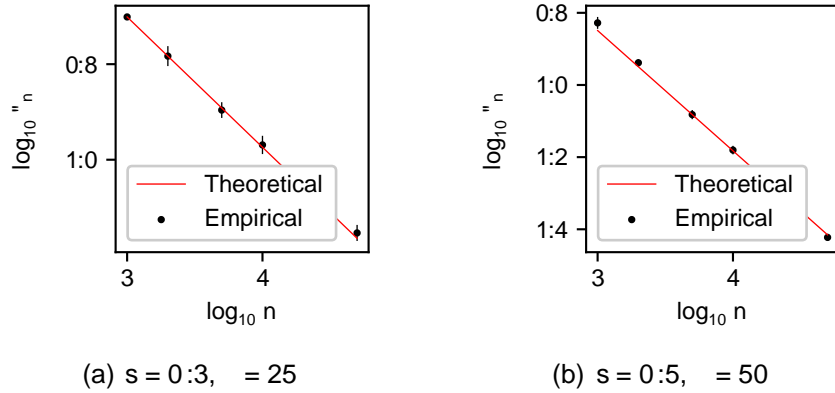


Figure 5: The dots are \log^n computed for different sample sizes n . The red line has a slope $\frac{2s}{4s+1}$. Standard error bars are plotted over 10 repetitions. See the text for more details.

ℓ_1 -ECE versus ℓ_2 -ECE. To confirm the effectiveness of using the ℓ_2 -ECE estimator $T_{m,n}^d$, we compare it with a plug-in ℓ_1 -ECE estimator defined as

$$T_{m,n}^{\ell_1} := \sum_{j=1}^K \frac{|I_{m,j}|}{n} \frac{1}{|I_{m,j}|} \sum_{i \in I_{m,j}} (Y_i - Z_i) :$$

At the moment, it is unknown how to debias this estimator. We use the optimal binning parameter $m = m_n = \lfloor n^{2/(4s+K-1)} \rfloor$ for both ℓ_1 and ℓ_2 estimators; because it is unknown what the ℓ_1 -optimal binning scheme is. Also, to isolate the effect of debiasing, we compare the biased ℓ_2 estimator $T_{m,n}^b$ as well, with the same number of bins. In Figure 4, we see the ℓ_2 estimators consistently outperform the ℓ_1 estimator, regardless of debiasing. While it is a common practice to use a plug-in estimator of ℓ_1 -ECE, our result suggests T-Cal compares favorably to it.

Minimum Detection Rate. We perform an experiment to support the result on the minimum detection rate of T-Cal, presented in Theorem 3. For each n , we find the largest integer, denoted $m(n)$, such that the type II error against $H_1 : (Z; Y) \sim P_{1;m(n)}$ is less than 0.05. We compute $\beta_n := \ell_2\text{-ECE}_{P_{1;m(n)}}(f) = k_{L^2} m(n)^s$ (a lower bound on the minimum detection rate) and plot $\log \beta_n$ versus $\log n$ in Figure 5. We see that the logarithm decreases as n grows, with the slope $\frac{2s}{4s+1}$ predicted by Theorem 3.

4.2 Results on Empirical Datasets

To verify the performance of adaptive T-Cal empirically, we apply it to the probability predictions output by deep neural networks trained on several datasets. Since our goal is to test calibration, we calculate the probabilities predicted by pre-trained models on the test sets. As in (Guo et al., 2017; Kumar et al., 2019; Nixon et al., 2019, etc), we binarize the test labels by taking the top-1 confidence as the new probability prediction,

	DenseNet 121		ResNet 50		VGG-19	
	$\hat{\epsilon}_1$ -ECE	Calibrated?	$\hat{\epsilon}_1$ -ECE	Calibrated?	$\hat{\epsilon}_1$ -ECE	Calibrated?
No Calibration	2.02%	reject	2.23%	reject	2.13%	reject
Platt Scaling	2.32%	reject	1.78%	reject	1.71%	reject
Poly. Scaling	1.71%	reject	1.29%	reject	0.90%	accept
Isot. Regression	1.16%	reject	0.62%	reject	1.13%	accept
Hist. Binning	0.97%	reject	1.12%	reject	1.28%	reject
Scal. Binning	1.94%	reject	1.21%	reject	1.67%	reject

Table 1: The values of the empirical $\hat{\epsilon}_1$ -ECE (Guo et al., 2017) and the testing results, via adaptive T-Cal and multiple binomial testing, of models trained on CIFAR-10.

and the labels as the results of the top-1 classification, i.e. $\hat{Z} = \max_{1 \leq k \leq K} [Z]_k$ and $\hat{Y} = l$ (correctly classified by the top-1 prediction). This changes the problem of detecting the full-class mis-calibration to testing the mis-calibration of a binary classifier. Hence, we choose $K = 2$ for adaptive T-Cal in the experiments below. We refer readers to (Gupta and Ramdas, 2022) for more details about binarization via the top-1 prediction.

CIFAR-10. For the CIFAR-10 dataset, the models are DenseNet 121, ResNet 50, and VGG-19. We first apply the adaptive test directly to the 10,000 uncalibrated probability predictions output by each model, with the false detection rate controlled at the level $\alpha = 0.05$.

For every choice of the number of bins m , we estimate the critical value by taking the upper 5% quantile of the values of the test statistic over 3000 bootstrap re-samples of the probability predictions. The labels are also chosen randomly, following Bernoulli distributions with the probability prediction as the success probability. We also provide the values of the standard empirical $\hat{\epsilon}_1$ -ECE calculated with Guo et al. (2017)'s approach for the reader's reference, and with 15 equal-width bins.

We then test the probability predictions of these three models calibrated by several post-calibration methods: Platt scaling (Platt, 1999), polynomial scaling, isotonic regression (Zadrozny and Elkan, 2002), histogram binning (Zadrozny and Elkan, 2001), and scaling-binning (Kumar et al., 2019). To this end, we split the original dataset of 10,000 images into 2 sets of sizes 2,000 and 8,000. The first set is used to calibrate the model, and the second is used to perform adaptive T-Cal and calculate the empirical $\hat{\epsilon}_1$ -ECE. In polynomial scaling, we use polynomials of order 3 to do regression on all the prediction-label pairs $(Z_i; Y_i)$, and truncate the calibrated prediction values into the interval $[0; 1]$. We set the binning scheme in both histogram binning and scaling binning as 15 equal-mass bins. Our implementation is adapted from Kumar et al. (2019).

Since the recalibrated probability predictions output by the latter two methods belong to a finite set, we use a test based on the binomial distribution. See Appendix B.4 for details. For completeness, we also provide the debiased empirical $\hat{\epsilon}_1$ -ECE values (Kumar et al., 2019) for models calibrated by the two discrete methods, see the details in Table 4, Appendix C.1.

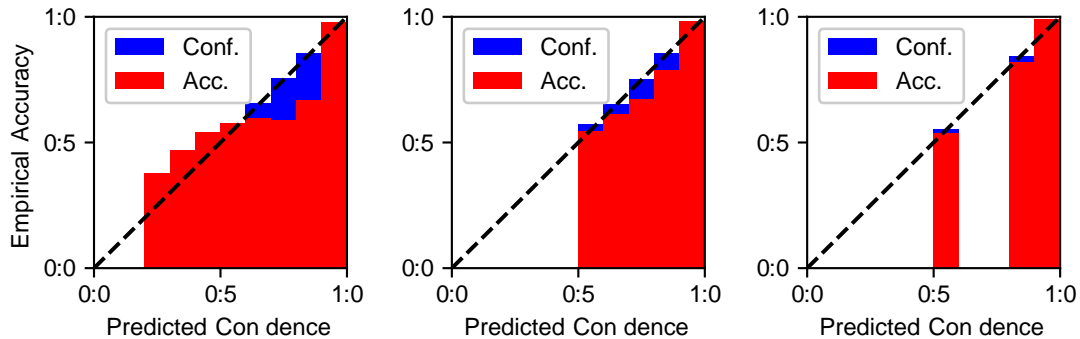


Figure 6: The reliability diagrams for VGG-19, trained on CIFAR-10, calibrated by Platt scaling (left), polynomial scaling (middle), and histogram binning (right). The bins (bars) containing less than 10 data points, where the sample noise dominates, are omitted for clarity. The dashed lines correspond to perfect calibration.

The results are listed in Table 1, where we use “accept” to denote that the test does not reject. The models with smaller empirical χ^2 -ECE are more likely to be accepted, by adaptive T-Cal and by multiple binomial testing, as perfectly calibrated. This can be further illustrated by the three empirical reliability diagrams given in Figure 6, where the model’s predictions calibrated by Platt scaling (left) are visually more “mis-calibrated” than those calibrated by polynomial scaling (middle) and histogram binning (right).

CIFAR-100. We perform the same experimental procedure for three models pre-trained on the CIFAR-100 dataset: MobileNet-v2, ResNet 56, and ShuffleNet-v2 (Chen, 2021). The test set provided by CIFAR-100 is split into two parts, containing 2,000 and 8,000 images, respectively. Since the regression functions of models trained on the larger CIFAR-100 dataset can be more complicated than those of models trained on CIFAR-10, we set the polynomial degree as ν in polynomial scaling.

The results are listed in Table 2. The values of the debiased empirical χ^2 -ECE (Kumar et al., 2019) for the two discrete calibration methods are provided in Table 5, Appendix C.1. The results roughly align with the magnitude of the empirical ECE value.

However, as can be observed in the column corresponding to ResNet 56, this trend is certainly not monotone. The calibrated ResNet 56 with the empirical ECE 18% is accepted while the calibrated ResNet 56 with a smaller value: 5% is rejected. Furthermore, the test results reveal that models with relatively large (or small) empirical χ^2 -ECE values may not necessarily be poorly (or well) calibrated since the χ^2 -ECE values measured can be highly dominated by the sample noise.

ImageNet. We repeat the above experiments on models pre-trained on the ImageNet dataset. We examine three pre-trained models provided in the torchvision package in PyTorch: DenseNet 161, ResNet 152, and EfficientNet-b7. We split the validation set of 50,000

	MobileNet-v2		ResNet 56		Shu eNet-v2	
	$\hat{\epsilon}_1$ -ECE	Calibrated?	$\hat{\epsilon}_1$ -ECE	Calibrated?	$\hat{\epsilon}_1$ -ECE	Calibrated?
No Calibration	11.87%	reject	15.2%	reject	9.08%	reject
Platt Scaling	1.40%	accept	1.84%	accept	1.34%	accept
Poly. Scaling	1.69%	reject	1.91%	reject	1.81%	accept
Isot. Regression	1.76%	accept	2.33%	reject	1.38%	accept
Hist. Binning	1.66%	reject	2.44%	reject	2.77%	reject
Scal. Binning	1.85%	reject	1.57%	reject	1.65%	accept

Table 2: The values of the empirical $\hat{\epsilon}_1$ -ECE (Guo et al., 2017) and the testing results, via adaptive T-Cal and multiple binomial testing, of models trained on CIFAR-100.

images into a calibration set and a test set of sizes 10,000 and 40,000, respectively. We use polynomials of degree 5 in polynomial scaling.

The results are listed in Table 3. The values of the debiased empirical $\hat{\epsilon}_2$ -ECE (Kumar et al., 2019) are provided in Table 6, in Appendix C.1. As can be seen, the test results here generally align with the empirical ECE values.

Remark 8 While it is hard to verify the Hölder smoothness of reg_k in these empirical datasets, we believe that some level of justification would be possible assuming (1) enough differentiability on f (which we think to be true for common neural net architectures and activation functions), (2) Y being deterministic given X (which is reasonable for low-noise datasets such as ImageNet), and (3) the set of inputs $C_k = \{x \in X : (Y|X = x) = e_k\}$ corresponding to each class being a Lipschitz domain, locally the graph of a Lipschitz continuous function. Given these assumptions, each coordinate of the regression function $\text{reg}_k(z)$ can be written as $[\text{reg}_k(z)]_k = E[Y_k | f(X) = z] = P_Z(C_k \cap \{f(X) = z\}) / P_Z(\{f(X) = z\})$; assuming the denominator is strictly positive. Then, Hölder continuity may follow from the inverse function theorem and the coarea formula (Federer, 2014, Theorem 3.2.3), which expresses the measure of a level set as an integral of the Jacobian.

However, this argument still requires making the essentially unverifiable Assumption (3) from the above paragraph. In some cases, this assumption can be viewed as reasonable: for instance, one may reasonably think that image manifolds for classes in ImageNet are locally Lipschitz; and thus Assumption (3) may hold. Under such conditions, this rough argument may provide an idea of why Hölder smoothness could be reasonable for certain predictive models such as neural net architectures.

When the Hölder condition does not hold, we still have the type I error guarantee, but may not have type II error control. We expect that the Hölder condition might be drastically violated when there is obvious discontinuity of the predictors/regression functions (e.g., a decision tree/random forest trained on data having discrete features).

	DenseNet 161		ResNet 152		EfficientNet-b7	
	$\hat{\gamma}$ -ECE	Calibrated?	$\hat{\gamma}$ -ECE	Calibrated?	$\hat{\gamma}$ -ECE	Calibrated?
No Calibration	5.67%	reject	4.99%	reject	2.82%	reject
Platt Scaling	1.58%	reject	1.41%	reject	1.90%	reject
Poly. Scaling	0.62%	accept	0.64%	accept	0.71%	accept
Isot. Regression	0.63%	reject	0.80%	reject	1.06%	reject
Hist. Binning	0.46%	reject	1.26%	reject	0.88%	reject
Scal. Binning	1.55%	reject	1.40%	reject	1.97%	reject

Table 3: The values of the empirical $\hat{\gamma}$ -ECE (Guo et al., 2017) and the testing results, via adaptive T-Cal and multiple binomial testing, of models trained on ImageNet.

5 Lower Bounds for Detecting Mis-calibration

To complement our results on the performance of the plug-in tests proposed earlier, we now show some fundamental lower bounds for detecting mis-calibration. We also provide a reduction that allows us to test calibration via two-sample tests and sample splitting. We show that this has a minimax optimal performance, but empirically does not perform as well as our previous test.

5.1 Impossibility for General Continuous Mis-calibration Curves

In Proposition 9, we show that detecting mis-calibration is impossible, even when the sample size is arbitrarily large unless the mis-calibration curve res_β has some level of smoothness. Intuitively, if the mis-calibration curve can be arbitrarily non-smooth, then it can oscillate between positive and negative values with arbitrarily high frequency, and these oscillations cannot be detected from a finite sample.

In this regard, one needs to be careful when concluding the quality of calibration from a finite sample. If we only assume that the mis-calibration curve is a continuous function of the probability predictions, then it is impossible to tell apart calibrated and mis-calibrated models. Further, for more complex models such as deep neural networks, one expects the predicted probabilities to be able to capture larger and larger classes of functions; thus this result is even more relevant for modern large-scale machine learning.

Let $\mathcal{P}_1^{\text{cont}}(\epsilon; p)$ be the family of probability distributions P over $(Z; Y)$ such that $\hat{\gamma}_p\text{-ECE}_P(f) \leq \epsilon$ and every entry of the mis-calibration curve $\text{res}_{\beta, P}$ is continuous. This is a larger set of distributions than $\mathcal{P}_1(\epsilon; p; s)$ in (3), because we only assume continuity, not Hölder smoothness. Denote the corresponding minimax type II error by $R_n^{\text{cont}}(\epsilon; p)$, namely

$$R_n^{\text{cont}}(\epsilon; p) := \inf_{\mathcal{A}_n(\epsilon)} \sup_{P \in \mathcal{P}_1^{\text{cont}}(\epsilon; p)} P(\text{reject} = 0):$$

This has the same interpretation as before, namely, it is the best possible false negative rate for detecting mis-calibration for data distributions belonging to $\mathcal{P}_1^{\text{cont}}(\epsilon; p)$, in a worst-case sense.

Proposition 9 (Impossibility of detecting mis-calibration) Let $\epsilon_0 = 0:1$. For any level $\alpha \in (0; 1)$, the minimax type II error $R_n^{\text{cont}}(\epsilon_0; p)$ for testing the null hypothesis of calibration at level α against the hypothesis $P \in \mathcal{P}_1^{\text{cont}}(\epsilon_0; p)$ of general continuous mis-calibration curves satisfies $R_n^{\text{cont}}(\epsilon_0; p) \geq 1 - \alpha$ for all n .

In words, this result shows that for a certain fixed calibration error ϵ_0 , and for a fixed false positive rate $\alpha > 0$, the false negative rate is at least $1 - \alpha$. Thus, it is not possible to detect mis-calibration in this setting. The choice of the constant $\epsilon_0 = 0:1$ is arbitrary and can be replaced by any other constant; the result holds with minor modifications to the proof.

The proof can be found in Appendix A.7. We make a few remarks on related results. While Example 3.2 of Kumar et al. (2019) demonstrates that related to earlier results on probability distribution and density estimation (Mann and Wald, 1942) using a binned estimator of ECE can arbitrarily underestimate the calibration error, we show a fundamental failure not due to binning, but instead due to the finite sample size. Also, our result echoes Theorem 3 of Gupta et al. (2020) which states that asymptotically perfect calibration is only possible for probability predictors with a countable support; but this does not overlap with Proposition 9 as our conclusion is about the impossibility of detecting mis-calibration.

5.2 Hölder Alternatives

As is customary in nonparametric statistics (Ingster, 1987; Low, 1997; Györfi et al., 2002; Ingster, 2012), we consider testing against Hölder continuous alternatives; or, differently put, detecting mis-calibration when the mis-calibration curves are Hölder continuous. This excludes the pathological examples where the mis-calibration curves oscillate widely that were discussed in Section 5.1; but still allow a very rich class of possible mis-calibration curves, including non-smooth ones.

Theorem 10 states that, for a K -class classification problem and for alternatives with a Hölder smoothness parameter s , the mis-calibration of a model can be detected only when the calibration error is of order $(n^{-2s/(4s+K-1)})$. In other words, the smallest possible calibration error that can be detected using a sample of size n is of order $n^{-2s/(4s+K-1)}$.

Testing calibration of a probability predictor in our nonparametric model leads to rates that are slower than the parametric case $n^{-1/2}$. This is because $2s/(4s+K-1) < 1/2$ for $s > 0$ and $K \geq 2$. The rate becomes even slower as the number of classes K grows. This indicates that evaluating model calibration on a small-sized dataset can be problematic. Further, it suggests that multi-class calibration may be even harder to achieve.

This rate is what one may expect based on results for similar problems in nonparametric hypothesis testing (Ingster, 2012), with $K - 1$ interpreted as the dimension. Specifically, the rate is equal to the minimum separation rate in two-sample goodness-of-fit testing for densities on $\mathcal{K} - 1$. This connection to two-sample testing will be made clear in Section 6.

Theorem 10 (Lower bound for detecting mis-calibration) Given a level $\alpha \in (0; 1)$ and $\beta \in (0; 1)$, consider the hypothesis testing problem (4), in which we test the calibration of the K -class probability predictor f assuming $(s; L)$ -Hölder continuity of mis-calibration curves as defined in (62). There exists $c_{\text{lower}} > 0$ depending only on

$(p; s; L; K; \epsilon)$ such that, for any $p > 0$, the minimum ϵ -ECE of f , i.e. $\epsilon_n(p; s)$, required to have a test with a false positive rate (type I error) at most ϵ and with a true positive rate (power) at least $1 - \epsilon$ satisfies $\epsilon_n(p; s) \geq c_{\text{lower}} n^{-2s/(4s+K-1)}$ for all n .

See Appendix A.8 for the proof. The proofs of both Proposition 9 and Theorem 10 are based on Ingster's method, also known as the chi-squared or Ingster-Suslina method (Ingster, 1987, 2012). Informally, Ingster's method states that if we can select alternative distributions with an average likelihood ratio to a null distribution close to unity, then no test with a fixed level can control the minimax type II error below a certain threshold.

Remark 11 (Hölder smoothness assumption) Since the residual function $\text{res}_{\epsilon, p}$ depends on the unknown joint distribution P of $(Z; Y)$, the Hölder continuity of the map $z \mapsto [\text{res}_{\epsilon}(z)]_k$ for each $k \in \{1, \dots, K\}$ is in general an assumption that we need to make. When the residual map res_{ϵ} does not satisfy the Hölder assumption, we still have the false detection rate control in Theorem 3 and 5, but we cannot guarantee the true detection rate control and the lower bound in Theorem 10. Extending our approach beyond the Hölder assumption may be possible in future work, inspired by works in nonparametric hypothesis testing that study Besov spaces of functions (Ingster, 2012).

6 Reduction to Two-sample Goodness-of-fit Testing

To further put our work in context in the literature on nonparametric hypothesis testing, in this section we carefully examine the connections between the problem of testing calibration, and a well-known problem in that area. Specifically, we describe a novel randomization scheme that allows us to reduce the null hypothesis of perfect calibration to a hypothesis of equality of two distributions—making a strong connection to the problem of two-sample goodness of fit testing. In other words, we can use a calibrated probabilistic classifier and randomization to generate two samples from an identical distribution. For a mis-calibrated classifier the same scheme will generally result in two samples from two different distributions.

If a classifier is perfectly calibrated, then its class probability predictions will match the true prediction-conditional class probabilities. Therefore, randomly sampling labels according to the classifier's probability predictions will yield a sample from the empirical distribution. We rely on sample splitting to obtain two samples: the empirical and the generated one. Then we can use any classical test to check if the two samples are generated from the same distribution. As we will show, the resulting test has a theoretically optimal detection rate. However, due to the sample splitting step, its empirical performance is inferior to the test based on the debiased plug-in estimator from Section 3.

We split our sample into two parts. For $i \in \{1, \dots, n\}$, we generate random variables Y_i following the categorical distribution $\text{Cat}(Z_i)$ over classes $Y = \{1, \dots, K\}$, with a K -class probability distribution $Z_i = f(X_i)$ predicted by the classifier f . These Y_i are independent of each other and of Y_1, \dots, Y_{n-2c} , due to the sample splitting step. For each $k \in \{1, \dots, K\}$, define

$$V_k := \frac{1}{n} \sum_{i=1}^n Z_i : [Y_i]_k = 1; \quad \bar{V}_k := \frac{1}{n} \sum_{i=1}^n Z_i : [Y_i]_k \neq 1$$

and

$$W_k := \frac{1}{n} \sum_{i=1}^n Z_i : [Y_i]_k = 1; \quad \frac{1}{2} + \frac{1}{n} \quad \text{for } k = 1, \dots, K$$

By construction, V_k is an i.i.d. sample from the distribution on the probability simplex Δ_{K-1} with a density³

$$V_k(z) := \frac{\mathbb{R}_{\Delta_{K-1}} [\text{reg}_f(z)]_k}{\int_{\Delta_{K-1}} [\text{reg}_f(z)]_k dP_Z(z)} = \frac{[\text{reg}_f(z)]_k}{E[Y]_k}$$

with respect to P_Z . Similarly, W_k is an i.i.d. sample from the distribution on Δ_{K-1} with a density

$$W_k(z) := \frac{\mathbb{R}_{\Delta_{K-1}} [z]_k}{\int_{\Delta_{K-1}} [z]_k dP_Z(z)} = \frac{[z]_k}{E[Z]_k}$$

Now we consider testing the null hypothesis

$$H_0 : V_k = W_k \text{ for all } k \in \{1, \dots, K\}$$

against the complement of H_0 . We claim that if we use an appropriate test for this null hypothesis, with an additional procedure to rule out "easily detectable" alternatives, then we can obtain a test that attains the optimal rate specified in Theorem 10.

We describe the main idea of this reduction. A formal result can be found in Theorem 12. Let $P \in \mathcal{P}_1(\cdot; p; s)$ and assume $\epsilon = \frac{1}{n} \epsilon^{2s/(4s+K-1)}$.⁴ The squared distance between V_k and W_k in $L^2(P_Z)$ is

$$\int_{\Delta_{K-1}} \left(\frac{[\text{reg}_f(z)]_k}{E[Y]_k} - \frac{[z]_k}{E[Z]_k} \right)^2 dP_Z(z) = \int_{\Delta_{K-1}} \frac{[\text{reg}_f(z)]_k^2}{E[Y]_k^2} + \frac{[z]_k^2}{E[Z]_k^2} - \frac{2[\text{reg}_f(z)]_k [z]_k}{E[Y]_k E[Z]_k} dP_Z(z) \quad (11)$$

Further,

$$\int_{\Delta_{K-1}} [z]_k [\text{reg}_f(z)]_k dP_Z(z) = E_P[[Z]_k E[Y | Z]_k] = E_P[[Z]_k [Y | Z]_k] \quad (12)$$

Since $E[Y | Z]_k = E[[Z]_k [Y | Z]_k] = 0$ and $\text{Var}([Y | Z]_k) = \text{Var}([Z]_k [Y | Z]_k) = 1$ under H_0 , we can detect mis-calibration for the alternatives $P \in \mathcal{P}_1(\cdot; p; s)$ such that $E_P[[Y | Z]_k] = \frac{1}{n} \epsilon^{1/2}$ or $E_P[[Z]_k [Y | Z]_k] = \frac{1}{n} \epsilon^{1/2}$ by rejecting the null hypothesis H_0 of calibration if $\frac{1}{n} \sum_{i=1}^n [Y_i | Z_i]_k \geq c n^{1/2}$ or $\frac{1}{n} \sum_{i=1}^n [Z_i]_k [Y_i | Z_i]_k \geq c n^{1/2}$ for some $c > 0$. For the remaining alternatives, choose $k_0 \in \{1, \dots, K\}$ such that

$$\frac{1}{(E[Y]_{k_0})^2} \int_{\Delta_{K-1}} [\text{reg}_f(z)]_{k_0}^2 dP_Z(z) \geq \frac{\epsilon^2}{K (E[Y]_{k_0})^2} = \frac{1}{n} \epsilon^{\frac{4s}{4s+K-1}}:$$

3. We assume that the densities V_k and W_k are well defined, and in particular that $E[Y]_k > 0$ and $E[Z]_k > 0$ for every $k \in \{1, \dots, K\}$. This follows from Assumption 2, which will be introduced later in the section.

4. Here we use the notation $\epsilon(\cdot)$ to include the adaptive case. See Corollary 13 and Remark 14 for further details.

Then, $k \frac{V}{k} \frac{W}{k} k_{L^2(P_Z)}^2$ is at least

$$\tau \left(n^{\frac{4s}{4s+K-1}} \right) + \frac{2E[Z - Y]_{k_0} E[[Z]_{k_0} [Y - Z]_{k_0}]}{(E[Y]_{k_0})^2 E[Z]_{k_0}} = \tau \left(n^{\frac{4s}{4s+K-1}} \right):$$

Since $|V_{k_0}|, |W_{k_0}| = \binom{n}{k_0}$ with high probability, the power of the test can be controlled using standard results on two sample testing. The full procedure is described in Algorithm 3.

In general, for a positive integer $d > 0$, we allow using an arbitrary deterministic two-sample testing procedure $TS_d : ([0; 1]^d)^{n_1} \times ([0; 1]^d)^{n_2} \rightarrow \{0, 1\}$, which takes in two d -dimensional samples V_1, \dots, V_{n_1}, g and W_1, \dots, W_{n_2}, g and outputs "1" if and only if the null hypothesis is rejected. The two samples V_1, \dots, V_{n_1}, g and W_1, \dots, W_{n_2}, g are sampled i.i.d. from distributions with densities f_1 and f_2 , respectively, with respect to an appropriate probability measure on $[0; 1]^d$. Further, it is assumed that $f_1 - f_2$ is $(s; L)$ -Holder continuous for a Holder smoothness parameters $s > 0$ and a Holder constant $L > 0$. Given $\epsilon \in (0; 1)$ and $\delta \in (0; 1)$, the two-sample test is required to satisfy, for some $c_{TS} > 0$ depending on $(s; L; d; \epsilon; \delta)$ and on $\mu; \nu$ from Assumption 1 to be introduced next,

$$\begin{aligned} P(TS_d(V_1, \dots, V_{n_1}; W_1, \dots, W_{n_2}) = 1) & \leq \epsilon & \text{if } f_1 = f_2; \\ P(TS_d(V_1, \dots, V_{n_1}; W_1, \dots, W_{n_2}) = 0) & \leq c_{TS} \frac{n_1 \wedge n_2}{\log \log(n_1 \wedge n_2)} \frac{2s}{4s+d} & \text{if } k \|f_1 - f_2\|_{L^2(\cdot)} \geq \epsilon. \end{aligned} \quad (13)$$

There are a number of such tests proposed in prior work, see e.g., Ingster (2012); Arias-Castro et al. (2018); Kim et al. (2022) and Appendix B.2. Our general approach allows using any of these. It is also known that there are adaptive tests TS^{ad} that do not require knowing the Holder smoothness parameters. In the adaptive setting, the best-known minimum required separation in this general dimensional situation is

$$k \|f_1 - f_2\|_{L^2(\cdot)} \geq c_{ad} \frac{n_1 \wedge n_2}{\log \log(n_1 \wedge n_2)} \frac{2s}{4s+d} \quad (14)$$

for some $c_{ad} > 0$. See Appendix B.2 for examples of TS and TS^{ad} .

Next, we state an additional assumption required in our theorem. See Appendix B.2 for more discussion. Assumption 2 guarantees that every class appears in the dataset. This is reasonable in many practical settings, as classes that do not appear can be omitted.

Assumption 2 (Lower bounded class probability) There exists a constant $d_c > 0$ such that $E[Y]_k > d_c$ for all $k \in \{1, \dots, K\}$.

Our result is as follows.

Theorem 12 (Optimal calibration test via sample splitting) Suppose $\epsilon \in (0; 1)$ and let τ_n^{split} be the test described in Algorithm 3. Assume the Holder smoothness parameter s is known. Under Assumption 1 and 2, we have

1. False detection rate control. For every P for which f is perfectly calibrated, i.e., for $P \in \mathcal{P}_0$, the probability of falsely claiming mis-calibration is at most ϵ , i.e., $P(\tau_n^{split} = 1) \leq \epsilon$.

Algorithm 3 Sample splitting calibration test $\frac{\text{split}}{n}$

Input: Probability predictor $f : X \rightarrow \mathcal{K}$; i.i.d. sample $(X_i; Y_i)_{i=1}^n \in X \times Y$; $\alpha \in (0, 1)$; false detection rate $\beta \in (0, 1)$; true detection rate $\gamma \in (0, 1)$; Hölder smoothness s ; minimax optimal two-sample density test TS

Procedure: $Z_i = f(X_i)$ for $i = 1, \dots, n$; independently sample $\mathcal{Y}_i \sim \text{Cat}(Z_i)$ for $i = 1, \dots, n$

for $k = 1$ to K do
 $T_{1;k} = \frac{1}{n} \sum_{i=1}^n [Y_i = Z_i]_k$, $T_{2;k} = \frac{1}{n} \sum_{i=1}^n [Z_i = \mathcal{Y}_i]_k$
 $V_k = \frac{1}{n} \sum_{i=1}^n [Y_i = Z_i]_k - \frac{\beta}{2}$, $W_k = \frac{1}{n} \sum_{i=1}^n [Z_i = \mathcal{Y}_i]_k - \frac{\beta}{2}$
 $b_k = \frac{1}{2} \left(\frac{1}{n} \sum_{i=1}^n [Y_i = Z_i]_k - \frac{\beta}{2} \right) - \frac{1}{2} \left(\frac{1}{n} \sum_{i=1}^n [Z_i = \mathcal{Y}_i]_k - \frac{\beta}{2} \right) - \text{TS}_{\frac{3K}{n}, \frac{\beta}{2}}(V_k; W_k)$

end for

Output: Reject H_0 if $\frac{\text{split}}{n} := \max_{k \in \mathcal{K}} b_k \geq \alpha$

2. True detection rate control. There exists $c_{\text{split}} > 0$ depending on $(s; L; K; \alpha; \beta; \gamma; \beta; \gamma)$ such that the power (true positive rate) is bounded $\mathbb{P}(\frac{\text{split}}{n} = 1) \geq 1 - \frac{1}{c_{\text{split}} n^{2s/(4s+K-1)}}$ for every $P \in \mathcal{P}_1(\cdot; p; s)$ i.e., when f is mis-calibrated with an α -ECE of at least $\frac{1}{c_{\text{split}} n^{2s/(4s+K-1)}}$.

The proof is in Appendix A.9. Theorem 10 and Theorem 12 together imply that the minimax optimal detection rate for calibration is $\frac{1}{n} \alpha(p; s) \geq \frac{1}{n^{2s/(4s+K-1)}}$. By replacing TS with an adaptive test TS^{ad} , we obtain an adaptive version of the test $\frac{\text{split}}{n}$.

Corollary 13 (Adaptive test via sample splitting) Suppose $\alpha \in (0, 1)$ and let $\frac{\text{ad-s}}{n}$ be the test described in Algorithm 3 with TS replaced by an adaptive two-sample test TS^{ad} . Under Assumption 1 and 2, we have

1. False detection rate control. For every P for which f is perfectly calibrated, i.e., for $P \in \mathcal{P}_0$, the probability of falsely claiming mis-calibration is at most α , i.e., $\mathbb{P}(\frac{\text{ad-s}}{n} = 1) \leq \alpha$.
2. True detection rate control. There exists $c_{\text{ad-s}} > 0$ depending on $(s; L; K; \alpha; \beta; \gamma; \beta; \gamma)$ such that the power (true positive rate) is bounded $\mathbb{P}(\frac{\text{ad-s}}{n} = 1) \geq 1 - \frac{1}{c_{\text{ad-s}} (n \log n)^{2s/(4s+K-1)}}$ for every $P \in \mathcal{P}_1(\cdot; p; s)$ i.e., when f is mis-calibrated with an α -ECE of at least $\frac{1}{c_{\text{ad-s}} (n \log n)^{2s/(4s+K-1)}}$.

Remark 14 (Adaptation cost and optimality) Spokoiny (1996); Ingster (2000) develop an adaptive chi-squared test for one-dimensional goodness-of-fit testing which can adapt to an unknown Hölder smoothness parameter while only losing a $(\log \log n)^{s/(4s+1)}$ factor in the separation rate. The test was proven to be minimax optimal in the adaptive setting. Arias-Castro et al. (2018) extend the adaptive test to a general dimension and attain an adaptive test at the cost of $\log n$ factor. Kim et al. (2022) provide a stronger analysis for their permutation test and reduce the adaptation cost to $(\log \log n)^{2s/(4s+d)}$. To our knowledge, the minimax optimality of these adaptive tests in general dimensions is so far not established.

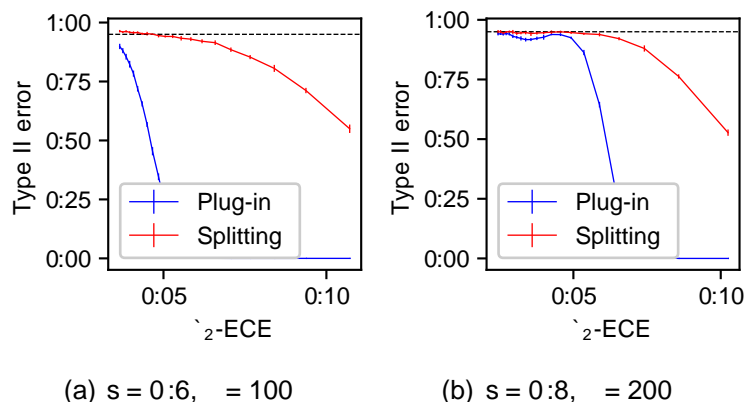


Figure 7: Type II error comparison for $\hat{T}_{m;n}$ and \hat{T}_n^{split} . The horizontal dashed line indicates a type II error of $1 - \alpha = 0.95$. Since \hat{T}_n^{split} relies on sampling splitting, its effective sample size is much smaller than that of the plug-in test. This results in higher type II errors as can be seen in the figure. Standard error bars are plotted over 10 repetitions.

While the adaptive test in Corollary 13 requires an additional factor of $(\log \log n)^{2s-(4s+K-1)}$ in the separation rate, Theorem 5 requires a factor of $(\log n)^{s-(4s+K-1)}$. This gap comes from the requirement (14) which we borrow from Kim et al. (2022). Theorem 6.1 and Lemma C.1 of Kim et al. (2022) develops combinatorial concentration inequalities to improve a polynomial dependency on n in the separation rate to a logarithmic dependency. This results in the $(\log \log n)^{2s-(4s+K-1)}$ factor in their adaptive test. Since our proof of Theorem 3 uses a quadratic tail bound from Chebyshev's inequality, the adaptation cost in Theorem 5 is $(\log n)^{s-(4s+K-1)}$. Currently, it appears challenging to improve the polynomial dependence for $\hat{T}_{m;n}^d$, due to its complicated conditional structure.

6.1 Comparison with the Debiased Plug-in Test

We compare the empirical performances of the debiased plug-in test $\hat{T}_{m;n}$ and the sample splitting test \hat{T}_n^{split} . As described in Section 4.1, we study the type II error against the fixed alternative where the mis-calibration is specified as $H_1 : (Z; Y) \sim P_{1;m}$, for various values of m . We use a sample size of $n = 20,000$ and pairs of Hölder smoothness and scaling parameter indicated in Figure 7. The critical value for $\alpha = 0.05$ and the corresponding type II error are estimated via 1,000 Monte Carlo simulations. For the sample splitting test \hat{T}_n^{split} , we use the chi-squared two-sample test of Arias-Castro et al. (2018).

Since \hat{T}_n^{split} relies on sample splitting and discards some of the observations, its effective sample size is smaller than that of the debiased plug-in test. For this reason, we find that $\hat{T}_{m;n}$ outperforms the sample splitting test by a large margin. While the sample splitting test reveals a theoretically interesting connection to two-sample density testing, it appears empirically suboptimal.

7 Conclusion

This paper studied the problem of testing model calibration from a finite sample. We analyzed the plug-in estimator of χ^2 -ECE(f)² as a test statistic for calibration testing. We discovered that the estimator needs debiasing and becomes minimax optimal when the number of bins is chosen appropriately. We also provided an adaptive version of the test, which can be used without knowing the Hölder smoothness parameters. We tested T-Cal with a broad range of experiments, including several neural net architectures and post-hoc calibration methods.

On the theoretical side, we provided an impossibility result for testing calibration against general continuous alternatives. Assuming that the calibration curve is α -Hölder-smooth, we derived a lower bound of $(n^{-2s/(4s+K-1)})$ on the calibration error required for a model to be distinguished from a perfectly calibrated one. We also discussed a reduction to two-sample testing and showed that the resulting test also matches the lower bound.

Interesting future directions include (1) developing a testing framework for comparing calibration of predictive models, (2) extending the theoretical result to χ^p -ECE with $p > 2$ and other calibration concepts such as top(k), within- k , and marginal calibration, (3) developing a minimax estimation theory for calibration error, and (4) establishing local rates for testing calibration.

Acknowledgments and Disclosure of Funding

This work was supported in part by the NSF TRIPODS 1934960, NSF DMS 2046874 (CAREER), ARO W911NF-20-1-0080, DCIST, and NSF CAREER award CIF-1943064, and Air Force Office of Scientific Research Young Investigator Program (AFOSR-YIP) #FA9550-20-1-0111 award. We thank the editors and associate editor for their work in handling our manuscript; and the reviewers for their very thorough reading and many helpful suggestions that have significantly improved our work. We are grateful to a number of individuals for their comments, discussion, and feedback; in particular, we would like to thank Sivaraman Balakrishnan, Chao Gao, Chirag Gupta, Lucas Janson, Adel Javanmard, Edward Kennedy, Ananya Kumar, Yuval Kluger, Mohammed Mehrabi, Alexander Podkopaev, Yury Polyanskiy, Mark Tygert, and Larry Wasserman.

Appendix A. Proofs

Notations. For completeness and to help the reader, we present (and in some cases recall) some notations. We will use the symbols $:=$ or $=:$ to define quantities in equations. We will occasionally use bold font for vectors. For an integer $d \geq 1$, we denote $\mathbb{I}^d := \{1, \dots, d\}$ and $1_d := (1, 1, \dots, 1) \in \mathbb{R}^d$. For a vector $v \in \mathbb{R}^d$, we will sometimes write $[v]_i$ for the i -th coordinate of v , for any $i \in \mathbb{I}^d$. The minimum of two scalars $a, b \in \mathbb{R}$ is denoted by $\min(a, b)$ or $a \wedge b$; their maximum is denoted by $\max(a, b)$ or $a \vee b$. We denote the d -dimensional Lebesgue measure by Leb_d . For a function $h : \mathbb{R}^d \rightarrow \mathbb{R}$, $1 \leq p < \infty$, and a measure μ on \mathbb{R}^d , we let $\|h\|_{L^p(\mu)} := (\int |h|^p d\mu)^{1/p}$. When $\mu = \text{Leb}_d$, we omit μ and write $\|h\|_{L^p}$. If $p = \infty$, then $\|h\|_{L^p} := \text{ess sup}_{x \in \mathbb{R}^d} |h(x)|$. We denote the p -norm of $x = (x_1, \dots, x_d) \in \mathbb{R}^d$ by $\|x\|_p := (\sum_{i=1}^d |x_i|^p)^{1/p}$. When p is unspecified, $\|x\|$ stands for $\|x\|_2$.

For two sequences $(a_n)_{n \geq 1}$ and $(b_n)_{n \geq 1}$ with $b_n \neq 0$, we write $a_n \sim b_n$ if $0 < \liminf_n a_n/b_n \leq \limsup_n a_n/b_n < \infty$. When the index n is self-evident, we may omit it above. We use the Bachmann-Landau asymptotic notations $\mathcal{O}(\cdot)$; $\mathcal{O}(\cdot)$ to hide constant factors in inequalities and use $\mathcal{O}(\cdot)$; $\mathcal{O}(\cdot)$ to also hide logarithmic factors. For a Lebesgue measurable set $A \subset \mathbb{R}^d$, we denote by $1_A : \mathbb{R}^d \rightarrow \{0, 1\}$ its indicator function where $1_A(x) = 1$ if $x \in A$ and $1_A(x) = 0$ otherwise. For a real number $s \in \mathbb{R}$, we denote the largest integer less than or equal to s by $\lfloor s \rfloor$. Also, the smallest integer greater than or equal to s is denoted by $\lceil s \rceil$.

For an integer $d \geq 1$, a vector $j = (j_1, \dots, j_d) \in \mathbb{N}^d$ is called a multi-index. We write $|j| := j_1 + \dots + j_d$. For a vector $x = (x_1, \dots, x_d) \in \mathbb{R}^d$ and a multi-index $j = (j_1, \dots, j_d) \in \mathbb{N}^d$, we write $x^j := x_1^{j_1} \dots x_d^{j_d}$. For a sufficiently smooth function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we denote its partial derivative of order $j = (j_1, \dots, j_d) \in \mathbb{N}^d$ by $f^{(j)} := \partial_1^{j_1} \dots \partial_d^{j_d} f$. A pure partial derivative with respect to an individual coordinate $i \in \mathbb{I}^d$ is also denoted as $\partial_i f$. For two sets S, T , a map $f : S \rightarrow T$, and a subset $S^0 \subset S$, we denote by $f(S^0)$ the image of S^0 under f . The support of a function $f : S \rightarrow \mathbb{R}$ is the set of points $\text{supp} f := \{a \in S : f(a) \neq 0\}$.

The uniform distribution on a compact set $S \subset \mathbb{R}^d$ is denoted by $\text{Unif}(S)$. The binomial distribution with $n \in \mathbb{N}$ trials and success probability $p \in [0, 1]$ is denoted by $\text{Bin}(n; p)$, and write $\text{Ber}(p) := \text{Bin}(1; p)$. For an integer $d \geq 2$, we let $\Delta_{d-1} := \{z = (z_1, \dots, z_d) \in [0, 1]^d : z_1 + \dots + z_d = 1\}$ be the $(d-1)$ -dimensional probability simplex. We denote the multinomial distribution with $n \in \mathbb{N}$ trials and class probability vector $p \in \Delta_{d-1}$ by $\text{Multi}(n; p)$, and write $\text{Cat}(p) := \text{Multi}(1; p)$. For a joint distribution $(X, Y) \in \mathcal{P}$, we will write P_X, P_Y for the marginal distributions of X, Y , respectively. For a distribution Q and a random variable $Z \in Q$, we will denote expectations of functions of Z with respect to Q as $E f(Z)$, $E_Z f(Z)$, $E_Q f(Z)$, or $E_Z^Q f(Z)$. We abbreviate almost surely by "a.s.", and almost everywhere by "a.e."

A.1 Proof of Theorem 3

We first state a Lemma used in the proof. This lemma generalizes Lemma 3 in Arias-Castro et al. (2018) from the uniform measure on the cube to a general probability measure on the probability simplex. See Section 3.2.2 of Ingster (2012) for a discussion of how such results connect to geometric notions like Kolmogorov diameters.

Lemma 15 For $m \geq 2$, let $B_m = \{B_1, \dots, B_{m-1}\}$ be the partition of \mathcal{K} defined in Appendix B.3, and μ be a probability measure on \mathcal{K} such that $\mu(B_i) > 0$ for all $i \in [m-1]$. For any continuous function $h : \mathcal{K} \rightarrow \mathbb{R}$, define

$$W_m[h] := \sum_{i=1}^{m-1} \frac{\int_{B_i} h(z) d\mu(z)}{\mu(B_i)} \mathbf{1}_{B_i}.$$

There are $b_1, b_2 > 0$ depending on $(K; s; L)$ such that for every $h \in H_{\mathcal{K}}(s; L)$,

$$\|W_m[h]\|_{L^2(\mu)} \leq b_1 \|h\|_{L^2(\mu)} \leq b_2 m^{-s}.$$

In other words,

$$\sum_{i=1}^{m-1} \frac{\int_{B_i} [h(Z) - \int_{B_i} h(z) d\mu(z)]^2 d\mu(z)}{\mu(B_i)} \leq b_1 \int_{\mathcal{K}} [h(z)]^2 d\mu(z) \leq b_2 m^{-s}.$$

The proof can be found in Appendix A.2.

Overview of the proof. The proof follows the classical structure of upper bound arguments in nonparametric hypothesis testing, see e.g., Arias-Castro et al. (2018); Kim et al. (2022) for recent examples. We compute or bound the mean and variance $\mathbb{E}_{P_0} T_{m;n}^d$ under null distributions $P_0 \in \mathcal{P}_0$ and alternative distributions $P_1 \in \mathcal{P}_{s;L;K}$ with a large ECE. Using Lemma 15, we can find a lower bound on $\mathbb{E}_{P_1} T_{m;n}^d - \mathbb{E}_{P_0} T_{m;n}^d$. The variances $\text{Var}_{P_0}(T_{m;n}^d)$ and $\text{Var}_{P_1}(T_{m;n}^d)$ can be also upper bounded. We argue that the mean difference $\mathbb{E}_{P_1} T_{m;n}^d - \mathbb{E}_{P_0} T_{m;n}^d$ is significantly larger than the square root of the variances $\text{Var}_{P_0}(T_{m;n}^d)$ and $\text{Var}_{P_1}(T_{m;n}^d)$. The conclusion follows from Chebyshev's inequality.

Proof Let $N_i := |I_{m;i}|$ for each $i \in [m-1]$ and $I := \{i \in [m-1] : N_i \geq 1\}$. Also write $N := (N_1, \dots, N_{m-1})^\top$, $Z := (Z_1, \dots, Z_n)^\top$, $\bar{Y} := \bar{Y} - Z$, and $\bar{Y}_j := Y_j - Z_j$ for all $j \in [n]$. By Assumption 1,

$$\sum_{i \in I} \int_{\mathcal{K}} k_{\text{res}}(z) k^2 dP_Z(z) \leq \sum_{i \in I} \int_{\mathcal{K}} k_{\text{res}}(z) k^2 dz;$$

where the latter integral is with respect to the uniform measure $\text{Unif}(\mathcal{K})$ on \mathcal{K} . Therefore, we may assume $P_Z = \text{Unif}(\mathcal{K})$ by merging μ with c . We prove the theorem for $p = 2$. Then, the general case follows since $P_1(\cdot; p; s) \leq P_1(\cdot; 2; s)$ for all $p \geq 2$.

Let $P_0 \in \mathcal{P}_0$ and $P_1 \in \mathcal{P}_{s;L;K}$ be a null and an alternative distribution over $(Z; Y)$, respectively. Write $\mu := \mathbb{E}_{P_1} f$. Under P_0 and conditioned on Z , recalling $T_{m;n}^d$ from (6),

$$\mathbb{E}_{P_0} [T_{m;n}^d | Z] = \frac{1}{n} \sum_{i \in I} \frac{1}{N_i} \sum_{j \in I_{m;i}} h(\bar{Y}_{j_1} - \bar{Y}_{j_2} | Z) = 0$$

because $\mathbb{E}_{P_0} [\bar{Y}_{j_1} - \bar{Y}_{j_2} | Z] = \mathbb{E}_{P_0} [\bar{Y}_{j_1} | Z] - \mathbb{E}_{P_0} [\bar{Y}_{j_2} | Z] = 0$ for all $j_1 \in I_{m;i_1}, j_2 \in I_{m;i_2}$. Therefore,

$$\mathbb{E}_{P_0} [T_{m;n}^d] = \mathbb{E}_{P_0} [\mathbb{E}_{P_0} [T_{m;n}^d | Z]] = 0: \tag{15}$$

Also,

$$\begin{aligned} \text{Var}_{P_0}(T_{m;n}^d | Z) &= \frac{1}{n^2} \sum_{i=1}^m \frac{1}{N_i^2} \text{Var}_{P_0} \left(\sum_{j_1 < j_2 \in I_{m;i}} \bar{Y}_{j_1} \bar{Y}_{j_2} | Z \right) \\ &= \frac{1}{n^2} \sum_{i=1}^m \frac{4}{N_i^2} \sum_{j_1 < j_2 \in I_{m;i}} \text{Var}_{P_0} \left(\bar{Y}_{j_1} \bar{Y}_{j_2} | Z \right); \end{aligned}$$

Here we used that the cross terms in the expansion of $\text{Var}_{P_0} \left(\sum_{j_1 < j_2 \in I_{m;i}} \bar{Y}_{j_1} \bar{Y}_{j_2} | Z \right)$ vanish since for $j_1, j_2, j_3 \in I_{m;i}$ that are pairwise different,

$$\begin{aligned} \text{Cov}_{P_0} \left(\bar{Y}_{j_1} \bar{Y}_{j_2}, \bar{Y}_{j_1} \bar{Y}_{j_3} | Z \right) &= E_{P_0} \left[\bar{Y}_{j_2} \bar{Y}_{j_1} \bar{Y}_{j_3} | Z \right] - E_{P_0} \left[\bar{Y}_{j_1} \bar{Y}_{j_2} | Z \right] E_{P_0} \left[\bar{Y}_{j_1} \bar{Y}_{j_3} | Z \right] \\ &= E_{P_0} \left[\bar{Y}_{j_2} | Z \right] E_{P_0} \left[\bar{Y}_{j_1} \bar{Y}_{j_3} | Z \right] - E_{P_0} \left[\bar{Y}_{j_1} \bar{Y}_{j_2} | Z \right] E_{P_0} \left[\bar{Y}_{j_1} \bar{Y}_{j_3} | Z \right] = 0; \end{aligned}$$

and for $j_1, j_2, j_3, j_4 \in I_{m;i}$ that are pairwise different, $\text{Cov}_{P_0} \left(\bar{Y}_{j_1} \bar{Y}_{j_2}, \bar{Y}_{j_3} \bar{Y}_{j_4} | Z \right) = 0$ by the independence of $\bar{Y}_{j_1}, \bar{Y}_{j_2}, \bar{Y}_{j_3}, \bar{Y}_{j_4}$ given Z . Further, since $\text{Var}_{P_0} \left(\bar{Y}_{j_1} \bar{Y}_{j_2} | Z \right) = K^2$ for all $j_1 < j_2 \in I_{m;i}$,

$$\text{Var}_{P_0}(T_{m;n}^d | Z) = \frac{1}{n^2} \sum_{i=1}^m \frac{2K^2 N_i (N_i - 1)}{N_i^2} = 2K^2 n^{-2} \sum_{i=1}^m \binom{N_i - 1}{1} = 2K^2 n^{-2} \sum_{i=1}^m (N_i - 1);$$

Thus, by the law of total variance,

$$\begin{aligned} \text{Var}_{P_0}(T_{m;n}^d) &= E_{P_0}[\text{Var}_{P_0}(T_{m;n}^d | Z)] + \text{Var}_{P_0}(E_{P_0}[T_{m;n}^d | Z]) = 2K^2 n^{-2} \sum_{i=1}^m P_0(N_i - 1) \\ &= 2K^2 m^{K-1} n^{-2} P_0(N_1 - 1); \end{aligned}$$

Since N_1 follows a Binomial distribution with n trials and success probability m^{K+1} , and as $(1-x)^{n-1} = 1 - (n-1)x$ for any $x \in [0, 1]$, we see that

$$P_0(N_1 - 1) = 1 - \frac{1}{m^{K+1}} = 1 + \frac{n-1}{m^{K+1}} = 1 + \frac{n^2}{m^{2(K+1)}}; \quad (16)$$

Therefore, denoting β below,

$$\text{Var}_{P_0}(T_{m;n}^d) = 2K^2 m^{K-1} n^{-2} \left(1 + \frac{n^2}{m^{2(K+1)}} \right) =: \beta; \quad (17)$$

Under P_1 , we have

$$\begin{aligned} E_{P_1}[T_{m;n}^d | Z] &= \frac{1}{n} \sum_{i=1}^m \frac{1}{N_i} \sum_{j_1 < j_2 \in I_{m;i}} E_{P_1} \left(\bar{Y}_{j_1} \bar{Y}_{j_2} | Z \right) \\ &= \frac{1}{n} \sum_{i=1}^m \frac{1}{N_i} \sum_{j_1 < j_2 \in I_{m;i}} \text{res}(Z_{j_1}) > \text{res}(Z_{j_2}); \end{aligned}$$

since for each $i \in \{1, \dots, m\}$, $E_{P_1}[\bar{Y}_{j_1} > \bar{Y}_{j_2} | Z] = E_{P_1}[\bar{Y}_{j_1} | Z] > E_{P_1}[\bar{Y}_{j_2} | Z] = \text{res}_f(Z_{j_1}) > \text{res}_f(Z_{j_2})$ for all $j_1 \in \{1, \dots, m\}$. Moreover,

$$\begin{aligned} E_{P_1}[T_{m,n}^d | N] &= E_{P_1}[E_{P_1}[T_{m,n}^d | Z] | N] \\ &= \frac{1}{n} \sum_{i=1}^m \frac{N_i(N_i - 1)}{N_i} E_{P_1}[\text{res}_f(Z) | Z \in B_i] > E_{P_1}[\text{res}_f(Z) | Z \in B_i] \\ &= \frac{1}{n} \sum_{i=1}^m (N_i - 1) k E_{P_1}[\text{res}_f(Z) | Z \in B_i]^2 \end{aligned} \quad (18)$$

and

$$\begin{aligned} E_{P_1}[T_{m,n}^d] &= E_{P_1}[E_{P_1}[T_{m,n}^d | N]] \\ &= \frac{1}{n} \sum_{i=1}^m E_{P_1}[(N_i - 1)(N_i - 1)] k E_{P_1}[\text{res}_f(Z) | Z \in B_i]^2 \end{aligned}$$

For any $x \in [0, 1]$, we have $1 - nx + \frac{n}{2} x^2 \geq \frac{n}{3} x^3 - (1 - x)^n - 1 - nx + \frac{n}{2} x^2$: Applying the inequality for $x = \frac{1}{m^{(K-1)}}$, we derive

$$\begin{aligned} \frac{1}{4} \frac{n}{m^{K-1}} \wedge \frac{n^2}{m^{2(K-1)}} &\leq \frac{n}{2} \frac{1}{m^{2(K-1)}} - \frac{n}{3} \frac{1}{m^{3(K-1)}} - \frac{n}{m^{K-1}} + 1 \\ E_{P_1}[(N_i - 1)(N_i - 1)] &= \frac{n}{m^{K-1}} + 1 - \frac{1}{m^{K-1}} - \frac{n}{2} \frac{1}{m^{2(K-1)}} - \frac{n^2}{m^{2(K-1)}} \end{aligned}$$

Also, since $(1 - x)^n \leq 1$ for $x \in [0, 1]$, we have

$$E_{P_1}[(N_i - 1)(N_i - 1)] \geq \frac{n}{m^{K-1}}$$

Therefore,

$$\frac{1}{4} \frac{n}{m^{K-1}} \wedge \frac{n^2}{m^{2(K-1)}} \leq E_{P_1}[(N_i - 1)(N_i - 1)] \leq \frac{n}{m^{K-1}} \wedge \frac{n^2}{m^{2(K-1)}} \quad (19)$$

By Lemma 15, and as $\| \cdot \| = \left(\sum_{k=1}^K E_{P_1}[[\text{res}_f(Z)]_k^2] \right)^{1/2} = \left(\sum_{k=1}^K (E_{P_1}[[\text{res}_f(Z)]_k^2])^{1/2} \right)^2$ by the Cauchy-Schwarz inequality,

$$\begin{aligned} \sum_{i=1}^m \frac{1}{m^{(K-1)}} k E_{P_1}[\text{res}_f(Z) | Z \in B_i]^2 &= \sum_{i=1}^m \sum_{k=1}^K \frac{1}{m^{(K-1)}} E_{P_1}[\text{res}_f(Z)_k^2 | Z \in B_i] \\ &= \sum_{k=1}^K W_m^2 [[\text{res}_f(Z)]_k] \sum_{k=1}^K b_k E_{P_1}[\text{res}_f(Z)_k^2]^{1/2} b_k m^{-s/2} \\ &= \sum_{k=1}^K b_k^2 E_{P_1}[\text{res}_f(Z)_k^2] \sum_{k=1}^K b_k b_k E_{P_1}[\text{res}_f(Z)_k^2]^{1/2} m^{-s} = b_1^{2n} \sum_{k=1}^K b_k b_k m^{-s} \quad (20) \end{aligned}$$

By (19) and (20), defining below,

$$E_{P_1}[T_{m;n}^d] = \frac{1}{4} (b_1^{2n} - 2^p \bar{K} b_1 b_2 m^{-s}) \wedge \frac{n}{m^{K-1}} =: \quad (21)$$

Moreover, we find

$$\begin{aligned} \text{Var}_{P_1}(T_{m;n}^d | N) &= \frac{1}{n^2} \sum_{i=2}^X \frac{1}{N_i^2} \text{Var}_{P_1} \left[\sum_{j_1 < j_2 \leq m_i} \bar{Y}_{j_1} \bar{Y}_{j_2} \right] N_i^2 \\ &= \frac{1}{n^2} \sum_{i=2}^X \frac{1}{N_i^2} \left[2N_i(N_i - 1) \text{Var}_{P_1}(\bar{Y}_1 \bar{Y}_2 | Z_1; Z_2) + 2N_i(N_i - 1)(N_i - 2) \text{Cov}_{P_1}(\bar{Y}_1 \bar{Y}_2; \bar{Y}_1 \bar{Y}_3 | Z_1; Z_2; Z_3) \right] \\ &= \frac{1}{n^2} \sum_{i=2}^X \left[\frac{2K^2 N_i(N_i - 1)}{N_i^2} + \frac{4K^2 N_i(N_i - 1)(N_i - 2)}{N_i^2} \right] k E_{P_1}[\text{res}_q(Z) | Z \geq B_i]^2 \quad (22) \end{aligned}$$

Further, by equations (16), (19) and since $k E_{P_1}[\text{res}_q(Z) | Z \geq B_i]^2 = E_{P_1}[\text{res}_q(Z)^2 | Z \geq B_i]$,

$$\begin{aligned} E_{P_1}[\text{Var}_{P_1}(T_{m;n}^d | N)] &= \frac{1}{n^2} \sum_{i=1}^{m^{K-1}} \left[2K^2 P_1(N_i \geq 2) + 4K^2 E_{P_1}[(N_i - 1)(N_i - 2)] \right] k E_{P_1}[\text{res}_q(Z) | Z \geq B_i]^2 \\ &= 2 + 4K^2 m^{K-1} n^{-2} E_{P_1}[(N_1 - 1)(N_1 - 2)] \sum_{i=1}^{m^{K-1}} E_{P_1}[\text{res}_q(Z)^2 | Z \geq B_i] \\ &= 2 + 4K^2 n^{-2} (n \wedge m^{K-1}); \end{aligned} \quad (22)$$

Also, from (18),

$$\begin{aligned} \text{Var}_{P_1}(E_{P_1}[T_{m;n}^d | N]) &= \frac{1}{n^2} \sum_{i=1}^{m^{K-1}} \frac{1}{N_i^2} \text{Var}_{P_1} \left[\sum_{j_1 < j_2 \leq m_i} \bar{Y}_{j_1} \bar{Y}_{j_2} \right] N_i^2 \\ &= \frac{1}{n^2} \text{Var}_{P_1} [I(N_1 \geq 2)] m^{2(K-1)}. \end{aligned}$$

Writing $x = m^{-(K-1)}$, we have

$$\text{Var}_{P_1} [I(N_1 \geq 2)] = nx(1-x) + (1-x)^n - (1-x)^{2n} = 2nx(1-x)^n;$$

Using that $1 - nx = (1-x)^n$, we find

$$\begin{aligned} \text{Var}_{P_1} [I(N_1 \geq 2)] &= nx(1-x) + (1-x)^n [1 - (1-x)^n - 2nx] \\ &= nx(1-x) - nx = \frac{n}{m^{K-1}}. \end{aligned}$$

Similarly,

$$\begin{aligned} \text{Var}_{P_1}[I(N_i = 1)(N_i - 1)] &= nx + (1 - x)^n - 1[1 - (1 - x)^n - 2nx] \\ &= nx + (1 - (n - 1)x)(-nx) - n^2x^2 = \frac{n^2}{m^{2(K-1)}}: \end{aligned}$$

Therefore,

$$\text{Var}_{P_1}(E_{P_1}[T_{m;n}^d | N]) = \frac{m^{K-1}}{n} \wedge 1^{-4}: \quad (23)$$

By equations (22), (23), and the law of total variance, defining ϵ^2 below,

$$\begin{aligned} \text{Var}_{P_1}(T_{m;n}^d) &= \text{Var}_{P_1}(E_{P_1}[T_{m;n}^d | N]) + E_{P_1}[\text{Var}_{P_1}(T_{m;n}^d | N)] \\ &= \frac{m^{K-1}}{n} \wedge 1^{-4} + 4K^2 \epsilon^2 (n^{-1} \wedge m^{-(K-1)}) + \frac{m^{K-1}}{n} \wedge 1^{-4} =: \epsilon^2: \end{aligned} \quad (24)$$

Recalling that $m = bn^{2s+(4s+K-1)c}$, we choose $c > 0$ such that $\epsilon < cn^{2s+(4s+K-1)}$ implies

$$\begin{aligned} \frac{r}{2} + \frac{r}{2} &= K(m^{\frac{K-1}{2}} n^{-1} \wedge m^{\frac{K-1}{2}}) + \frac{r}{4} K n^{-\frac{1}{2}} \wedge m^{\frac{K-1}{2}} + \frac{r}{2} \frac{m^{K-1}}{n} \wedge 1^{-\frac{1}{2}} \epsilon^2 \\ &= \frac{1}{4} (b_1^{2s+2} - 2^p \overline{K} b_1 b_2 m^{-s}) \wedge 1 \wedge \frac{n}{m^{K-1}} \end{aligned}$$

for all large enough n . For ϵ , δ , and η from (17), (21), and (24), this gives

$$\rho = \epsilon + \delta =: \eta: \quad (25)$$

By equations (15), (17), the definition of $T_{m;n}$ from Algorithm (1), and Chebyshev's inequality,

$$P_0(T_{m;n} = 1) = P_0(T_{m;n}^d = \rho) = \frac{\text{Var}_{P_0}(T_{m;n}^d)}{2\epsilon} =: \eta: \quad (26)$$

By equations (21), (24), (25), and Chebyshev's inequality,

$$\begin{aligned} P_1(T_{m;n}^d < \rho) &= P_1(T_{m;n}^d - E_{P_1}[T_{m;n}^d] < -\rho) \\ &= P_1(T_{m;n}^d - E_{P_1}[T_{m;n}^d] > \rho) = \frac{\text{Var}_{P_1}(T_{m;n}^d)}{2\epsilon} =: \eta: \end{aligned} \quad (27)$$

By the above arguments, Theorem 3 holds for all $n \geq N$, where $N \geq 2N_+$ depends on $(s; L; K; \epsilon; \delta; \eta)$: If we require c to further satisfy $c \geq N^{2s+(4s+K-1)}$, then the family $P_1(\cdot; p; s)$ is empty for $n < N$ given $\epsilon < cn^{2s+(4s+K-1)} > 1$. Therefore, Theorem 3 becomes vacuously true for $n < N$, and thereby true for all $n \geq N$. This finishes the proof. ■

A.2 Proof of Lemma 15

We state and prove Lemma 16 and 17 which we use in the proof of Lemma 15.

Lemma 16 Fix $h \in H_K(s; L)$ and $z_0 \in \mathbb{C}^{K-1}$. Let u be the $(d_{se} - 1)$ -th order Taylor series of h at z_0 . There is L^0 depending on $(K; s; L)$ such that

$$\|h(z) - u(z)\| \leq L^0 k \|z - z_0\|^s \tag{28}$$

for all $z \in \mathbb{C}^{K-1}$:

Proof Let $\alpha = \alpha_K$ as in (6.2). Recall that for a multi-index $\alpha = (\alpha_1; \dots; \alpha_{K-1}) \in \mathbb{Z}^{K-1}$, we write $(z - z_0)^\alpha = \prod_{j=1}^{K-1} (z_j - z_{0j})^{\alpha_j}$. By a Taylor series expansion, there exists $t \in [0; 1]$ such that

$$\begin{aligned} h(z) - u(z) &= \sum_{|\alpha| \geq d_{se}} \frac{(h^{(\alpha)}(z_0))}{|\alpha|!} (z - z_0)^\alpha \\ &+ \sum_{|\alpha| = d_{se}} \frac{(h^{(\alpha)}(t(z) + (1-t)z_0))}{|\alpha|!} (z - z_0)^\alpha : \end{aligned}$$

Then, $h(z) - u(z)$ equals

$$\sum_{|\alpha| = d_{se}} \frac{(h^{(\alpha)}(t(z) + (1-t)z_0)) - (h^{(\alpha)}(z_0))}{|\alpha|!} (z - z_0)^\alpha :$$

By the triangle inequality and the s -Hölder continuity of h ,

$$\begin{aligned} \|h(z) - u(z)\| &\leq \sum_{|\alpha| = d_{se}} \frac{\| (h^{(\alpha)}(t(z) + (1-t)z_0)) - (h^{(\alpha)}(z_0)) \|}{|\alpha|!} \| (z - z_0)^\alpha \| \\ &\leq \sum_{|\alpha| = d_{se}} \frac{L (k \|z - z_0\|)^{s_{d_{se}} + 1}}{|\alpha|!} k \|z - z_0\|^{d_{se}} \\ &\leq L^0 k \|z - z_0\|^s \end{aligned}$$

■

Lemma 17 Let P_q^K be the class of polynomials of degree at most q . There are $a_1; a_2 > 0$ depending on $(K; q)$ such that

$$\|W_m[v]\|_{L^2(\cdot)} \leq a_1 \|v\|_{L^2(\cdot)} \tag{29}$$

for every $v \in P_q^K$ and $m \geq a_2$.

Proof If (29) does not hold, then we can find a sequence $\{m_k\}_{k=1}^\infty$ increasing to infinity and a sequence of polynomials $\{v_k\}_{k=1}^\infty \in P_q^K$ such that $kW_{m_k}[v_k]k_{L^2(\cdot)} < \frac{1}{k} kv_kk_{L^2(\cdot)}$. Dividing v_k by $kv_kk_{L^2(\cdot)}$, we may assume $kv_kk_{L^2(\cdot)} = 1$. Now $\{v_k\}_{k=1}^\infty \subset P_q^K : kv_kk_{L^2(\cdot)} = 1$ is compact in the topology induced by the norm $k_{L^2(\cdot)}$, due to the Heine-Borel theorem because it is closed and bounded. Thus, we can find a convergent subsequence $\{v_{k_j}\}_{j=1}^\infty$. Denote the limit by v_1 . On one hand,

$$W_{m_k}[v_1]_{L^2(\cdot)} = W_{m_k}[v_1 - v_k]_{L^2(\cdot)} + W_{m_k}[v_k]_{L^2(\cdot)} \leq kv_1 - v_kk_{L^2(\cdot)} + \frac{1}{k} \rightarrow 0:$$

Here, we used that

$$kW_m[v]k_{L^2(\cdot)}^2 = \sum_{i=1}^{m^{\kappa-1}} \frac{\int_{B_i} v(z)^2 d(z)}{|B_i|} = \sum_{i=1}^{m^{\kappa-1}} \int_{B_i} v(z)^2 d(z) = kvk_{L^2(\cdot)}^2: \quad (30)$$

On the other hand, $kW_{m_k}[v_1]k_{L^2(\cdot)} \neq kv_1k_{L^2(\cdot)} = 1$ since $W_{m_k}[v_1] \neq v_1$ a.e. and $kW_{m_k}[v_1]k_{L^1} \leq kv_1k_{L^1}$. The conclusion follows due to the contradiction. \blacksquare

Now we proceed with the proof of Lemma 15. Since α_1 does not depend on m , for sufficiently large m , we can choose ϵ such that $m\epsilon$ is an integer and $m\epsilon \leq \alpha_2$. Partition $\mathcal{K} \setminus \mathcal{K}_1$ into $M := (m\epsilon)^{\kappa-1}$ simplices as described in Appendix B.3 and call them $\mathcal{B}_1, \dots, \mathcal{B}_M$. By this construction, we can ensure that each $\mathcal{B}_j, j \in [M]$, consists of $r^{\kappa-1}$ different simplices B_i (also defined in Appendix B.3). For $j \in [M]$, let u_j be the $(d\epsilon - 1)$ -th order Taylor expansion of h at an arbitrary vertex of \mathcal{B}_j . Define $u := \sum_{j=1}^M u_j \mathbb{1}_{\mathcal{B}_j}$. By equation (28),

$$\|h(z) - u(z)\|_{L^0} \leq \text{diam}(\mathcal{B}_1)^s = L^0 \text{diam}(\mathcal{K} \setminus \mathcal{K}_1)^s \frac{r}{m}^s =: bm^{-s} \quad (31)$$

for all $z \in \mathcal{K} \setminus \mathcal{K}_1$. Therefore, by (30) and Lemma 16,

$$\begin{aligned} kW_m[h]k_{L^2(\cdot)} &\leq kW_m[u]k_{L^2(\cdot)} + kW_m[h - u]k_{L^2(\cdot)} \\ &\leq kW_m[u]k_{L^2(\cdot)} + \|u - h\|_{L^2(\cdot)} \leq kW_m[u]k_{L^2(\cdot)} + bm^{-s}. \end{aligned}$$

Note that

$$kW_m[u]k_{L^2(\cdot)}^2 = \sum_{j=1}^M W_m[u_j \mathbb{1}_{\mathcal{B}_j}]_{L^2(\cdot)}^2;$$

and that Lemma 17 with $q = d\epsilon - 1$ can be applied to \mathcal{B}_j and its $r^{\kappa-1}$ sub-simplices to get

$$W_m[u_j \mathbb{1}_{\mathcal{B}_j}]_{L^2(\cdot)}^2 \leq a_1^2 \|u_j \mathbb{1}_{\mathcal{B}_j}\|_{L^2(\cdot)}^2:$$

Thus,

$$kW_m[u]k_{L^2(\cdot)}^2 \leq \sum_{j=1}^M a_1^2 \|u_j \mathbb{1}_{\mathcal{B}_j}\|_{L^2(\cdot)}^2 = a_1^2 \|u\|_{L^2(\cdot)}^2:$$

In conclusion, combining the above inequalities, and by (31)

$kW_m[h]k_{L^2(\cdot)} \leq a_1 k u k_{L^2(\cdot)} + b m^{-s} \leq a_1 (k h k_{L^2(\cdot)} + b m^{-s}) + b m^{-s} =: b_1 k h k_{L^2(\cdot)} + b_2 m^{-s}$.
This finishes the proof of Lemma 15.

A.3 Proof of Remark 4

We follow the same strategy in Appendix A.1. For null distributions P_0 and alternative distributions P_1 such that $\int \rho - ECE_{P_0}(f) \leq c_0 n^{-2s/(4s+K-1)}$ and $\int := \int \rho - ECE_{P_1}(f) \geq c_1 n^{-2s/(4s+K-1)}$, we show the mean difference $E_{P_1}[T_{m,n}^d] - E_{P_0}[T_{m,n}^d]$ is larger than $\text{Var}_{P_0}(T_{m,n}^d)^{1/2}$ and $\text{Var}_{P_1}(T_{m,n}^d)^{1/2}$.

While $E_{P_0}[T_{m,n}^d] = 0$ under the null hypothesis of perfect calibration, we now have

$$\begin{aligned} E_{P_0}[T_{m,n}^d] &= \frac{1}{n} \sum_{i=1}^{m^K-1} E_{P_0}[I(N_i-1)(N_i-1)] E_{P_0}[\text{res}_\#(Z) | Z \in B_i] k^2 \\ &= (m^{-(K-1)} \wedge m^{-2(K-1)} n)^{\sum_{i=1}^{m^K-1}} E_{P_0}[\text{res}_\#(Z) k^2 | Z \in B_i] \\ &\geq c_0^2 n^{-\frac{4s}{4s+K-1}} (1 \wedge m^{-(K-1)} n). \end{aligned} \tag{32}$$

Therefore,

$$E_{P_1}[T_{m,n}^d] - E_{P_0}[T_{m,n}^d] \geq c_0^2 n^{-\frac{4s}{4s+K-1}} (1 \wedge m^{-(K-1)} n) =: \delta.$$

By the equation (24),

$$\text{Var}_{P_0}(T_{m,n}^d) \leq 2 + 5K^2 c_0^2 n^{-\frac{4s}{4s+K-1}} (n^{-1} \wedge m^{-(K-1)}) =: \epsilon^2.$$

Similar to (25), we can choose large enough $b_1 > 0$ such that

$$\delta \geq \epsilon.$$

The conclusion follows from Chebyshev's inequality as in (26) and (27).

A.4 Proof of Theorem 5

By the union bound, for $P \in \mathcal{P}_0$,

$$P\left(\bigcup_{b=1}^B \bar{B} = 1\right) \leq \sum_{b=1}^B P\left(\bar{B} = 1\right) \leq B \cdot \frac{1}{2^{b_0}}$$

There exists $b_0 \geq 1; \dots; B$ such that $2^{b_0-1} < (n \log n)^{2s/(4s+K-1)} \leq 2^{b_0}$. Let $m_0 = 2^{b_0}$ and repeat the argument in the proof of Theorem 3. The condition (25) for type II error control is now changed to

$$\frac{2}{-K m_0^{\frac{K-1}{2}} n^{-1} \log n} + \frac{1}{-K m_0^{\frac{K-1}{2}} n^{-1}} \frac{1}{4} (b_1^2 n^{2s} + 2^p \bar{K} b_1 b_2 m_0^{s''});$$

which is satisfied when $c_{\text{ad}}(n) = \frac{p}{\log n} \frac{2s-(4s+K-1)}{2s-(4s+K-1)}$ for a sufficiently large $c_{\text{ad}} > 0$. Assuming $c_{\text{ad}}(n) = \frac{p}{\log n} \frac{2s-(4s+K-1)}{2s-(4s+K-1)}$ and $P \geq P_1(\cdot; p; s)$, we have

$$P(\hat{\alpha}_n = 1) \geq P(m_{0;n} = 1) \geq 1 - \epsilon.$$

This finishes the proof.

A.5 Proof of Proposition 7

Overview of the proof. We repeat the computation in Appendix A.1. However, due to the bias term, now the mean difference $E_{P_1}[T_{m;n}^b] - E_{P_0}[T_{m;n}^b]$ cannot be lower bounded by a positive number. Instead, we prove that $E_{P_0}[T_{m;n}^b] \geq E_{P_1}[T_{m;n}^b]$ holds for all large enough n .

Proof We use the same notations as in Appendix A.1. Since

$$\begin{aligned} E_{P_0}[T_{m;n}^b | Z] &= \frac{1}{n} \sum_{i=1}^m \frac{1}{N_i} \sum_{j \in \mathcal{I}_{m,i}} \left(\frac{1}{4} \sum_{j \in \mathcal{I}_{2l_{m,i}}} \left(\frac{h_j}{N_j} \sum_{i=1}^m \bar{Y}_j^2 \right) + \sum_{j_1 \in \mathcal{I}_{2l_{m,i}}} \sum_{j_2 \in \mathcal{I}_{2l_{m,i}}} \bar{Y}_{j_1} \bar{Y}_{j_2} \right) \\ &= \frac{1}{n} \sum_{i=1}^m \frac{1}{N_i} \sum_{j \in \mathcal{I}_{2l_{m,i}}} Z_j^2; \end{aligned}$$

we have

$$E_{P_0}[T_{m;n}^b | N] = E_{P_0}[E_{P_0}[T_{m;n}^b | Z] | N] = \frac{1}{n} \sum_{i=1}^m E_{P_0} \left[\sum_{j \in \mathcal{I}_{2l_{m,i}}} Z_j^2 \right] = \sum_{i=1}^m B_i; \quad (33)$$

Similarly,

$$\begin{aligned} E_{P_1}[T_{m;n}^b | Z] &= \frac{1}{n} \sum_{i=1}^m \frac{1}{N_i} \sum_{j \in \mathcal{I}_{2l_{m,i}}} \left(\frac{1}{4} \sum_{j \in \mathcal{I}_{2l_{m,i}}} \left(\frac{h_j}{N_j} \sum_{i=1}^m \bar{Y}_j^2 \right) + \sum_{j_1 \in \mathcal{I}_{2l_{m,i}}} \sum_{j_2 \in \mathcal{I}_{2l_{m,i}}} \bar{Y}_{j_1} \bar{Y}_{j_2} \right) \\ &= \frac{1}{n} \sum_{i=1}^m \frac{1}{N_i} \sum_{j \in \mathcal{I}_{2l_{m,i}}} \left(\text{reg}_f(Z_j) + \text{reg}_f(Z_j)^2 + \text{res}_f(Z_j)^2 \right) + \sum_{j_1 \in \mathcal{I}_{2l_{m,i}}} \sum_{j_2 \in \mathcal{I}_{2l_{m,i}}} \text{res}_f(Z_{j_1}) \text{res}_f(Z_{j_2}); \end{aligned}$$

and thus

$$\begin{aligned} E_{P_1}[T_{m;n}^b | N] &= E_{P_1}[E_{P_1}[T_{m;n}^b | Z] | N] \\ &= \frac{1}{n} \sum_{i=1}^m E_{P_1} \left[\text{reg}_f(Z) + \text{reg}_f(Z)^2 + \text{res}_f(Z)^2 \right] + (N_i - 1) E_{P_1} [\text{res}_f(Z) | Z] \sum_{i=1}^m B_i; \end{aligned} \quad (34)$$

Since $Z \sim \text{Unif}([0; 1])$ under both P_0 and P_1 , the equations (33) and (34) imply

$$\begin{aligned} E_{P_0}[T_{m;n}^b | N] &\geq E_{P_1}[T_{m;n}^b | N] \\ &= \frac{1}{n} \sum_{i=1}^m E[\text{res}_f(Z)(2Z - 1) | Z] \sum_{i=1}^m B_i + (N_i - 1) E[\text{res}_f(Z) | Z] \sum_{i=1}^m B_i^2 \\ &= \frac{1}{n} \sum_{i=1}^m E[\text{res}_f(Z)(2Z - 1) | Z] \sum_{i=1}^m B_i + 2k \sum_{i=1}^m B_i^2. \end{aligned}$$

Here we used that $E[\text{res}_\#(Z) | Z \in B_i]^2 \leq k^2 k_{L,1}^2 m^{-2s}$ for all $i \in [m]$ and $P_{i \geq 2l} (N_i = 1) \leq n^{-2l}$. Taking total expectation,

$$\begin{aligned} E_{P_0}[T_m^b; n] - E_{P_1}[T_m^b; n] &= \frac{1}{n} \sum_{i=1}^m E[\text{res}_\#(Z)(2Z - 1) | Z \in B_i] \leq k^2 k_{L,1}^2 m^{-2s} \\ &= \frac{1}{n} \sum_{i=1}^m P(N_i = 1) E[\text{res}_\#(Z)(2Z - 1) | Z \in B_i] \leq k^2 k_{L,1}^2 m^{-2s} \\ &= \frac{1}{n} P(N_1 = 1) \sum_{i=1}^m E[\text{res}_\#(Z)(2Z - 1) | Z \in B_i] \leq k^2 k_{L,1}^2 m^{-2s}. \end{aligned} \tag{35}$$

From (9), we see that $E[\text{res}_\#(Z)(2Z - 1) | Z \in B_i] = 0$ for all $i \in [m]$. Thus,

$$\begin{aligned} \sum_{i=1}^m E[\text{res}_\#(Z)(2Z - 1) | Z \in B_i] &= \sum_{i=\frac{m}{4}+1}^m E[\text{res}_\#(Z)(2Z - 1) | Z \in B_i] \\ &= \frac{1}{4} \sum_{i=\frac{m}{4}+1}^m E[\text{res}_\#(Z) | Z \in B_i] = \frac{1}{32} k^2 k_{L,1}^2 m^{1-s}. \end{aligned} \tag{36}$$

Combining (35) and (36), we find

$$E_{P_0}[T_m^b; n] - E_{P_1}[T_m^b; n] \geq \frac{1}{32} P(N_1 = 1) k^2 k_{L,1}^2 m^{1-s} n^{-1} - k^2 k_{L,1}^2 m^{-2s}. \tag{37}$$

Since $m = bn^{2(4s+1)}$ and $\frac{2}{4s+1} < 1$, we find

$$\lim_{n \rightarrow \infty} P(N_1 = 1) = \lim_{n \rightarrow \infty} \frac{1}{n!} = 1.$$

Also, we have $m^{1-s} n^{-1} = n^{(1-6s)/(4s+1)}$ and $m^{-2s} = n^{-4s/(4s+1)}$ with $\frac{1-6s}{4s+1} > \frac{4s}{4s+1}$. In conclusion, the RHS of (37) is positive for all large enough n . ■

A.6 Ingster's method

Lemma 18 (Ingster's method for the lower bound) Let $P_0 \in \mathcal{P}_0$ and $P_1, \dots, P_M \in \mathcal{P}_1(\cdot; p; s)$ be probability distributions on $\mathcal{K} \times \mathcal{Y}$, and suppose that P_1, \dots, P_M are absolutely continuous with respect to P_0 . For an i.i.d. sample $f = (Z_i; Y_i)_{i=1}^n$ from P_0 , define the average likelihood ratio between P_1, \dots, P_M and P_0 as

$$L_n := \frac{1}{M} \sum_{i=1}^M \prod_{j=1}^n \frac{dP_i}{dP_0}(Z_j; Y_j)$$

If $E_{P_0}[L_n^2] \leq 1 + \epsilon$, then the minimax type II error (false negative rate) for testing $H_0 : P \in \mathcal{P}_0$ against $H_1 : P \in \mathcal{P}_1(\cdot; p; s)$ at level α satisfies $R_n(\cdot; p; s) \geq \alpha$ and the minimum separation rate to ensure type II error at most α obeys $\rho_n(\cdot; p; s) \leq \epsilon$.

The proof follows from the results of Ingster (1987, 2012); see also Lemma G.1 in Kim et al. (2022) for a very clear statement. By definition, it holds that

$$\begin{aligned} R_n(\epsilon; p; s) &= \inf_{\mathcal{P}_n(\epsilon)} \sup_{\mathcal{P}_2(\mathcal{P}_1(\epsilon; p; s))} E_P[1] \inf_{\mathcal{P}_n(\epsilon)} \frac{1}{M} \sum_{i=1}^M E_{P_i}[1] \\ &= \inf_{\mathcal{P}_n(\epsilon)} E_{P_0}[1] + \frac{1}{M} \sum_{i=1}^M E_{P_i}[1] - E_{P_0}[1] \\ &= 1 + \inf_{\mathcal{P}_n(\epsilon)} \frac{1}{M} \sum_{i=1}^M E_{P_i}[1] - E_{P_0}[1] \end{aligned}$$

where the last inequality holds because $E_{P_0}[1] = 1$. Further,

$$\frac{1}{M} \sum_{i=1}^M E_{P_i}[1] - E_{P_0}[1] = E_{P_0} \left[\frac{1}{M} \sum_{i=1}^M E_{P_i}[1] - 1 \right] = E_{P_0}[L_n]$$

by a change of variables and the Cauchy-Schwarz inequality. Therefore, we have $R_n(\epsilon; p; s) \leq \sqrt{E_{P_0}[L_n^2]}$. Finally, since $R_n(\epsilon; p; s)$ is non-increasing, we find $R_n(\epsilon; p; s) \leq \sqrt{E_{P_0}[L_n^2]}$.

A.7 Proof of Proposition 9

Overview of the proof. We construct distributions P_1, \dots, P_M over $(Z; Y)$ under which the predictor f has an ϵ -ECE of at least $\epsilon_0 = 0.1$. We can choose the mis-calibration curves of P_1, \dots, P_M to be orthogonal in L^2 , so that the cross terms in the expansion of $E_{P_0}[L_n^2]$ cancel out. By choosing M sufficiently large, we can ensure that $E_{P_0}[L_n^2]$ is at most $1 + (1 - \epsilon_0)^2$. The conclusion follows from Lemma 18.

Proof We prove Proposition 9 for the binary case. The generalization to the multi-class case follows the same argument and is omitted. The construction in this proof is inspired by Ingster (1987, 2000); Burnashev (1979). Let P_0 be a null distribution over $(Z; Y) \in [0; 1] \times \{0, 1\}$ defined as follows: the distribution of the predicted probabilities follows $Z \stackrel{P_0}{\sim} \text{Unif}([0; 1])$ and $P_0(Y = 1 | Z = z) = z$ for all $z \in [0; 1]$. Under P_0 , the probability predictor f is perfectly calibrated. For each $i \in [M]$, let

$$g_i(u) := \begin{cases} u + \frac{q}{3} \sin(2i\pi(u - \frac{1}{4})) & u \in [\frac{1}{4}, \frac{3}{4}]; \\ u & u \notin [\frac{1}{4}, \frac{3}{4}]; \end{cases}$$

and define P_i as follows: $Z \stackrel{P_i}{\sim} \text{Unif}([0; 1])$ and $P_i(Y = 1 | Z = z) = g_i(z)$ for all $z \in [0; 1]$.

It can be verified that $0 \leq g_i(u) \leq 1$ for all $u \in [0, 1]$. Since $\rho_i \leq 1$, for all $i \in [M]$, the ρ_i -ECE of the probability predictor f under P_i is lower bounded as

$$\begin{aligned} \rho_i\text{-ECE}_{P_i}(f) &= \int_0^1 g_i(u) \rho_i(u) du \\ &= 2 \int_{\frac{1}{4}}^{\frac{3}{4}} \frac{u(1-u)}{3} \sin(2i u) du \quad 0:1: \end{aligned}$$

Thus we know that $P_i \geq P_1^{\text{cont}}(\rho_i; p)$ for all $i \in [M]$. Now, observe that

$$L_n = \frac{1}{M} \prod_{i=1}^M \prod_{j=1}^n \frac{dP_i}{dP_0}(Z_j; Y_j) = \frac{1}{M} \prod_{i=1}^M \prod_{j=1}^n \frac{1 - Y_j + (2Y_j - 1)g_i(Z_j)}{1 - Y_j + (2Y_j - 1)Z_j}$$

and thus, for a random variable $(Z; Y) \sim P_0$, and defining $A_{a,b}$ below

$$\begin{aligned} E_{P_0}[L_n^2] &= \frac{1}{M^2} \prod_{a,b \in [M]} E_{P_0} \left[\prod_{j=1}^n \frac{1 - Y_j + (2Y_j - 1)g_a(Z_j)}{1 - Y_j + (2Y_j - 1)Z_j} \frac{1 - Y_j + (2Y_j - 1)g_b(Z_j)}{1 - Y_j + (2Y_j - 1)Z_j} \right] \\ &= \frac{1}{M^2} \prod_{a,b \in [M]} E_{P_0} \left[\frac{1 - Y + (2Y - 1)g_a(Z)}{1 - Y + (2Y - 1)Z} \frac{1 - Y + (2Y - 1)g_b(Z)}{1 - Y + (2Y - 1)Z} \right]^n \\ &=: \frac{1}{M^2} \prod_{a,b \in [M]} E_{P_0}[A_{a,b}]^n : \end{aligned}$$

In the second line, we have used the independence of the observations. If $a = b$, then

$$\begin{aligned} E_{P_0}[A_{a,b}] &= \int_0^1 u \frac{g_a(u)^2}{u^2} + (1-u) \frac{(1-g_a(u))^2}{(1-u)^2} du = 1 + \int_0^1 \frac{(g_a(u) - u)^2}{u(1-u)} du \\ &= 1 + \int_{\frac{1}{4}}^{\frac{3}{4}} \frac{1}{3} \sin^2(2a u) \frac{1}{4} du = \frac{13}{12} : \end{aligned}$$

If $a \neq b$, then

$$\begin{aligned} E_{P_0}[A_{a,b}] &= \int_0^1 u \frac{g_a(u)g_b(u)}{u^2} + (1-u) \frac{(1-g_a(u))(1-g_b(u))}{(1-u)^2} du \\ &= 1 + \int_0^1 \frac{(g_a(u) - u)(g_b(u) - u)}{u(1-u)} du \\ &= 1 + \int_{\frac{1}{4}}^{\frac{3}{4}} \frac{1}{3} \sin(2a u) \frac{1}{4} \sin(2b u) \frac{1}{4} du = 1 : \end{aligned}$$

Therefore,

$$E_{P_0}[L_n^2] = \frac{1}{M^2} \prod_{a,b \in [M]} \left(\frac{13}{12} \right)^n + (M^2 - M) : :$$

Choose a large enough $M \geq N_+$ such that $M \left[\left(\frac{13}{12} \right)^n - 1 \right] = (1 + \epsilon)^2$. Then,

$$\frac{1}{M^2} \prod_{a,b \in [M]} \left(\frac{13}{12} \right)^n + (M^2 - M) \leq 1 + (1 + \epsilon)^2 ;$$

and the result follows by Lemma 18. ■

A.8 Proof of Theorem 10

Overview of the proof. We construct m^{K-1} distributions under which the predictor f has an ϵ -ECE of $(n^{-2s/(4s+K-1)})$. The mis-calibration curves are constructed by linearly combining bump functions with disjoint supports. By properly scaling them, we can guarantee Hölder continuity. Also, the mis-calibration curves are chosen to be "almost" orthogonal in L^2 , so that the cross terms in the expansion of $\mathbb{E}_{P_0}[L_n^2]$ are small. We use Lemma 18 to conclude.

Proof The proof is inspired by the lower bound arguments in Arias-Castro et al. (2018). For $m := dn^{2/(4s+K-1)}$ and $2f - 1g^{[m]^{K-1}}$, we define alternative distributions $P \geq P_1(\cdot; p; s)$ with $\epsilon := c_{\text{lower}} n^{-2s/(4s+K-1)}$ and use Lemma 18 to prove $\mathbb{E}_n(p; s) \geq c_{\text{lower}} n^{-2s/(4s+K-1)}$. Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ be the function from (8). It can be verified that ψ is infinitely differentiable and its derivatives of every order are bounded. For $\mathbf{z} = (z_1, \dots, z_K) \geq 7! (z_1, \dots, z_{K-1}) \geq 7!$, we see that $[\frac{1}{2K}, \frac{1}{K}]^{K-1} \subset (\mathbf{z}_{K-1} \setminus [\frac{1}{2K}, 1]^K)$. For each $\mathbf{j} = (j_1, \dots, j_{K-1}) \geq 2 [m]^{K-1}$, define $\psi_{\mathbf{j}} : \mathbb{R}^{K-1} \rightarrow \mathbb{R}$ by

$$\psi_{\mathbf{j}}(x_1, \dots, x_{K-1}) := m^{-s} \prod_{k=1}^{K-1} \psi(2Kx_k - 1) \mathbf{1}_{j_k+1} : \quad (38)$$

Then, each $\psi_{\mathbf{j}}$ is supported on the cube

$$\text{supp}(\psi_{\mathbf{j}}) = \prod_{k=1}^{K-1} \left[\frac{j_k - 1 + m}{2Km}, \frac{j_k + m}{2Km} \right] :$$

The sets $\text{supp}(\psi_{\mathbf{j}})$ are disjoint for different indices $\mathbf{j} \geq 2 [m]^{K-1}$, and we have

$$\left| \bigcup_{\mathbf{j} \geq 2 [m]^{K-1}} \text{supp}(\psi_{\mathbf{j}}) \right| \leq \frac{1}{2K} \cdot \frac{1}{K}^{K-1} \cdot \left(\frac{1}{2K} \cdot 1 \right)^{K-1} :$$

Let $c_{\mathbf{j}} := \log(1 + (1 - \epsilon)^{2s/(4s+K-1)})$ and

$$:= \max_{t \geq 0, \dots, d_{\text{seg}}} \binom{t}{L-1}^{K-1} \frac{1}{2K} \wedge \frac{L(2K)^{d_{\text{seg}}}}{2^{K-1}} \wedge \frac{L(2K)^{d_{\text{seg}}+1}}{4} \wedge c_{\mathbf{j}} \frac{(2K)^{\frac{K-1}{2}}}{2^{K-1}} : \quad (39)$$

By the definition of $c_{\mathbf{j}}$ in (39), we see that

$$k \leq k_{L-1}^{K-1} \frac{1}{2K} ; \quad (40)$$

$$P \frac{1}{K-1} (2K)^{d_{\text{seg}}} \max_{t \geq 0, \dots, d_{\text{seg}}} \binom{t}{L-1}^{K-1} \frac{L}{2} ; \quad (41)$$

$$2(2K)^{dse-1} \max_{t \in \{0, \dots, dse\}} \binom{K-1}{L^1} \frac{L}{2}; \quad (42)$$

and

$$4^{2K} (2K)^{K+1} k k_{L^2}^{2(K-1)} c^2; \quad 1: \quad (43)$$

For each $2 \leq f \leq 1g^{[m]^{K-1}}$, define $g : \mathbb{R}^K \rightarrow \mathbb{R}^K$ by

$$g(z) := z + \sum_{j \in 2[m]^{K-1}} \binom{K-1}{j} (z)^A (1; 1; 0; \dots; 0)^j;$$

Then we have $\binom{K-1}{j} \leq k-1$ for all $2 \leq f \leq 1g^{[m]^{K-1}}$. This is because

$$[g(z)]_k = [z]_k + \sum_{j \in 2[m]^{K-1}} \binom{K-1}{j} (z)^A \sum_{j \in 2[m]^{K-1}} \binom{K-1}{j} (z)^A = 1$$

for all $z \in \mathbb{R}^K$ and

$$[g(z)]_k \geq [z]_k + \sum_{j \in 2[m]^{K-1}} \binom{K-1}{j} (z)^A \frac{1}{2^k} m^s k k_{L^1}^{K-1};$$

$$[g(z)]_k \geq [z]_k; \quad \text{if } z \in [\frac{1}{2^k}; 1]^K; k \geq 1; 2g;$$

$$0; \quad \text{otherwise.}$$

for all $z \in \mathbb{R}^K$ and $k \geq 1; \dots; K$ by (40).

Next, we claim that each coordinate function of the mapping $z \mapsto g(z) - z$ belongs to $H_K(s; L)$, i.e., for all multi-indices $2 \leq j \leq dse-1$ and $x_1, x_2 \in \mathbb{R}^K$,

$$\sum_{j \in 2[m]^{K-1}} \binom{K-1}{j} (x_1) \binom{K-1}{j} (x_2) \leq L k x_1 - x_2 k^{s d se+1}; \quad (44)$$

If $m k x_1 - x_2 k \leq 1$, then from the mean value theorem

$$\binom{K-1}{j} (x_1) \binom{K-1}{j} (x_2) \leq \frac{1}{K-1} \max_{\substack{2 \leq j \leq dse \\ j \in 2[m]^{K-1}}} \binom{K-1}{j} L^1 k x_1 - x_2 k; \quad (45)$$

By the definition of j in (38), for any $0 = (0; \dots; 0; 1) \in 2 \leq N^{K-1}$ with $j \in 2[m]^{K-1}$,

$$\binom{K-1}{j} L^1 = m^{j \cdot 0} s (2K)^{j \cdot 0} \sum_{k=1}^{K-1} \binom{K-1}{L^1} m^{dse-s} (2K)^{dse} \max_{t \in \{0, \dots, dse\}} \binom{K-1}{L^1} t; \quad (46)$$

Plugging (46) into (45), and using (41), we reach

$$\binom{K-1}{j} (x_1) \binom{K-1}{j} (x_2) \leq \frac{1}{K-1} m^{dse-s} (2K)^{dse} \max_{t \in \{0, \dots, dse\}} \binom{K-1}{L^1} t k x_1 - x_2 k$$

$$\leq \frac{L}{2} m^{dse-s} k x_1 - x_2 k \leq \frac{L}{2} k x_1 - x_2 k^{s d se+1};$$

If $m k x_1 \quad x_2 k > 1$, using (42), we similarly get

$$\begin{aligned} & \binom{(\cdot)}{j}(x_1) \quad \binom{(\cdot)}{j}(x_2) \quad 2 \max_{j \in \{1, \dots, d\}} \binom{(\cdot)}{j} \quad L^1 \\ & 2 m^{d s - 1} (2K)^{d s - 1} \max_{t \in \{0, \dots, d\}} \binom{(\cdot)}{t} \quad L^1 \quad \frac{L}{2} m^{d s - 1} \quad \frac{L}{2} k x_1 \quad x_2 k^{s d + 1} : \end{aligned}$$

Thus, for any $j \in [m]^{K-1}$ and $x_1, x_2 \in (\mathbb{R}^{K-1})$, it holds that

$$\binom{(\cdot)}{j}(x_1) \quad \binom{(\cdot)}{j}(x_2) \geq \frac{L}{2} k x_1 \quad x_2 k^{s d + 1} : \quad (47)$$

Given $x_1, x_2 \in (\mathbb{R}^{K-1})$, there can be two cases: (1) there exist $j_1 \in [m]^{K-1}$ such that

$$\sum_{j \in [m]^{K-1}} \binom{(\cdot)}{j}(x_1) \quad \binom{(\cdot)}{j}(x_2) \quad A = \sum_{j_1} \binom{(\cdot)}{j_1}(x_1) \quad \binom{(\cdot)}{j_1}(x_2) ;$$

or (2) there exist distinct $j_1, j_2 \in [m]^{K-1}$ such that $x_1 \in \text{supp}(\binom{(\cdot)}{j_1})$ and $x_2 \in \text{supp}(\binom{(\cdot)}{j_2})$. In the first case, (44) directly follows from (47). In the second case, choose a point x_3 on the line segment connecting x_1 and x_2 such that $x_3 \in \text{supp}(\binom{(\cdot)}{j_1}) \cap \text{supp}(\binom{(\cdot)}{j_2})$. Such a point exists since $\text{supp}(\binom{(\cdot)}{j_1}), \text{supp}(\binom{(\cdot)}{j_2})$ are open and $\text{supp}(\binom{(\cdot)}{j_1}) \cap \text{supp}(\binom{(\cdot)}{j_2}) \neq \emptyset$. For any $z \in \mathbb{R}^{K-1}$ with $j = j_1$, we have

$$\begin{aligned} & \sum_{j \in [m]^{K-1}} \binom{(\cdot)}{j}(x_1) \quad \binom{(\cdot)}{j}(x_2) = \sum_{j_1} \binom{(\cdot)}{j_1}(x_1) \quad \binom{(\cdot)}{j_2}(x_2) \\ & = \sum_{j_1} \binom{(\cdot)}{j_1}(x_1) \quad \binom{(\cdot)}{j_1}(x_3) + \sum_{j_2} \binom{(\cdot)}{j_2}(x_3) \quad \binom{(\cdot)}{j_2}(x_2) \\ & \quad \binom{(\cdot)}{j_1}(x_1) \quad \binom{(\cdot)}{j_1}(x_3) + \binom{(\cdot)}{j_2}(x_3) \quad \binom{(\cdot)}{j_2}(x_2) \\ & \geq \frac{L}{2} k x_1 \quad x_3 k^{s d + 1} + \frac{L}{2} k x_3 \quad x_2 k^{s d + 1} \geq \frac{L}{2} k x_1 \quad x_2 k^{s d + 1} : \end{aligned}$$

The second inequality holds because of (47), and the last inequality holds because $k x_3 \geq k x_1 \quad x_2 k$ and $s d + 1 > 0$. This finishes the proof of (44).

Now, let P_0 and P , $\mathbb{P}^{[m]^{K-1}}$, be the distributions of $(Z; Y) \in (\mathbb{R}^{K-1} \times \mathbb{Y})$ characterized by

$$Z \sim P_0 \text{ Unif}(\mathbb{R}^{K-1}); \quad Y | Z = z \sim P_0 \text{ Cat}(z) \text{ for all } z \in \mathbb{R}^{K-1};$$

and

$$Z \sim P \text{ Unif}(\mathbb{R}^{K-1}); \quad Y | Z = z \sim P \text{ Cat}(g(z)) \text{ for all } z \in \mathbb{R}^{K-1};$$

We have $P_0 \leq P$ by definition. For $Z \sim \text{Unif}(\mathbb{R}^{K-1})$ and $j_0 := 1_{\mathbb{R}^{K-1}} \in [m]^{K-1}$,

$$\begin{aligned} \mathbb{P}\text{-ECE}_P(f)^p &= \mathbb{E} \sum_{k=1}^K j[g(Z) - Z]_k^p = 2^p \mathbb{E} \sum_{j \in [m]^{K-1}} j(\binom{(\cdot)}{j})(Z) \\ &= 2(K-1)! p m^{K-1} \int_{\mathbb{R}^{K-1}} j \cdot j_0(x) j^p dx : \end{aligned}$$

Further,

$$\begin{aligned} m^{K-1} \int_{\mathbb{R}^{K-1}} \int_{\mathbb{R}} j_0(x) j^p dx &= m^{ps} \int_{\mathbb{R}} \int_{\mathbb{R}^{K-1}} (m(2Kx_k - 1)) j^p dx_k \\ &= m^{ps} (2K)^{K+1} k_{Lp}^{(K-1)p} \end{aligned} \quad (48)$$

Thus, we have

$$\hat{p}\text{-ECE}_P(f) = (2(2K)^{K+1} (K-1)!)^{\frac{1}{p}} m^{-s} k_{Lp}^{K-1} c_{\text{lower}} n^{-\frac{2s}{4s+K-1}} \quad (49)$$

for some $c_{\text{lower}} > 0$ because

$$\lim_{n \rightarrow \infty} \frac{m^s}{n^{\frac{2s}{4s+K-1}}} = \lim_{n \rightarrow \infty} \frac{1}{n^{\frac{2}{4s+K-1}}} e^{s \ln n^{\frac{2s}{4s+K-1}}} = 1:$$

From (44) and (49), we see that $P \geq P_1(\cdot; p; s)$ with $\epsilon = c_{\text{lower}} n^{-\frac{2s}{4s+K-1}}$ for all $\epsilon \leq \frac{1}{2} \frac{1}{m^{K-1}}$.

The final step is to apply Lemma 18. Given n i.i.d. observations $f(Z_i; Y_i) : i \in [n]$, the average likelihood ratio between P , $\frac{1}{2} \frac{1}{m^{K-1}}$, and P_0 is

$$L_n = \frac{1}{2^{m^{K-1}}} \prod_{i=1}^n \frac{[g_1(Z_i)]_{\arg\max_k [Y_i]_k}}{[Z_i]_{\arg\max_k [Y_i]_k}}$$

Let Z_1, Z_2 be independent random variables uniformly drawn from $\frac{1}{2} \frac{1}{m^{K-1}}$, and $(Z; Y) \sim P_0$. Then,

$$\begin{aligned} E_{P_0}[L_n^2] &= E_{Z_1; Z_2} E_{P_0} \prod_{i=1}^n \frac{[g_1(Z_i)]_{\arg\max_k [Y_i]_k}}{[Z_i]_{\arg\max_k [Y_i]_k}} \frac{[g_2(Z_i)]_{\arg\max_k [Y_i]_k}}{[Z_i]_{\arg\max_k [Y_i]_k}} \\ &= E_{Z_1; Z_2} E_{P_0} \frac{[g_1(Z)]_{\arg\max_k [Y]_k} [g_2(Z)]_{\arg\max_k [Y]_k}}{[Z]_{\arg\max_k [Y]_k}^2} \end{aligned}$$

Moreover,

$$\begin{aligned} E_{P_0} \frac{[g_1(Z)]_{\arg\max_k [Y]_k} [g_2(Z)]_{\arg\max_k [Y]_k}}{[Z]_{\arg\max_k [Y]_k}^2} &= \int_{\mathbb{R}^{K-1}} \int_{\mathbb{R}} \frac{[g_1(z)]_k [g_2(z)]_k}{[z]_k^2} dz \\ &= 1 + \int_{\mathbb{R}^{K-1}} \int_{\mathbb{R}} \frac{[g_1(z) - z]_k [g_2(z) - z]_k}{[z]_k} dz \\ &= 1 + \int_{\mathbb{R}^{K-1}} \int_{\mathbb{R}} \frac{1}{j} \left(\frac{1}{j} \right) (z) A \frac{1}{j} \left(\frac{1}{j} \right) (z) A \frac{1}{[z]_0} + \frac{1}{[z]_1} dz \end{aligned}$$

where the integral $\int_{\mathbb{R}^{K-1}}$ is over $\text{Unif}(\mathbb{R}^{K-1})$. Since $\{ \text{supp}(j) : j \in [m]^{K-1} \}$ is a collection of pairwise disjoint sets,

$$\begin{aligned} & \int_{\mathbb{R}^{K-1}} \sum_{j \in [m]^{K-1}} \mathbb{1}_{\text{supp}(j)}(z) \left(\frac{1}{[z]_0} + \frac{1}{[z]_1} \right) dz \\ &= \int_{\mathbb{R}^{K-1}} \sum_{j \in [m]^{K-1}} \mathbb{1}_{\text{supp}(j)}(z) \left(\frac{1}{[z]_0} + \frac{1}{[z]_1} \right) dz: \end{aligned}$$

The random variable $(\frac{1}{[z]_0}, \frac{1}{[z]_1})_{j \in [m]^{K-1}}$ is also uniformly distributed on $[0, 1]^{[m]^{K-1}}$. Hence,

$$E_{P_0}[L_n^2] = E \left[\sum_{j \in [m]^{K-1}} \int_{\mathbb{R}^{K-1}} \mathbb{1}_{\text{supp}(j)}(z) \left(\frac{1}{[z]_0} + \frac{1}{[z]_1} \right) dz \right]^2 : \quad (50)$$

Since $\text{supp}(j) \subset (\frac{1}{2K}, 1]^K$, we have $\frac{1}{[z]_0} + \frac{1}{[z]_1} \leq 4K$ for every $z \in \mathbb{R}^{K-1}$ such that $\mathbb{1}_{\text{supp}(j)}(z) > 0$ for at least one $j \in [m]^{K-1}$. Thus,

$$\begin{aligned} & \sum_{j \in [m]^{K-1}} \int_{\mathbb{R}^{K-1}} \mathbb{1}_{\text{supp}(j)}(z) \left(\frac{1}{[z]_0} + \frac{1}{[z]_1} \right) dz \\ & \leq \sum_{j \in [m]^{K-1}} \int_{\mathbb{R}^{K-1}} \mathbb{1}_{\text{supp}(j)}(z) 4K dz \\ & = 4K \sum_{j \in [m]^{K-1}} \int_{\mathbb{R}^{K-1}} \mathbb{1}_{\text{supp}(j)}(z) dz = 4 \cdot 2^K m^{K-1} \int_{\mathbb{R}^{K-1}} j_0(x)^2 dx: \quad (51) \end{aligned}$$

The last equality is because $\text{Unif}(\mathbb{R}^{K-1})$ has density $(K-1)!$ with respect to $\text{Leb}_{\mathbb{R}^{K-1}}$, when projected to \mathbb{R}^{K-1} . Also, by (43) and (48),

$$4 \cdot 2^K m^{K-1} \int_{\mathbb{R}^{K-1}} j_0(x)^2 dx = 4 \cdot 2^K m^{2s} (2K)^{K+1} k_{L^2}^{2(K-1)} c_{\cdot}^2 m^{2s-1}:$$

By (50) and that $(1+x)^n = \exp(nx)$ for all $x \in (-1, 1]$,

$$E_{P_0}[L_n^2] = E \left[\sum_{j \in [m]^{K-1}} \int_{\mathbb{R}^{K-1}} \mathbb{1}_{\text{supp}(j)}(z) \left(\frac{1}{[z]_0} + \frac{1}{[z]_1} \right) dz \right]^2 : \quad (52)$$

Since $\{j_1 : j_2 \in [m]^{K-1}\}$ is a set of i.i.d. random variables drawn from $\text{Unif}(\mathcal{J})$, we have, with $\cosh(x) := [\exp(x) + \exp(-x)]/2$,

$$\begin{aligned} & E \left[\exp \left(\sum_{j_2 \in [m]^{K-1}} X_{j_2} \right) \right] \\ &= E \left[\prod_{j_2 \in [m]^{K-1}} \exp(X_{j_2}) \right] \\ &= \prod_{j_2 \in [m]^{K-1}} E \left[\exp(X_{j_2}) \right] \\ &= \prod_{j_2 \in [m]^{K-1}} \int_{\mathcal{Z}} \cosh \left(\sum_{k=1}^K (z_k)^2 \frac{1}{[z]_0} + \frac{1}{[z]_1} \right) dz \end{aligned}$$

Similarly to (51) and (52), we have

$$\int_{\mathcal{Z}} \sum_{k=1}^K (z_k)^2 \frac{1}{[z]_0} + \frac{1}{[z]_1} dz \leq c_{\mathcal{Z}}^2 m^{2s} K^{+1} n^{-1} \quad (53)$$

for each $j_2 \in [m]^{K-1}$. Using that $\cosh(x) \leq 1 + x^2 e^{x^2}$ for $x \in [-1, 1]$,

$$\begin{aligned} & E \left[\exp \left(\sum_{j_2 \in [m]^{K-1}} X_{j_2} \right) \right] \\ & \leq \prod_{j_2 \in [m]^{K-1}} \int_{\mathcal{Z}} \cosh \left(\sum_{k=1}^K (z_k)^2 \frac{1}{[z]_0} + \frac{1}{[z]_1} \right) dz \\ & \leq \prod_{j_2 \in [m]^{K-1}} \int_{\mathcal{Z}} \exp \left(\sum_{k=1}^K (z_k)^2 \frac{1}{[z]_0} + \frac{1}{[z]_1} \right) dz \leq A \end{aligned}$$

Again from (53), it follows that

$$\begin{aligned} & E \left[\exp \left(\sum_{j_2 \in [m]^{K-1}} X_{j_2} \right) \right] \\ & \leq \exp \left(c_{\mathcal{Z}}^4 m^{4s} K^{+1} n^2 \right) \leq 1 + (1 + c_{\mathcal{Z}}^4 m^{4s} K^{+1} n^2)^2 \end{aligned}$$

In conclusion, $\mathbb{P}_n(p; s) \geq \mathbb{P}_{\text{lower}}(n^{2s=(4s+K-1)})$ by Lemma 18. ■

A.9 Proof of Theorem 12

Overview of the proof. Under $\mathbb{P} \in \mathcal{P}_0$, we prove that $T_{1;k}$ and $T_{2;k}$ have zero mean, and their variances are bounded by unity. By rejecting H_0 when $|jT_{1;k}| \geq \frac{3K=n}{3}$ or $|jT_{2;k}| \geq \frac{3K=n}{3}$, we can filter out distributions $\mathbb{P} \in \mathcal{P}_1(p; s)$ such that $E_{\mathbb{P}}[T_{1;k}] = (n^{-1/2})$ or $E_{\mathbb{P}}[T_{2;k}] = (n^{-1/2})$. For the remaining cases, we compute $\mathbb{V}_k^{\mathbb{P}} \mathbb{W}_k^{\mathbb{P}} \mathbb{K}_{L^2(\mathcal{P}_{\mathcal{Z}})}$ and show it is lower bounded by $(n^{2s=(4s+K-1)})$. We conclude using the minimax optimality of the two-sample test TS

Proof We prove the theorem for $p = 2$. Then, the general case follows since $P_1(\cdot; p; s) \leq P_1(\cdot; 2; s)$ for all $p \geq 2$. Assume $P \in \mathcal{P}_0$. By the union bound,

$$P\left(\frac{\text{split}}{n} = 1\right) \leq \sum_{k=1}^K P\left(|T_{1;k}| \geq \frac{r}{3K} \frac{!}{n}\right) + P\left(|T_{2;k}| \geq \frac{r}{3K} \frac{!}{n}\right) + P\left(\text{TS}_{\frac{3K}{2}}(V_k; W_k) = 1\right) \quad \#$$

Moreover, for all $k \in \{1, \dots, K\}$,

$$E_P[Y - Z]_k = E_P[E_P[[Y - Z]_{k|Z}]] = E_P[[E_P[Y|Z] - Z]_k] = 0$$

and $\text{Var}_P([Y - Z]_k) = E_P[[Y - Z]_k^2] \leq 1$: Thus, by Chebyshev's inequality

$$P\left(|T_{1;k}| \geq \frac{r}{3K} \frac{!}{n}\right) \leq \frac{n}{3K} \text{Var}_P(T_{1;k}) = \frac{n}{3K} \text{Var}_P([Y - Z]_k) \leq \frac{1}{3K}.$$

Similarly, we have

$$P\left(|T_{2;k}| \geq \frac{r}{3K} \frac{!}{n}\right) \leq \frac{n}{3K} \text{Var}_P(T_{2;k}) = \frac{n}{3K} \text{Var}_P([Z]_k [Y - Z]_k) \leq \frac{1}{3K}.$$

From (13), we know that $P(\text{TS}_{\frac{3K}{2}}(V_k; W_k) = 1) \leq \frac{1}{3K}$: Therefore,

$$P\left(\frac{\text{split}}{n} = 1\right) \leq \sum_{k=1}^K \frac{1}{3K} + \frac{1}{3K} + \frac{1}{3K} = \frac{3}{3K} = \frac{1}{K}.$$

Let $P \in \mathcal{P}_1(\cdot; p; s)$ and suppose that, for some $k \in \{1, \dots, K\}$,

$$|E_P[T_{1;k}]| \geq |E_P[Y - Z]_k| \geq \frac{1}{n} \left(\frac{r}{3K} + \frac{1}{p} \right) \quad (54)$$

By Chebyshev's inequality,

$$P\left(|T_{1;k} - E_P[T_{1;k}]| \geq \frac{1}{n} \left(\frac{r}{3K} + \frac{1}{p} \right)\right) \leq \frac{n \text{Var}_P(T_{1;k})}{\left(\frac{r}{3K} + \frac{1}{p} \right)^2} \leq \frac{1}{\left(\frac{r}{3K} + \frac{1}{p} \right)^2}.$$

Note that (54) and $|T_{1;k} - E_P[T_{1;k}]| \leq \frac{1}{n}$ imply

$$|T_{1;k} - E_P[T_{1;k}]| \geq |T_{1;k} - E_P[T_{1;k}]| \geq \frac{r}{3K}.$$

Therefore,

$$P\left(\frac{\text{split}}{n} = 1\right) \leq P\left(|T_{1;k}| \geq \frac{r}{3K}\right) \leq P\left(|T_{1;k} - E_P[T_{1;k}]| \geq \frac{r}{3K}\right) \leq \frac{1}{\left(\frac{r}{3K} + \frac{1}{p} \right)^2} \leq \frac{1}{\left(\frac{r}{3K} \right)^2}.$$

The same conclusion can be drawn when

$$|E_P[T_{2;k}]| \geq |E_P[[Z]_k [Y - Z]_k]| \geq \frac{1}{n} \left(\frac{r}{3K} + \frac{1}{p} \right)$$

for some $k_0 \in \{1, \dots, K\}$.

Now it remains to prove the claim for $P \in \mathcal{H}(\cdot; p; s)$ such that

$$|E_P[Y - Z]_{k_0}| \leq \frac{1}{n} \left(\frac{3K}{n} + p \right)^{\frac{1}{2}} \quad (55)$$

for every $k \in \{1, \dots, K\}$. Since

$$\|E_P(f)\|^2 = \sum_{k=1}^K \int [re_{\mathfrak{F}}(z)]_k^2 dP_Z(z) \quad (56)$$

we can choose $k_0 \in \{1, \dots, K\}$ such that $\int [re_{\mathfrak{F}}(z)]_{k_0}^2 dP_Z(z) \geq \frac{1}{K}$. Choose $c_{\text{split}}^0 > 0$ such that, for d_c from Assumption 2 and c_{ts} from (13),

$$\frac{(c_{\text{split}}^0)^2}{K} \geq \frac{4}{d_c^3} \left(\frac{3K}{n} + p \right)^{\frac{1}{2}} + c_{\text{ts}}^2 \frac{d_c}{8} \frac{4s}{4s+K-1} \quad (56)$$

There exists $N \in \mathbb{N}_+$ such that for all $n \geq N$,

$$\frac{1}{n} \left(\frac{3K}{n} + p \right)^{\frac{1}{2}} \leq \frac{d_c}{2}, \quad \frac{2}{e} \frac{d_c n}{8} \geq \frac{1}{2} \quad (57)$$

Let $c_{\text{split}} = c_{\text{split}}^0 - N^{2s-(4s+K-1)}$: If $n < N$, then $\mathcal{P}_1(\cdot; p; s)$ is empty since $c_{\text{split}} n^{2s-(4s+K-1)} > 1$, so the claim is vacuously true. Assume $n \geq N$. By (55), (57), and Assumption 2,

$$E_P[Z]_{k_0} - E_P[Y]_{k_0} \leq |E_P[Y - Z]_{k_0}| \leq \frac{d_c}{2} \quad (58)$$

By (11), (12), (55), and (58), $k_{k_0}^V \int_{k_0}^W k_{k_0}^2 dL^2(P_Z)$ is lower bounded by

$$\frac{1}{(E_P[Y]_{k_0})^2} \int [re_{\mathfrak{F}}(z)]_{k_0}^2 dP_Z(z) + \frac{2E_P[Z - Y]_{k_0} E_P[[Z]_{k_0}[Y - Z]_{k_0}]}{(E_P[Y]_{k_0})^2 E_P[Z]_{k_0}} \geq \frac{1}{K} \frac{4}{d_c^3} \left(\frac{3K}{n} + p \right)^{\frac{1}{2}} \geq \frac{1}{n} \quad (59)$$

Further by $c_{\text{split}}^0 n^{2s-(4s+K-1)}$ and (56),

$$\frac{1}{K} \frac{4}{d_c^3} \left(\frac{3K}{n} + p \right)^{\frac{1}{2}} \geq \frac{1}{n} \geq \frac{(c_{\text{split}}^0)^2}{K} \frac{4}{d_c^3} \left(\frac{3K}{n} + p \right)^{\frac{1}{2}} \frac{1}{n^{4s+K-1}} + c_{\text{ts}}^2 \frac{d_c n}{8} \frac{4s}{4s+K-1} \quad (59)$$

In conclusion,

$$k_{k_0}^V \int_{k_0}^W k_{k_0}^2 dL^2(P_Z) \geq c_{\text{ts}} \frac{d_c n}{8} \frac{4s}{4s+K-1} \quad (59)$$

Note that $\frac{V_{k_0}}{3K} + \frac{W_{k_0}}{2}$ is s -Holder since it is a linear combination of two s -Holder functions $z \mapsto [res(z)]_{k_0}$ and $z \mapsto [z]_{k_0}$, possibly with different Holder constants. Thus by (13), we have

$$P \left(TS_{\frac{1}{3K}, \frac{1}{2}}(V_{k_0}; W_{k_0}) = 1 \mid |V_{k_0}| = v; |W_{k_0}| = w \right) \leq \frac{1}{2} \quad (60)$$

given that $v; w \leq \frac{d_c n}{8}$. For convenience, assume that is even (if required, drop an observation). Since $|V_{k_0}| \sim \text{Bin}(\frac{n}{2}; E_P[Y]_{k_0})$ and $|W_{k_0}| \sim \text{Bin}(\frac{n}{2}; E_P[Z]_{k_0})$, we find

$$P \left(|V_{k_0}| < \frac{n E_P[Y]_{k_0}}{4} \right) \leq \frac{2}{e} \frac{e^{-\frac{n E_P[Y]_{k_0}}{4}}}{4}; \quad P \left(|W_{k_0}| < \frac{n E_P[Z]_{k_0}}{4} \right) \leq \frac{2}{e} \frac{e^{-\frac{n E_P[Z]_{k_0}}{4}}}{4}$$

by Chernoff's inequality (Exercise 2.3.2 of Vershynin (2018)). Therefore,

$$\begin{aligned} P \left(|V_{k_0}| < \frac{d_c n}{8} \text{ or } |W_{k_0}| < \frac{d_c n}{8} \right) &= P \left(|V_{k_0}| < \frac{d_c n}{8} \right) + P \left(|W_{k_0}| < \frac{d_c n}{8} \right) \\ &= P \left(|V_{k_0}| < \frac{n E_P[Y]_{k_0}}{4} \right) + P \left(|W_{k_0}| < \frac{n E_P[Z]_{k_0}}{4} \right) \\ &\leq \frac{2}{e} \frac{e^{-\frac{n E_P[Y]_{k_0}}{4}}}{4} + \frac{2}{e} \frac{e^{-\frac{n E_P[Z]_{k_0}}{4}}}{4} \leq 2 \frac{2}{e} \frac{e^{-\frac{d_c n}{8}}}{2}; \end{aligned}$$

and thus

$$P \left(|V_{k_0}|; |W_{k_0}| \geq \frac{d_c n}{8} \right) \geq \frac{1}{2}; \quad (61)$$

The last inequality holds by (56). Finally by (60) and (61),

$$\begin{aligned} P \left(TS_{\frac{1}{3K}, \frac{1}{2}}(V_{k_0}; W_{k_0}) = 1 \right) &= \int_{\mathbb{R}^2} P \left(TS_{\frac{1}{3K}, \frac{1}{2}}(V_{k_0}; W_{k_0}) = 1 \mid |V_{k_0}| = v; |W_{k_0}| = w \right) dP_{|V_{k_0}|; |W_{k_0}|}(v; w) \\ &\geq \int_{\mathbb{R}^2} \frac{1}{2} \mathbb{1}_{\left(|V_{k_0}| \geq \frac{d_c n}{8}; |W_{k_0}| \geq \frac{d_c n}{8} \right)} dP_{|V_{k_0}|; |W_{k_0}|}(v; w) \end{aligned}$$

where the integral is with respect to the joint distribution of $|V_{k_0}|$ and $|W_{k_0}|$. This proves that

$$P \left(\frac{\text{split}}{n} = 1 \right) \geq P \left(TS_{\frac{1}{3K}, \frac{1}{2}}(V_{k_0}; W_{k_0}) = 1 \right) \geq \frac{1}{2}; \quad \blacksquare$$

A.10 Proof of Corollary 13

The proof of Theorem 12 can be repeated by replacing (13) with (14). The only difference is (59), where we used $n^{4s-(4s+K-1)} = n^{-1}$. In the adaptive setting, we instead need $(n = \log \log n)^{4s-(4s+K-1)} = n^{-1}$. This inequality holds for all large $n \in \mathbb{N}_+$, say $n \geq N^0$. Then, we can define c_{ad-s} to be larger than $(N^0 = \log \log N^0)^{2s-(4s+K-1)}$, so that $P_1(\cdot; p; s)$ becomes empty when $\epsilon < N^0$.

Appendix B. Background

B.1 Hölder Continuity on the Probability Simplex

To define derivatives (and thus the class of functions we study) on Σ_{K-1} , a coordinate chart $\phi: \Sigma_{K-1} \rightarrow \mathbb{R}^{K-1}$ has to be specified. For example, we can consider the canonical projection $\phi_K: (z_1, \dots, z_K) \mapsto (z_1, \dots, z_{K-1}, z_{K+1}, \dots, z_K)$. The definition of Hölder smoothness below depends on the choice of ϕ . We assume $\phi = \phi_K$, but all conclusions and proofs remain the same for any choice of the coordinate chart.

For an integer $d \geq 1$, a vector $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}^d$ is called a multi-index. We write $|\alpha| := \alpha_1 + \dots + \alpha_d$. For a sufficiently smooth function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, we denote its partial derivative of order $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}^d$ by $f^{(\alpha)} := \partial^{\alpha_1} \dots \partial^{\alpha_d} f$. For a Hölder smoothness parameter $s > 0$ and a Hölder constant $L > 0$, let $H_K(s; L)$ be the class of $(s; L)$ -Hölder continuous functions $g: \Sigma_{K-1} \rightarrow \mathbb{R}$ satisfying, for all $x_1, x_2 \in \Sigma_{K-1}$ and multi-indices $\alpha \in \mathbb{N}^d$ with $|\alpha| = d - 1$,

$$|g^{(\alpha)}(x_1) - g^{(\alpha)}(x_2)| \leq L \|x_1 - x_2\|^{s+1} \quad (62)$$

In particular, $H_K(1; L)$ denotes all L -Lipschitz functions.]

B.2 Two-sample Goodness-of-fit Tests

Here we state and slightly extend the results of Arias-Castro et al. (2018); Kim et al. (2022). For $d \in \mathbb{N}_+$, let μ be a measure on $[0, 1]^d$ which is absolutely continuous with respect to Leb_d and satisfies

$$\left| \frac{d\mu}{d\text{Leb}_d} - u \right| \leq \epsilon \quad (63)$$

almost everywhere for some constants $\epsilon, u > 0$. For $n_1, n_2 \in \mathbb{N}_+$, suppose we have two samples $V_1, \dots, V_{n_1} \stackrel{i.i.d.}{\sim} \mu$ and $W_1, \dots, W_{n_2} \stackrel{i.i.d.}{\sim} \mu$ sampled from the distributions on $[0, 1]^d$ with densities f_1 and f_2 , respectively, with respect to μ . We also assume $f_1 - f_2$ is $(s; L)$ -Hölder continuous for a Hölder smoothness parameter $s > 0$ and a Hölder constant $L > 0$.

For $m \in \mathbb{N}_+$ and $i = (i_1, \dots, i_d) \in [m]^d$, let $R_{m;i} := \prod_{k=1}^d \frac{i_k - 1}{m}$.

$$v_{i;m} := \frac{1}{n_1} \sum_{j=1}^{n_1} \mathbb{1}_{V_j \in R_{m;i}};$$

$$w_{i;m} := \frac{1}{n_2} \sum_{j=1}^{n_2} \mathbb{1}_{W_j \in R_{m;i}};$$

The unnormalized chi-squared statistic is defined by

$$\chi_{m;n_1;n_2}^2 := \sum_{i \in [m]^d} (n_2 v_{i;m} - n_1 w_{i;m})^2;$$

Theorem 19 (Chi-squared test, Arias-Castro et al. (2018)) Consider the two-sample goodness-of-fit testing problem described above. Assume the Hölder smoothness parameter s is known and let $m = b(n_1 \wedge n_2)^{2/(4s+d)}$. For any $\epsilon \in (0, 1)$ and $\delta \in (0, 1)$, there

exist $c > 0$ depending on $(d; L)$ and $c_{ts} > 0$ depending on $(s; d; L; \rho; \mu; \gamma)$ such that for $n := n_1 n_2 (n_1 + n_2) + c n_1 n_2 m^{d-2}$,

$$\begin{aligned} P(U_{m; n_1; n_2} \leq c_{ts}) &= 1 && \text{if } f_1 = f_2; \\ P(U_{m; n_1; n_2} \leq c_{ts}) &< 1 && \text{if } k f_1 - f_2 k_{L^2(\cdot)} \geq c_{ts} (n_1 \wedge n_2)^{\frac{2s}{4s+d}}; \end{aligned}$$

For $m \in \mathbb{N}_+$, let $k_m : ([0; 1]^d)^4 \rightarrow \mathbb{R}$ be the kernel

$$k_m(v_1; v_2; w_1; w_2) = \frac{1}{m} \sum_{i \in [m]^d} 1_{R_{m,i}}(v_1) 1_{R_{m,i}}(v_2) + 1_{R_{m,i}}(w_1) 1_{R_{m,i}}(w_2) - 1_{R_{m,i}}(v_1) 1_{R_{m,i}}(w_2) - 1_{R_{m,i}}(w_1) 1_{R_{m,i}}(v_2);$$

For the two samples $\{V_1, \dots, V_{n_1}\}$, $\{W_1, \dots, W_{n_2}\}$ described above, define

$$U_{m; n_1; n_2} := \frac{1}{n_1(n_1 - 1)n_2(n_2 - 1)} \sum_{i_1 \in [n_1], i_2 \in [n_2]} \sum_{j_1 \in [n_1], j_2 \in [n_2]} k_m(V_{i_1}; V_{i_2}; W_{j_1}; W_{j_2});$$

For any $\alpha \in (0; 1)$, the $1 - \alpha$ quantile $c_{1-\alpha; m; n_1; n_2}$ of the U-statistic $U_{m; n_1; n_2}$ can be found by the permutation procedure described in Section 2.1 of Kim et al. (2022).

Theorem 20 (Permutation test, Kim et al. (2022)) Consider the two-sample goodness-of-fit testing described above. Assume the Hölder smoothness parameter s is known and let $m = \lfloor (n_1 \wedge n_2)^{\frac{2s}{4s+d}} \rfloor$. For any $\alpha \in (0; 1)$ and $\beta \in (0; 1)$, there exists $c_{ts} > 0$ depending on $(s; d; L; \rho; \mu; \gamma)$ such that

$$\begin{aligned} P(U_{m; n_1; n_2} \leq c_{1-\alpha; m; n_1; n_2}) &= 1 - \alpha && \text{if } f_1 = f_2; \\ P(U_{m; n_1; n_2} \leq c_{1-\alpha; m; n_1; n_2}) &< 1 - \alpha && \text{if } k f_1 - f_2 k_{L^2(\cdot)} \geq c_{ts} (n_1 \wedge n_2)^{\frac{2s}{4s+d}}; \end{aligned}$$

Corollary 21 (Multi-scale permutation test, Kim et al. (2022)) Consider the two-sample goodness-of-fit testing problem described above. Let $\beta = d_d^2 \log_2(\frac{n_1 \wedge n_2}{\log \log(n_1 \wedge n_2)})$ and define

$$c_{\beta}^{\text{perm}} := \max_{b \in \{1, \dots, B\}} P(U_{2^b; n_1; n_2} \leq c_{1-\beta; 2^b; n_1; n_2});$$

For any $\alpha \in (0; 1)$ and $\beta \in (0; 1)$, there exists $c_{\text{ad}} > 0$ depending on $(d; L; \rho; \mu; \gamma)$ such that

$$\begin{aligned} P(c_{\beta}^{\text{perm}} = 1) &= 1 && \text{if } f_1 = f_2; \\ P(c_{\beta}^{\text{perm}} = 1) &< 1 && \text{if } k f_1 - f_2 k_{L^2(\cdot)} \geq c_{\text{ad}} \frac{n_1 \wedge n_2}{\log \log(n_1 \wedge n_2)}^{\frac{2s}{4s+d}}; \end{aligned}$$

Remark 22 (Comment on the proofs) Theorem 19, Theorem 20, and Corollary 21 are generalization of Theorem 4 of Arias-Castro et al. (2018), Proposition 4.6 of Kim et al. (2022), and Proposition 7.1 of Kim et al. (2022), respectively. The original statements are for $\mu = \text{Leb}_d$. The proofs in Arias-Castro et al. (2018); Kim et al. (2022) can be adapted with only minor differences. For example, equation (93) of Arias-Castro et al. (2018) is still true assuming (63) above. Also, we proved Lemma 15 in this paper, which is a generalization of Lemma 3 of Arias-Castro et al. (2018) to a general measure μ . Other parts of the proofs in Arias-Castro et al. (2018); Kim et al. (2022) can be repeated without any modification.

B.3 Binning Scheme for the Probability Simplex

Here we describe a binning scheme for the probability simplex $\Delta_{K-1} \subset \mathbb{R}^K$ into equal volume simplices, i.e., a refinement of the standard probability simplex.

Hypersimplex. For $u \geq 2$ and $v \in [u-1]$, define the $(u-1)$ -dimensional polytope hypersimplex $\Delta_{u-1,v}$ a generalization of the standard probability simplex that can have more vertices and edges as

$$\Delta_{u-1,v} := \{(x_1, \dots, x_u) \in [0, 1]^u : x_1 + \dots + x_u = v\}$$

Let $A_{u-1,v}$ be the Eulerian number:

$$A_{u-1,v} := \sum_{i=0}^{X^v} \binom{u-1}{i} (v-i)^{u-1}$$

It is known that the hypersimplex $\Delta_{u-1,v}$ can be partitioned into $A_{u-1,v}$ simplices (Stanley, 1977; Sturmfels, 1996) whose volumes are identical to that of the unit probability simplex Δ_{u-1} .

Construction. Let $m \in \mathbb{N}_+$ and $R_{m;i} := \prod_{k=1}^K [i_k - \frac{1}{m}, i_k + \frac{1}{m}]$ for $i = (i_1, \dots, i_K) \in [m]^K$. The hypercube $R_{m;i}$ has a positive-volume intersection with Δ_{K-1} when

$$m+1 \leq \sum_{k=1}^K i_k \leq m+K-1 \tag{64}$$

Suppose that (64) holds and we write $j = \prod_{k=1}^K i_k$. Then, the intersection is

$$\begin{aligned} \Delta_{K-1} \cap R_{m;i} &= \{(x_1, \dots, x_K) \in \Delta_{K-1} : x_1 + \dots + x_K = 1\} \\ &= \left\{ \frac{i_1 - 1}{m} + \frac{1}{m} \sum_{k=1}^K x_k \in [0, 1] : x_1 + \dots + x_K = m+K-j \right\} \end{aligned}$$

which is a $\frac{1}{m}$ -scaled and translated version of the hypersimplex $\Delta_{K-1, m+K-j}$. Recall that this hypersimplex can be further partitioned into $A_{K-1, m+K-j}$ simplices with the volume $\text{vol}(\Delta_{K-1}) = m^{K-1}$. Let S_i be the set of such $A_{K-1, m+K-j}$ simplices. Now, we define

$$B_m := \left[\frac{i_1 - 1}{m} + \frac{1}{m} : 2 \right] S_i$$

$i \in [m]^K$
 $m+1 \leq \sum_{j=1}^K i_j \leq m+K-1$

which is the collection of all simplices obtained from decomposing $\Delta_{K-1} \cap R_{m;i}$ for each $i \in [m]^K$ satisfying (64). Since each simplex in B_m has a volume $\text{vol}(\Delta_{K-1}) = m^{K-1}$, it follows that $\text{vol}(B_m) = m^K$. Noting that there are $\binom{j}{K-1}$ multi-indices $i \in [m]^K$ with $\sum i_j = j$, it also directly follows from Worpitzky's identity (Equation 6.37 in Graham et al. (1994))

that

$$\begin{aligned}
 jB_m &= \sum_{j=m+1}^{m+K-1} \binom{j-1}{K-1} A_{K-1; m+K-j-1} = \sum_{j=0}^{K-2} \binom{m+j}{K-1} A_{K-1; K-j-2} \\
 &= \sum_{j=0}^{K-2} \binom{m+j}{K-1} A_{K-1; j} = m^{K-1}.
 \end{aligned}$$

We naturally index the partition as $B_m = \{B_1, \dots, B_{m+K-1}\}$:

B.4 Calibration Test for Discrete Predictions

Testing calibration of discrete probability predictions has been studied in Cox (1958); Miller (1962); Harrell (2015). Here we describe a test based on testing multiple binomial parameters. This test is related to the chi-squared test from Miller (1962), but does not use the asymptotic distribution of the test statistic when choosing critical values. Let $f: \mathcal{X} \rightarrow [0, 1]$ be the range of a discrete-valued probability predictor f . For each $i \in [t]$, we let $N_i = |\{j \in [n] : f(X_j) = v_i\}|$, $M_i = |\{j \in [n] : f(X_j) = v_i; Y_j = 1\}|$, and $p_i = P(Y = 1 | f(X) = v_i)$. In this setting, the random variable M_i , under the null hypothesis of perfect calibration, follows the binomial distribution $\text{Binom}(N_i; p_i)$ with N_i trials and success probability p_i , given N_i . We use an exact binomial test to test the null hypothesis $H_{0,i} : p_i = v_i$ for each $i \in [t]$ and apply the Bonferroni correction to control the false detection rate under the null hypothesis $H_0 = \bigcap_{i=1}^t H_{0,i}$.

Appendix C. More Experiments

C.1 Debiased ECE from Kumar et al. (2019)

In this subsection, we provide additional experiments, to evaluate calibration methods whose outputs range over a finite set. When the predicted probabilities belong to a finite set, Kumar et al. (2019) propose a debiased version of the empirical₂-ECE, whose sample complexity required is smaller than that of the plug-in estimator; see Section 3.3 for more discussion. We calculate Kumar et al. (2019)'s debiased estimator along with calibration testing results. We use models trained on CIFAR-10 (Table 4), CIFAR-100 (Table 5), and ImageNet (Table 6). The values of the debiased₂-ECE estimator are typically very small, which is consistent with the experimental results in Kumar et al. (2019). As can be observed, there is no clear relation between the debiased empirical₂-ECE values and test results.

	DenseNet 121		ResNet 50		VGG-19	
	$\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$	Calibrated?	$\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$	Calibrated?	$\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$	Calibrated?
Hist. Bin.	0.02%	reject	0.02%	reject	0.05%	reject
Scal. Bin.	0.11%	reject	0.10%	reject	0.20%	reject

Table 4: The values of the debiased empirical $\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$ (Kumar et al., 2019) and the testing results, via multiple binomial testing, of models trained on CIFAR-10 with two discrete calibration methods.

	MobileNet-v2		ResNet 56		ShuffleNet-v2	
	$\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$	Calibrated?	$\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$	Calibrated?	$\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$	Calibrated?
Hist. Bin.	0.04%	reject	0.09%	reject	0.15%	reject
Scal. Bin.	0.04%	reject	0.03%	reject	0.02%	accept

Table 5: The values of the debiased empirical $\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$ (Kumar et al., 2019) and the testing results, via multiple binomial testing, of models trained on CIFAR-100 with two discrete calibration methods.

	DenseNet 161		ResNet 152		EfficientNet-b7	
	$\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$	Calibrated?	$\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$	Calibrated?	$\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$	Calibrated?
Hist. Bin.	0.01%	reject	0.03%	reject	0.02%	reject
Scal. Bin.	0.05%	reject	0.03%	reject	0.06%	reject

Table 6: The values of the debiased empirical $\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$ (Kumar et al., 2019) and the testing results, via multiple binomial testing, of models trained on ImageNet with two discrete calibration methods.

C.2 Estimation of ECE via DPE

The DPE $T_{m;n}^d$ is used as a test statistic in the main text, but it can also be interpreted as an estimate of $\hat{\epsilon}_2^{\text{ECE}^{\text{db}}}$. While it builds on ideas from nonparametric functional estimation (Remark 2), the convergence rate $n^{-\frac{4s}{4s+K-1}}$ does not directly follow here, since our construction involves an additional estimation step $P_Z(B_i)$ for $m_i, j=n$. If we use the same binning parameter $m = m_n^{\frac{2}{4s+K-1}}$, then the ratio between the standard deviation and the mean of the estimate $m_i, j=n$ is

$$\frac{\sqrt{\text{Var}(P_Z(B_i))}}{E(P_Z(B_i))} = \frac{1}{n^{(K-4s)/[2(4s+K-1)]}},$$

which is large for small s . For this reason, the DPE fails to achieve the parametric convergence rate when $\frac{K-1}{4}$. We experimentally support this claim in Figure 8, where we use

Figure 8: Estimation error of $T_{m;n}^d$, varying n , s over a log scale. The dashed line has slope $\frac{1}{2}$, which corresponds to the parametric convergence rate. The lines for $s = 0.4$ and 0.6 fail to achieve the parametric convergence rate, even though s is larger than the threshold value 0.25 where the parametric rate arises in standard nonparametric estimation.

$K = 2$ and $Z \sim \text{Unif}([0; 1])$, with a deterministic choice of $Y = 1$. Optimal estimation of $\int_0^1 \text{ECE}(f)^2$ remains an open problem.

C.3 Comparison of Critical Values

We compare three choices of critical values, where we sample $(Z_i; Y_i)_{i=1}^n$ according to the following rules.

1. Oracle Monte Carlo: $Z_i \sim P_Z, Y_i \sim \text{Cat}(Z_i)$.
2. Full bootstrapping (consistency resampling): $Z_i \sim \text{Unif}(f(Z_i)_{i=1}^n), Y_i \sim \text{Cat}(Z_i)$.
3. Y-only bootstrapping: $Z_i = Z_i, Y_i \sim \text{Cat}(Z_i)$.

Here, $f(Z_i; Y_i)_{i=1}^n$ is the original calibration sample, and all sampling is done independently. We repeat the above sampling procedure N times to create $f(Z_i^j; Y_i^j)_{i=1}^n, j \in [N]$ and compute the DPE $T_{m;n}^{d;j}$ test statistics for each $j \in [N]$. Denote their order statistics by $T_{m;n}^{d;(1)}, \dots, T_{m;n}^{d;(N)}$. For oracle Monte Carlo and Y-only bootstrapping, the DPEs $T_{m;n}^d; T_{m;n}^{d;1}; \dots; T_{m;n}^{d;N}$ are exchangeable under the null, so (assuming ties happen with zero probability)

$$P(T_{m;n}^d \leq T_{m;n}^{d;(j)}) = \frac{N+1-j}{N+1};$$

Figure 9: Type I error of the full/ \mathcal{Y} -only bootstrapping. We use $N = 19$; $\alpha = 0.05$ and compute the type I error from 10,000 oracle Monte Carlo trials. Standard error bars are plotted over 1,000 repetitions.

for any $j \in [N]$. In other words, the test that rejects when $I(T_m^d; n, T_m^{d;(j)}; n)$ has type I error $\frac{\alpha}{N+1}$. We do not have such a guarantee for consistency resampling.

Figure 9 compares the type I errors when using critical values based on full/ \mathcal{Y} -only bootstrapping. We see that there is no significant difference between the two, and they stay relatively close to the nominal level $\alpha = 0.05$. Since consistency resampling and \mathcal{Y} -only bootstrapping give randomized critical values, the true type I error is estimated by the sample mean of type I errors obtained from 10,000 independent trials.

C.4 Cross-tting

The idea of cross-tting (see e.g., Hajek, 1962; Schick, 1986; Newey and Robins, 2018; Kennedy, 2020, for related ideas) can be applied to the sample splitting test in Section 6. Specifically, we compute the two-sample test statistic again by swapping the role of $f(Z_i; Y_i)g_{i=1}^{bn=2c}$ and $f(Z_i; Y_i)g_{i=bn=2c+1}^n$, and use the average of two test statistics. The critical value for $\alpha = 0.05$ and the corresponding type II error are estimated via 1,000 oracle Monte Carlo simulations. In Figure 10, we compare the type II error of the cross-tted test statistics with other tests discussed in the main text. The experimental setup is identical to Figure 7. We observe that the cross-tting procedure does not significantly improve the power.

(a) $s = 0.6$, $n = 100$ (b) $s = 0.8$, $n = 200$

Figure 10: Type II error comparison for plug-in, sample-splitting, and cross-fitting test. The horizontal dashed line indicates a type II error of $1 - \alpha = 0.95$. Standard error bars are plotted over 10 repetitions.

References

- Anastasios N Angelopoulos, Stephen Bates, Emmanuel J Candès, Michael I Jordan, and Lihua Lei. Learn then test: Calibrating predictive algorithms to achieve risk control. arXiv preprint arXiv:2110.01052, 2021.
- Ery Arias-Castro, Bruno Pelletier, and Venkatesh Saligrama. Remember the curse of dimensionality: the case of goodness-of-fit testing in arbitrary dimension. *Journal of Non-parametric Statistics*, 30(2):448{471, 2018.
- Imanol Arrieta-Ibarra, Paman Gujral, Jonathan Tannen, Mark Tygert, and Cherie Xu. Metrics of calibration for probabilistic predictions. *Journal of Machine Learning Research*, 23(351):1{54, 2022.
- Peter C Austin and Ewout W Steyerberg. Graphical assessment of internal and external calibration of logistic regression models by using loess smoother. *Statistics in Medicine*, 33(3):517{535, 2014.
- Yu Bai, Song Mei, Huan Wang, and Caiming Xiong. Don't just blame over-parametrization for over-confidence: Theoretical analysis of calibration in binary classification. In *International Conference on Machine Learning PMLR*, 2021.
- Sivaraman Balakrishnan and Larry Wasserman. Hypothesis testing for high-dimensional multinomials: a selective review. *The Annals of Applied Statistics*, 12(2):727{749, 2018.
- Richard Berk. An impact assessment of machine learning risk forecasts on parole board decisions and recidivism. *Journal of Experimental Criminology*, 13(2):193{216, 2017.
- Piotr Berman, Sofya Raskhodnikova, and Grigory Yaroslavtsev. L_p -testing. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing* 2014.

- Thomas B Berrett, Ioannis Kontoyiannis, and Richard J Samworth. Optimal rates for independence testing via u-statistic permutation tests. *The Annals of Statistics*, 49(5): 2457{2490, 2021.
- J Eric Bickel. Some comparisons among quadratic, spherical, and logarithmic scoring rules. *Decision Analysis*, 4(2):49{65, 2007.
- P. J. Bickel and Y. Ritov. Estimating integrated squared density derivatives: Sharp best order of convergence estimates. *Sankhya: The Indian Journal of Statistics, Series A (1961-2002)*, 50(3):381{393, 1988.
- Lucien Birgé and Pascal Massart. Estimation of integral functionals of a density. *The Annals of Statistics*, 23(1):11{29, 1995.
- Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. arXiv preprint arXiv:1604.07316, 2016.
- Glenn W Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1):1{3, 1950.
- Jochen Brecker. Estimating reliability and resolution of probability forecasts through decomposition of the empirical score. *Climate Dynamics*, 39(3):655{667, 2012.
- Jochen Brecker and Leonard A Smith. Increasing the reliability of reliability diagrams. *Weather and Forecasting* 22(3):651{661, 2007.
- Lawrence D Brown and Mark G Low. Asymptotic equivalence of nonparametric regression and white noise. *The Annals of Statistics*, 24(6):2384{2398, 1996.
- Andreas Buja, Werner Stuetzle, and Yi Shen. Loss functions for binary class probability estimation and classification: structure and applications, 2005.
- MV Burnashev. On the minimax detection of an inaccurately known signal in a white gaussian noise background. *Theory of Probability & Its Applications*, 24(1):107{119, 1979.
- Cristina Butucea and Karine Tribouley. Nonparametric homogeneity tests. *Journal of Statistical Planning and Inference*, 136(3):597{639, 2006.
- T Tony Cai and Mark G Low. Optimal adaptive estimation of a quadratic functional. *The Annals of Statistics*, 34(5):2298{2325, 2006.
- Yafo Chen. Pytorch cifar models, 2021. URL <https://github.com/chenyafo/pytorch-cifar-models>.
- Julien Chhor and Alexandra Carpentier. Goodness-of-fit testing for Hölder-continuous densities: Sharp local minimax rates. arXiv preprint arXiv:2109.04346, 2021.
- David R Cox. Two further applications of a model for binary regression. *Biometrika*, 45 (3/4):562{565, 1958.

- A Philip Dawid. The well-calibrated bayesian. *Journal of the American Statistical Association*, 77(379):605{610, 1982.
- A Philip Dawid. Present position and potential developments: Some personal views: Statistical theory: The prequential approach. *Journal of the Royal Statistical Society. Series A (General)*, 147(2):278{292, 1984.
- Christina Dawkins, Thirukodikaval Nilakanta Srinivasan, and John Whalley. Calibration. In *Handbook of econometrics* volume 5, pages 3653{3703. Elsevier, 2001.
- Bruno De Finetti. Does it make sense to speak of "good probability appraisers". *The Scientist Speculates: An Anthology of Partly-baked Ideas* pages 257{364, 1962.
- Morris H DeGroot and Stephen E Fienberg. The comparison and evaluation of forecasters. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 32(1-2):12{22, 1983.
- Francis X Diebold and Roberto S Mariano. Comparing predictive accuracy. *Journal of Business & Economic Statistics* 13(3):134{144, 1995.
- Francis X Diebold, Todd A Gunther, and Anthony S Tay. Evaluating density forecasts with applications to nancial risk management. *International Economic Review*, 39(4): 863{883, 1998.
- David Donoho and Jiashun Jin. Higher criticism for large-scale inference, especially for rare and weak e ects. *Statistical Science* 30(1):1{25, 2015.
- David L Donoho and Michael Nussbaum. Minimax quadratic estimation of a quadratic functional. *Journal of Complexity*, 6(3):290{323, 1990.
- Amandine Dubois, Thomas Berrett, and Cristina Butucea. Goodness-of- t testing for Hölder continuous densities under local differential privacy. arXiv preprint arXiv:2107.02439, 2021.
- Sam Efromovich and Mark Low. On optimal adaptive estimation of a quadratic functional. *The Annals of Statistics*, 24(3):1106{1125, 1996.
- Andre Esteva, Brett Kuprel, Roberto A. Novoa, Justin M. Ko, Susan M. Swetter, Helen M. Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542:115{118, 2017.
- Herbert Federer. *Geometric measure theory* Springer, 2014.
- Christopher AT Ferro and Thomas E Fricker. A bias-corrected decomposition of the brier score. *Quarterly Journal of the Royal Meteorological Society* 138(668):1954{1960, 2012.
- Dean P Foster and Rakesh V Vohra. Asymptotic calibration. *Biometrika*, 85(2):379{390, 1998.
- Allan Franklin. Calibration. In *Can that be Right?*, pages 237{272. Springer, 1999.

- Ursula Garczarek. Classification rules in standardized partition spaces PhD thesis, Universität Dortmund, 2002.
- Martin Gebel. Multivariate calibration of classifier scores into the probability space PhD thesis, Universität Dortmund, 2009.
- Evarist Giré and Richard Nickl. A simple adaptive estimator of the integrated square of a density. *Bernoulli*, 14(1):47{61, 2008.
- Evarist Giré and Richard Nickl. *Mathematical foundations of infinite-dimensional statistical models* Cambridge University Press, 2021.
- Tilmann Gneiting and Matthias Katzfuss. Probabilistic forecasting. *Annual Review of Statistics and Its Application*, 1:125{151, 2014.
- Tilmann Gneiting, Fadoua Balabdaoui, and Adrian E Raftery. Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(2):243{268, 2007.
- Irving John Good. Rational decisions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 14(1):107{114, 1952.
- Ronald L. Graham, Donald E. Knuth, and Oren Patashnik. *Concrete Mathematics: A Foundation for Computer Science* Addison-Wesley Longman Publishing Co., Inc., USA, 2nd edition, 1994. ISBN 0201558025.
- Chuan Guo, Geo Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning* PMLR, 2017.
- Chirag Gupta and Aaditya Ramdas. Top-label calibration and multiclass-to-binary reductions. In *International Conference on Learning Representations 2022*. URL <https://openreview.net/forum?id=WqoBaaPHS-> .
- Chirag Gupta, Aleksandr Podkopaev, and Aaditya Ramdas. Distribution-free binary classification: prediction sets, confidence intervals and calibration. *Advances in Neural Information Processing Systems*33, 2020.
- Kartik Gupta, Amir Rahimi, Thalaiyasingam Ajanthan, Thomas Mensink, Cristian Sminchisescu, and Richard Hartley. Calibration of neural networks using splines. In *International Conference on Learning Representations 2021*. URL <https://openreview.net/forum?id=eQe8DEWNN2W>
- Laszlo Györfi, Michael Kohler, Adam Krzyżak, and Harro Walk. *A distribution-free theory of nonparametric regression* Springer, 2002.
- Jaroslav Hajek. Asymptotically Most Powerful Rank-Order Tests. *The Annals of Mathematical Statistics*, 33(3):1124{1147, 1962.
- David J Hand. *Construction and assessment of classification rules* Wiley, 1997.

- Frank E Harrell. Regression modeling strategies: with applications to linear models, logistic and ordinal regression, and survival analysis Springer, 2015.
- Trevor Hastie and Robert Tibshirani. Classification by pairwise coupling. The Annals of Statistics, 26(2):451 { 471, 1998.
- L Held, Kaspar Ru bach, and Fadoua Balabdaoui. A score regression approach to assess calibration of continuous probabilistic predictions. Biometrics, 66(4):1295{1305, 2010.
- Jorgen Hilden, J Dik F Habbema, and Beth Bjerregaard. The measurement of performance in probabilistic diagnosis. Methods of Information in Medicine, 17(4):227{237, 1978.
- David W Hosmer and Stanley Lemeshow. Goodness of fit tests for the multiple logistic regression model. Communications in Statistics-Theory and Methods 9(10):1043{1069, 1980.
- Yingxiang Huang, Wentao Li, Fima Macheret, Rodney A Gabriel, and Lucila Ohno-Machado. A tutorial on calibration measurements and calibration models for clinical prediction models. Journal of the American Medical Informatics Association, 27(4):621{ 633, 2020.
- Irina A Ingster, Yuri Iand Suslina. Nonparametric goodness-of-fit testing under Gaussian models Springer Science & Business Media, 2012.
- Yuri I Ingster. An asymptotic minimax testing of nonparametric hypotheses on the density of the distribution of an independent sample. Journal of Soviet Mathematics 33(1): 744{758, 1986.
- Yuri I Ingster. Minimax testing of nonparametric hypotheses on a distribution density in the L_p metrics. Theory of Probability & Its Applications , 31(2):333{337, 1987.
- Yuri I Ingster. Adaptive chi-square tests. Journal of Mathematical Sciences 99(2):1110{ 1119, 2000.
- Joan Ivanov, Jack V Tu, and C David Naylor. Ready-made, recalibrated, or remodeled? issues in the use of risk indexes for assessing mortality after coronary artery bypass graft surgery. Circulation , 99(16):2098{2104, 1999.
- Jiashun Jin and Zheng Tracy Ke. Rare and weak effects in large-scale inference: methods and phase diagrams. Statistica Sinica, pages 1{34, 2016.
- Ian T Jolliffe and David B Stephenson. Forecast verification: a practitioner's guide in atmospheric science John Wiley & Sons, 2012.
- Joseph B Kadane and Sarah Lichtenstein. A subjectivist view of calibration. Technical report, Decision Research, Eugene, OR, 1982.
- Daniel Kahneman and Amos Tversky. Prospect theory: an analysis of decision under risk. In Handbook of the Fundamentals of Financial Decision Making: Part J pages 99{127. World Scientific, 2013.

- Edward H Kennedy. Optimal doubly robust estimation of heterogeneous causal effects. arXiv preprint arXiv:2004.14497, 2020.
- Gideon Keren. Calibration and probability judgements: conceptual and methodological issues. *Acta Psychologica* 77(3):217{273, 1991.
- Gerard Kerkycharian and Dominique Picard. Estimating nonquadratic functionals of a density using haar wavelets. *The Annals of Statistics*, 24(2):485{507, 1996.
- Ilmun Kim, Sivaraman Balakrishnan, and Larry Wasserman. Minimax optimality of permutation tests. *The Annals of Statistics*, 50(1):225{251, 2022.
- Keiichi Kisamori, Motonobu Kanagawa, and Keisuke Yamazaki. Simulator calibration under covariate shift with kernels. In *International Conference on Artificial Intelligence and Statistics*, pages 1244{1253. PMLR, 2020.
- Jan Kodovsky and Jessica Fridrich. Calibration revisited. In *Proceedings of the 11th ACM workshop on Multimedia and security 2009*.
- Meelis Kull, Telmo M Silva Filho, and Peter Flach. Beyond sigmoids: how to obtain well-calibrated probabilities from binary classifiers with beta calibration. *Electronic Journal of Statistics*, 11(2):5052{5080, 2017.
- Meelis Kull, Miquel Perello Nieto, Markus Kangsepp, Telmo Silva Filho, Hao Song, and Peter Flach. Beyond temperature scaling: obtaining well-calibrated multi-class probabilities with dirichlet calibration. *Advances in Neural Information Processing Systems* 2019.
- Ananya Kumar, Percy Liang, and Tengyu Ma. Verified uncertainty calibration. In *Advances in Neural Information Processing Systems* 2019.
- Aviral Kumar, Sunita Sarawagi, and Ujjwal Jain. Trainable calibration measures for neural networks from kernel mean embeddings. In *International Conference on Machine Learning*. PMLR, 2018.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems* 2017.
- Beatrice Laurent. Efficient estimation of integral functionals of a density. *The Annals of Statistics*, 24(2):659{681, 1996.
- Erich L Lehmann and Joseph P Romano. *Testing statistical hypotheses* Springer Science & Business Media, 2005.
- Oleg V Lepski. On a problem of adaptive estimation in gaussian white noise. *Theory of Probability & Its Applications*, 35(3):454{466, 1991.
- Oleg V Lepski and Vladimir G Spokoiny. Optimal pointwise adaptive methods in nonparametric estimation. *The Annals of Statistics*, 25(6):2512{2546, 1997.

- Sarah Lichtenstein, Baruch Fischhoff, and Lawrence D Phillips. Calibration of probabilities: the state of the art. In *Decision Making and Change in Human Affairs: Proceedings of the Fifth Research Conference on Subjective Probability, Utility, and Decision Making* pages 275-324. Springer, 1977.
- Mark G Low. On nonparametric confidence intervals. *The Annals of Statistics*, 25(6): 2547-2554, 1997.
- HB Mann and Abraham Wald. On the choice of the number of class intervals in the application of the chi square test. *The Annals of Mathematical Statistics*, 13(3):306-317, 1942.
- Dimitrios Miliotis, Raffaele Camoriano, Pietro Michiardi, Lorenzo Rosasco, and Maurizio Filippone. Dirichlet-based gaussian processes for large-scale calibrated classification. In *Advances in Neural Information Processing Systems* 2018.
- Michael E Miller, Siu L Hui, and William M Tierney. Validation techniques for logistic regression models. *Statistics in Medicine*, 10(8):1213-1226, 1991.
- Michael E Miller, Carl D Langefeld, William M Tierney, Siu L Hui, and Clement J McDonald. Validation of probabilistic predictions. *Medical Decision Making*, 13(1):49-57, 1993.
- Robert G Miller. *Statistical prediction by discriminant analysis*. Springer, 1962.
- Jacob A Mincer and Victor Zarnowitz. The evaluation of economic forecasts. In *Economic forecasts and expectations: Analysis of forecasting behavior and performance* pages 3-46. NBER, 1969.
- Rajarshi Mukherjee, Eric Tchetgen Tchetgen, and James Robins. Lepski's method and adaptive estimation of nonlinear integral functionals of density. arXiv preprint arXiv:1508.00249, 2015.
- Jishnu Mukhoti, Viveka Kulharia, Amartya Sanyal, Stuart Golodetz, Philip Torr, and Puneet Dokania. Calibrating deep neural networks using focal loss. In *Advances in Neural Information Processing Systems* 2020.
- Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *AAAI Conference on Artificial Intelligence*, 2015.
- Whitney K Newey and James R Robins. Cross-fitting and fast remainder rates for semi-parametric estimation. arXiv preprint arXiv:1801.09138, 2018.
- Alexandru Niculescu-Mizil and Rich Caruana. Predicting good probabilities with supervised learning. In *International Conference on Machine Learning* 2005.
- Jeremy Nixon, Michael W Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran. Measuring calibration in deep learning. In *CVPR Workshops*, 2019.
- Christine Osborne. Statistical calibration: A review. *International Statistical Review / Revue Internationale de Statistique* 59(3):309-336, 1991.

- John Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *Advances in Large Margin Classifiers*, pages 61{74, 1999.
- C Radhakrishna Rao. Large sample tests of statistical hypotheses concerning several parameters with applications to problems of estimation. In *Mathematical Proceedings of the Cambridge Philosophical Society* volume 44, pages 50{57. Cambridge University Press, 1948.
- James Robins, Lingling Li, Eric Tchetgen, Aad van der Vaart, et al. Higher order influence functions and minimax estimation of nonlinear functionals. *Probability and statistics: essays in honor of David A. Freedman* 2:335{421, 2008.
- Leonard J Savage. Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, 66(336):783{801, 1971.
- Mark J Schervish. A general method for comparing probability assessors. *The Annals of Statistics*, 17(4):1856{1879, 1989.
- Anton Schick. On Asymptotically Efficient Estimation in Semiparametric Models. *The Annals of Statistics*, 14(3):1139 { 1151, 1986.
- Francoise Seillier-Moiseiwitsch and A Philip Dawid. On testing the validity of sequential probability forecasts. *Journal of the American Statistical Association*, 88(421):355{359, 1993.
- Nicolas Serrano. Calibration strategies to validate predictive models: is new always better? *Intensive Care Medicine*, 38(8):1246{1248, 2012.
- Nilay D Shah, Ewout W Steyerberg, and David M Kent. Big data and predictive analytics: recalibrating expectations. *JAMA* , 320(1):27{28, 2018.
- Yandi Shen, Chao Gao, Daniela Witten, and Fang Han. Optimal estimation of variance in nonparametric regression with random design. *The Annals of Statistics*, 48(6):3589{3618, 2020.
- Samuel D Silvey. The lagrangian multiplier test. *The Annals of Mathematical Statistics*, 30(2):389{407, 1959.
- David J Spiegelhalter. Probabilistic prediction in patient management and clinical trials. *Statistics in Medicine*, 5(5):421{433, 1986.
- Vladimir G Spokoiny. Adaptive hypothesis testing using wavelets. *The Annals of Statistics*, 24(6):2477{2498, 1996.
- Richard Stanley. Eulerian partitions of a unit hypercube. *Higher Combinatorics*, 31:49, 1977.
- Ewout W Steyerberg. *Clinical prediction models*. Springer, 2019.

- Ewout W Steyerberg, Andrew J Vickers, Nancy R Cook, Thomas Gerds, Mithat Gonen, Nancy Obuchowski, Michael J Pencina, and Michael W Kattan. Assessing the performance of prediction models: a framework for some traditional and novel measures. *Epidemiology (Cambridge, Mass.)* 21(1):128, 2010.
- Bernd Sturmfels. *Grobner Bases and Convex Polytopes*. American Mathematical Soc., 1996.
- Ambrus Tamás and Balázs Csáji. Exact distribution-free hypothesis tests for the regression function of binary classification via conditional kernel mean embeddings. *IEEE Control Systems Letters* 2021.
- Eric Tchetgen, Lingling Li, James Robins, and Aad van der Vaart. Minimax estimation of the integral of a power of a density. *Statistics & probability letters*, 78(18):3307{3311, 2008.
- Philip E Tetlock and Dan Gardner. *Superforecasting: The Art and Science of Prediction*. Random House, 2016.
- Sunil Thulasidasan, Gopinath Chennupati, Je A. Bilmes, Tanmoy Bhattacharya, and Sarah Ellen Michalak. On mixup training: improved calibration and predictive uncertainty for deep neural networks. In *Advances in Neural Information Processing Systems*, 2019.
- DB Toll, KJM Janssen, Y Vergouwe, and KGM Moons. Validation, updating and impact of clinical prediction rules: a review. *Journal of Clinical Epidemiology*, 61(11):1085{1094, 2008.
- Juozas Vaicenavicius, David Widmann, Carl Andersson, Fredrik Lindsten, Jacob Roll, and Thomas Schön. Evaluating model calibration in classification. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2019.
- Ben Van Calster and Andrew J Vickers. Calibration of risk prediction models: impact on decision-analytic performance. *Medical Decision Making*, 35(2):162{169, 2015.
- Ben Van Calster, Daan Nieboer, Yvonne Vergouwe, Bavo De Cock, Michael J Pencina, and Ewout W Steyerberg. A calibration hierarchy for risk models was needed: from utopia to empirical data. *Journal of Clinical Epidemiology*, 74:167{176, 2016.
- Ben Van Calster, David J McLernon, Maarten Van Smeden, Laure Wynants, and Ewout W Steyerberg. Calibration: the achilles heel of predictive analytics. *BMC medicine*, 17(1):1{7, 2019.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*. Cambridge University Press, 2018.
- Vladimir Vovk and Glenn Shafer. Good randomized sequential probability forecasting is always possible. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 67(5):747{763, 2005.

- Vladimir Vovk, Ivan Petej, Paolo Toccaceli, Alexander Gammerman, Ernst Ahlberg, and Lars Carlsson. Conformal calibrators. In *Conformal and Probabilistic Prediction and Applications*. PMLR, 2020.
- Lie Wang, Lawrence D Brown, T Tony Cai, and Michael Levine. Effect of mean on variance function estimation in nonparametric regression. *The Annals of Statistics*, 36(2):646–664, 2008.
- David Widmann, Fredrik Lindsten, and Dave Zachariah. Calibration tests in multi-class classification: a unifying framework. *Advances in Neural Information Processing Systems*, 2019.
- Robert L Winkler, Javier Munoz, José L Cervera, José M Bernardo, Gail Blattenberger, Joseph B Kadane, Dennis V Lindley, Allan H Murphy, Robert M Oliver, and David Ríos-Insua. Scoring rules and the evaluation of probabilities. *Test*, 5(1):1–60, 1996.
- Bianca Zadrozny and Charles Peter Elkan. Obtaining calibrated probability estimates from decision trees and naive bayesian classifiers. In *International Conference of Machine Learning*, 2001.
- Bianca Zadrozny and Charles Peter Elkan. Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining*, 2002.
- Jize Zhang, Bhavya Kaikhura, and T Yong-Jin Han. Mix-n-match: ensemble and compositional methods for uncertainty calibration in deep learning. In *International Conference on Machine Learning*. PMLR, 2020.
- Shengjia Zhao, Tengyu Ma, and Stefano Ermon. Individual calibration with randomized forecasting. In *International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2020.