

Optimizing Estimators of Squared Calibration Errors in Classification

Anonymous authors

Paper under double-blind review

Abstract

In this work, we propose a mean-squared error-based risk that enables the comparison and optimization of estimators of squared calibration errors in practical settings. Improving the calibration of classifiers is crucial for enhancing the trustworthiness and interpretability of machine learning models, especially in sensitive decision-making scenarios. Although various calibration (error) estimators exist in the current literature, there is a lack of guidance on selecting the appropriate estimator and tuning its hyperparameters. By leveraging the bilinear structure of squared calibration errors, we reformulate calibration estimation as a regression problem with independent and identically distributed (i.i.d.) input pairs. This reformulation allows us to quantify the performance of different estimators even for the most challenging calibration criterion, known as canonical calibration. Our approach advocates for a training-validation-testing pipeline when estimating a calibration error on an evaluation dataset. We demonstrate the effectiveness of our pipeline by optimizing existing calibration estimators and comparing them with novel kernel ridge regression-based estimators on standard image classification tasks.

1 Introduction

In the field of machine learning, classification tasks involve predicting discrete class labels for given instances (Bishop & Nasrabadi, 2006). As these models are increasingly being employed in critical applications such as healthcare (Haggenmüller et al., 2021), autonomous driving (Feng et al., 2020), weather forecasting (Gneiting & Raftery, 2005), and financial decision-making (Frydman et al., 1985), the need for reliable and interpretable predictions has become of critical importance. A key aspect of reliability in classification models is the calibration of their predicted probabilities (Murphy & Winkler, 1977; Hekler et al., 2023). Calibration refers to the alignment between predicted probabilities and the true likelihoods of outcomes, ensuring that predictions are not only accurate but also meaningful in terms of their confidence scores (Murphy, 1973). Despite the advancements in model architectures and learning algorithms, many modern classifiers, such as deep neural networks, are prone to producing overconfident predictions (Minderer et al., 2021). This overconfidence can be attributed to several factors, including the model’s complexity, training data limitations, and inherent biases in learning processes (Guo et al., 2017). Consequently, even models that achieve high accuracy might suffer from poor calibration, leading to potential misinterpretations and suboptimal decisions.

To quantify the extent to which a model is miscalibrated, calibration errors have been introduced (Naeini et al., 2015). However, their estimators are usually biased (Roelofs et al., 2022) and inconsistent (Vaicenavicius et al., 2019). Other calibration errors with unbiased estimators exist but they lack theoretical derivation and are difficult to interpret (Widmann et al., 2019; 2021; Marx et al., 2024). This, in turn, is highly problematic since we cannot quantify how reliable a model is if we either do not know how reliable the metric is or how to interpret it. In this work, we tackle the former problem. We propose mean-squared error risk minimization to find an optimal calibration estimator for any squared calibration error. This risk can be applied in any practical scenario to compare and select different estimators, and works for all notions of calibration, even for the notoriously difficult canonical calibration. Our **contributions** are as follows:

- We propose a novel risk applicable for squared calibration estimators in Section 3.1, which represents an optimization objective for comparing calibration estimators proposed by the literature.
- We formulate a calibration-evaluation pipeline based on our risk in Section 3.2.2, which allows to optimize calibration estimators used by the literature.
- We propose novel kernel ridge regression-based calibration estimators in Section 4, and compare these with optimized baselines on common image classification models in Section 5.

2 Background

In the following, we offer an extensive introduction into the background of this work. First, we give briefly measure theoretic preliminaries, followed by the different notions of calibration used for classification by the literature. We then introduce and discuss commonly used estimators for these notions.

2.1 Measure Theoretic Preliminaries

Throughout this work we use measure theoretic definitions to formalise our contribution, its assumptions, and its limitations. Specifically, we assume a measure space $(\Omega, \mathcal{F}, \mu)$ is given. We associate with a random variable $X: \Omega \rightarrow \mathcal{X}$ (a measurable function) a probability space $(\mathcal{X}, \mathcal{F}_X, \mathbb{P}_X)$ with σ -field $\mathcal{F}_X = \{X(A) \mid A \in \mathcal{F}\}$ and pushforward measure $\mathbb{P}_X = \mu \circ X^{-1}$. Further, we may also assume a second random variable $Y: \Omega \rightarrow \mathcal{Y}$ such that we have a joint probability space $(\mathcal{X} \times \mathcal{Y}, \mathcal{F}_{XY}, \mathbb{P}_{XY})$, where \mathcal{F}_{XY} and \mathbb{P}_{XY} are defined analogous as \mathcal{F}_X and \mathbb{P}_X . Further, if a random variable Y is discrete, we use \mathbb{P}_Y to represent both its probability measure and the associated probability vector in the simplex $\Delta^d := \{(p_1, \dots, p_k)^\top \in [0, 1]^d \mid \sum_{i=1}^d p_i = 1\}$, where d is the number of unique outcomes. In the same sense we may use $\mathbb{P}_{Y|X} := \frac{d\mathbb{P}_{XY}}{d\mathbb{P}_X}$ as a measurable function $\mathcal{X} \rightarrow \Delta^d$. Finally, we say a statement holds μ -almost surely (μ -a.s.) if it is true except on a set of μ -measure zero (Capiński & Kopp, 2004).

2.2 Mean-Squared Error Risk Minimization and Calibration

The mean-squared error (MSE) is one of the most common loss functions used in regression problems, see, e.g., (Efron, 1994; Schölkopf & Smola, 2002; Bishop & Nasrabadi, 2006; Goodfellow et al., 2016; Murphy, 2022; Bach, 2024). Its expected loss for a vector-valued sample $Y \sim \mathbb{P}_Y$ from a target distribution \mathbb{P}_Y and a prediction $c \in \mathbb{R}^d$ is defined by

$$L_{\text{MSE}}(c, \mathbb{P}_Y) := \mathbb{E}_{Y \sim \mathbb{P}_Y} [\|c - Y\|^2]. \quad (1)$$

It holds that the expectation of the target term in the difference is the unique minimizer, i.e.,

$$\mathbb{E}[Y] = \arg \min_{c \in \mathbb{R}^d} L_{\text{MSE}}(c, \mathbb{P}_Y). \quad (2)$$

Now, let's introduce another random variable X such that $(X, Y) \sim \mathbb{P}_{XY}$ follows a joint distribution and X can be used as input for a regression model $m: \mathcal{X} \rightarrow \mathbb{R}^d$ to approximate the target conditional distribution $m^*(x) := \mathbb{E}[Y \mid X = x]$. Then, the corresponding risk of m is defined by

$$\mathcal{R}_{\text{MSE}}(m) := \mathbb{E}_{X \sim \mathbb{P}_X} [L_{\text{MSE}}(m(X), \mathbb{P}_{Y|X})] = \mathbb{E}_{(X, Y) \sim \mathbb{P}_{XY}} [\|m(X) - Y\|^2]. \quad (3)$$

Given $m \neq m^*$, where the inequality holds for a set of positive probability mass, it follows that

$$\mathcal{R}_{\text{MSE}}(m) > \mathcal{R}_{\text{MSE}}(m^*). \quad (4)$$

However, we have the measure theoretic limitation that there might be a null set $A \in \mathcal{F}_X$ and some \tilde{m} such that $\tilde{m}(x) \neq m^*(x)$ for all $x \in A$ while $\mathcal{R}_{\text{MSE}}(\tilde{m}) = \mathcal{R}_{\text{MSE}}(m^*)$. This limitation of risk minimization is usually irrelevant in machine learning applications. However, in our case, it will require additional theoretical assumptions on the target conditional distribution to make risk minimization a usable tool for assessing calibration estimators.

The MSE can also be used for classification tasks with a one-hot encoded target in Equation 1, often referred to as Brier score (Brier, 1950). Then, Murphy (1973) introduces the concept of calibration for a classifier $f: \mathcal{X} \rightarrow \Delta^d$ by showing that

$$\mathcal{R}_{\text{MSE}}(f) = \mathbb{E}_{X \sim \mathbb{P}_X} [\|f(X) - \mathbb{P}_{Y|f(X)}\|^2] - \mathbb{E}_{X \sim \mathbb{P}_X} [\|\mathbb{P}_Y - \mathbb{P}_{Y|f(X)}\|^2] + \|\mathbb{P}_Y\|^2. \quad (5)$$

The first term on the right-hand side is usually referred to as the calibration term, and the second as sharpness term, coining Equation (5) as calibration-sharpness decomposition of the Brier score (Gneiting et al., 2007; Gruber & Buettner, 2022; Kuleshov & Deshpande, 2022; Gruber et al., 2024; Sun et al., 2024). In the calibration term, $f(X)$ is compared with the target distribution $\mathbb{P}(Y | f(X))$ given the full predicted vector. In current literature, this notion of calibration is referred to as canonical calibration (Vaicenavicius et al., 2019; Popordanoska et al., 2022b; Gupta & Ramdas, 2022; Gruber et al., 2024). Formally, we say the model f is **canonically calibrated** if and only if

$$\mathbb{P}(Y = i | f(X) = p) = p_i \text{ for all } p \in \Delta^d, i = 1 \dots d. \quad (6)$$

The corresponding L^2 **canonical calibration error** is defined as

$$\text{CCE}_2(f) := \sqrt{\mathbb{E}[\|f(X) - \mathbb{P}_{Y|f(X)}\|^2]}, \quad (7)$$

which is equal to the calibration term in Equation (5). It holds that $\text{CCE}_2(f) = 0$ if and only if f is canonically calibrated \mathbb{P}_X -almost surely.

However, canonical calibration errors are notoriously difficult to estimate and represent a calibration strictness which may not be necessary in practice (Vaicenavicius et al., 2019). Consequently, other notions of calibration have been proposed, which we discuss next.

2.3 Alternative Notions of Calibration

Besides canonical calibration, multiple notions of calibration have been introduced in the literature (Zadrozny & Elkan, 2002; Vaicenavicius et al., 2019; Kull et al., 2019; Gupta & Ramdas, 2022). Respective calibration errors assess the degree a classifier violates a given notion. In recent literature, the most common notion is top-label confidence calibration (Naeini et al., 2015; Guo et al., 2017; Joo et al., 2020; Kristiadi et al., 2020; Rahimi et al., 2020; Tomani et al., 2021; Minderer et al., 2021; Tian et al., 2021; Islam et al., 2021; Menon et al., 2021; Morales-Álvarez et al., 2021; Gupta et al., 2021; Wang et al., 2021; Fan et al., 2022; Dehghani et al., 2023; Chang et al., 2024). Here, we compare if the predicted top-label confidence $\max_{i \in \mathcal{Y}} f_i(X)$ matches the conditional accuracy $\mathbb{P}(Y = \arg \max_{i \in \mathcal{Y}} f_i(X) | \max_{i \in \mathcal{Y}} f_i(X))$ given the prediction. Formally, we say the classifier f is **top-label confidence calibrated** if and only if

$$\mathbb{P}\left(Y = \arg \max_i f_i(X) \mid \max_i f_i(X) = p\right) = p \text{ for all } p \in [0, 1]. \quad (8)$$

The corresponding L^2 **top-label confidence calibration error** is defined as

$$\text{TCE}_2(f) := \sqrt{\mathbb{E}\left[\left(\max_i f_i(X) - \mathbb{P}\left(Y = \arg \max_i f_i(X) \mid \max_i f_i(X)\right)\right)^2\right]}. \quad (9)$$

It holds that $\text{TCE}_2(f) = 0$ if and only if f is top-label confidence calibrated \mathbb{P}_X -almost surely. It is easier to estimate than CCE_2 since the target conditional distribution is only based on a scalar random variable independent of the number of classes. However, top-label confidence calibration is a weaker condition than canonical calibration since the implication $\text{CCE}_2(f) = 0 \implies \text{TCE}_2(f) = 0$ does not generally hold in the reverse direction (Gruber & Buettner, 2022).

Other notions of calibration exist, which also reduce the prediction to a scalar, such as class-wise calibration (Zadrozny & Elkan, 2002; Kull et al., 2019; Kumar et al., 2019). Gupta & Ramdas (2022) introduce further of such notions. In general, these notions are transformations of the full probability vectors to a lower dimensional space. Similar to top-label confidence calibration, this makes them easier to estimate than canonical calibration but also turns them into a weaker condition due to the information loss of the transformation (Vaicenavicius et al., 2019; Gruber & Buettner, 2022).

2.4 Calibration Estimators

We assume a dataset of i.i.d. samples $(X_1, Y_1), \dots, (X_n, Y_n) \sim \mathbb{P}_{XY}$ to estimate the calibration of a given classifier $f: \mathcal{X} \rightarrow \Delta^d$. The most common approach to estimate calibration errors based on scalar conditionals are binning schemes (Naeini et al., 2015; Guo et al., 2017; Minderer et al., 2021; Detlefsen et al., 2022). A prominent estimator is the so-called *expected calibration error* (ECE), which is a binning-based estimator of the L^1 top-label confidence calibration error (Guo et al., 2017). In essence, the conditional target distribution $\mathbb{P}(Y = \arg \max_i f_i(X) \mid \max_i f_i(X))$ is estimated via a histogram binning scheme, which places all top-label confidence predictions into mutually distinct bins $B_m := \{i \mid \arg \max_j f_j(X_i) \in I_m\}$, $m = 1, \dots, M$, based on a partition $\bigcup_m I_m = [0, 1]$. The analogous L^2 estimator is given by

$$\text{TCE}_2^{\text{bin}} := \sqrt{\sum_{m=1}^M \frac{|B_m|}{n} (\text{acc}(B_m) - \text{conf}(B_m))^2} \quad (10)$$

with $\text{acc}(B) = \frac{1}{|B|} \sum_{i \in B} \mathbf{1}_{Y_i = \arg \max_j f_j(X_i)}$ and $\text{conf}(B) = \frac{1}{|B|} \sum_{i \in B} \arg \max_j f_j(X_i)$ (Kumar et al., 2019). This estimator is primarily suitable for target distributions conditioned on a scalar random variable. The choice of bin intervals I_1, \dots, I_M is user-defined and the estimator only converges to TCE_2 for an adaptive scheme (Vaicenavicius et al., 2019). Patel et al. (2021) and Roelofs et al. (2022) propose approaches to automatically select appropriate bins. However, in practice, it remains an open challenge to definitively select the optimal choice for a specific dataset and classifier. Analogous binning-based estimators for class-wise calibration exist in the literature and share these limitations (Kumar et al., 2019; Nixon et al., 2019; Vaicenavicius et al., 2019).

Estimating canonical calibration is more difficult than other notions of calibration due to the target distribution $\mathbb{P}(Y \mid f(X))$ being conditioned on a vector-valued random variable. Popordanoska et al. (2022b) propose to use a kernel density ratio estimator, which is closely related to the Nadaraya-Watson-estimator (Bierens, 1996). The estimator for CCE_2 is given by

$$\text{CCE}_2^{\text{kde}} := \sqrt{\frac{1}{n} \sum_{i=1}^n \left\| f(X_i) - \frac{\sum_j k_{\text{dir}}(f(X_j); f(X_i)) e_{Y_j}}{\sum_j k_{\text{dir}}(f(X_j); f(X_i))} \right\|^2}, \quad (11)$$

where e_i refers to the unit vector with a 1 at index i and k_{dir} is chosen to be the Dirichlet kernel, which is specifically suited for the simplex space (Ouimet & Tolosana-Delgado, 2022). The authors also propose analogous kernel density based estimators for top-label and class-wise calibration errors, which we will denote as $\text{TCE}_2^{\text{kde}}$ and $\text{CWCE}_2^{\text{kde}}$. Even though the kernel density approach is advantageous compared to the binning approach for higher dimensions, it still shares some of its limitations. Popordanoska et al. (2022b) show that the estimator converges in the infinite data limit, but it is still not clear what is the optimal choice of kernel and kernel hyperparameters in the finite data regime.

To summarize, a lot of different approaches have been proposed by the literature to estimate calibration errors. However, it is not clear how to compare different estimators and pinpoint an optimal choice in a finite data setup in practice.

3 A Mean-Squared Error Risk for Calibration Estimators

In this section, we present our main contribution: A mean-squared error based risk, which can be applied to compare different calibration estimators in a practical, finite data setup. We first discuss its measure theoretic foundations in Section 3.1, and, then, propose a training-inference pipeline for estimating the calibration error in practice in Section 3.2. This pipeline is analogous to how model training, model selection, and test error evaluation is done in practice in machine learning (Bishop & Nasrabadi, 2006). All formulations are with respect to canonical calibration, since this is the most general case. Other notions, like top-label confidence calibration, can be derived by restricting the canonical case to binary classification. All missing proofs are presented in Appendix C.

3.1 Theoretical Definition and Properties

Note that for the CCE_2 calibration error it is sufficient to find a function $h^*: \Delta^d \times \Delta^d \rightarrow \mathbb{R}$ such that

$$h^*(p, p') = \langle p - \mathbb{P}_{Y|f(X)=p}, p' - \mathbb{P}_{Y|f(X)=p'} \rangle, \quad (12)$$

since from this follows that

$$\mathbb{E}_X [h^*(f(X), f(X))] = \text{CCE}_2. \quad (13)$$

Indeed, in a later section, we will discover that current estimators already implicitly use such a form. In general, we refer to a function $h: \Delta^d \times \Delta^d \rightarrow \mathbb{R}$ as **calibration estimation function**. We now propose a risk, which quantifies how close such a h is to h^* . To achieve this, we slightly modify the mean-squared error loss function of Equation (1) in the following.

Definition 1. For a prediction $c \in \mathbb{R}$, a target product measure $\mathbb{P}_Y \otimes \mathbb{P}_V$, and constants $p, p' \in \Delta^d$ we define the **calibration estimator loss** by

$$L_{\text{CE}}(c, \mathbb{P}_Y \otimes \mathbb{P}_V; p, p') := \mathbb{E}_{(Y,V) \sim \mathbb{P}_Y \otimes \mathbb{P}_V} [((p - e_Y, p' - e_V) - c)^2], \quad (14)$$

where e_i refers to the unit vector with a 1 at index i .

Similar to the mean-squared error, it holds that L_{CE} has an unique minimizer given by

$$\langle p - \mathbb{P}_Y, p' - \mathbb{P}_V \rangle = \arg \min_{c \in \mathbb{R}} L_{\text{CE}}(c, \mathbb{P}_Y \otimes \mathbb{P}_V; p, p'). \quad (15)$$

We use this definition of a novel loss to define the respective risk in the following.

Definition 2. For a calibration estimator function $h: \Delta^d \times \Delta^d \rightarrow \mathbb{R}$, we define the **calibration estimation risk** by

$$\mathcal{R}_{\text{CE}}(h) := \mathbb{E}_{X, X'} [L_{\text{CE}}(h(f(X), f(X')), \mathbb{P}_{Y|f(X)=f(X)} \otimes \mathbb{P}_{Y|f(X)=f(X')}; f(X), f(X'))] \quad (16)$$

with $X, X' \stackrel{iid}{\sim} \mathbb{P}_X$.

Similar to \mathcal{R}_{MSE} in Equation (3), we may also express \mathcal{R}_{CE} in a simpler form, since it holds

$$\mathcal{R}_{\text{CE}}(h) = \mathbb{E}_{X, X', Y, Y'} [((f(X) - e_Y, f(X') - e_{Y'}) - h(f(X), f(X')))^2] \quad (17)$$

with $(X, Y), (X', Y') \stackrel{iid}{\sim} \mathbb{P}_{XY}$. This formulation will be used in a later section to construct the empirical risk. Next, we establish that our proposed risk can distinguish the right solution almost surely.

Theorem 1. For any $h: \Delta^d \times \Delta^d \rightarrow \mathbb{R}$ for which $h = h^*$ does **not** hold $\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)}$ -almost surely we have that

$$\mathcal{R}_{\text{CE}}(h) > \mathcal{R}_{\text{CE}}(h^*). \quad (18)$$

This property would be sufficient in practice, if we were using h with arguments sampled from $\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)}$ to estimate the calibration error. However, as demonstrated by Equation (13), we predict CCE_2 via the same sample in both arguments. Mirroring the arguments results in a possible $\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)}$ -null set since only elements in the diagonal $\mathcal{D}(\Delta^d) := \{(p, p) \mid p \in \Delta^d\} \subset \Delta^d \times \Delta^d$ are considered. Consequently, the diagonal of an optimum h^* identified via \mathcal{R}_{CE} may "slip through" the $\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)}$ -a.s. guarantee in Theorem 1, and in turn may result in $\mathbb{E}[h^*(f(X), f(X))] \neq \text{CCE}_2$. To avoid such theoretical exceptions, we require additional theoretical assumptions, which we state in the following.

Theorem 2. Assume a function $h: \Delta^d \times \Delta^d \rightarrow \mathbb{R}$ is continuous in all points of the diagonal $\mathcal{D}(\Delta^d \setminus A)$ with A being a $\mathbb{P}_{f(X)}$ -null set, and the target as a function $\mathbb{P}_{Y|f(X)}: \Delta^d \rightarrow \Delta^d$ is continuous $\mathbb{P}_{f(X)}$ -almost surely. Further, assume the boundary of the support of $\mathbb{P}_{f(X)}$ does not involve a singular distribution, then it holds that

$$\mathcal{R}_{\text{CE}}(h) = \mathcal{R}_{\text{CE}}(h^*) \implies \mathbb{E}[h(f(X), f(X))] = \text{CCE}_2. \quad (19)$$

This result states under which conditions we can expect that an optimal risk indicates a truthful calibration estimation. We briefly discuss these conditions, which are of purely technical nature and should not influence practical results. First, the continuity of h is non-problematic since it is user-defined and infinitely many points of discontinuity are allowed (as long as they have no probability mass). The continuity of $\mathbb{P}_{Y|f(X)}$ may be considered the most relevant condition in practice, since we usually do not know about the nature of the target distribution. However, again, infinitely many points of discontinuity are allowed, as long as their overall probability mass is zero. The last assumption, namely that the boundary of the support is not allowed to be part of a singular distribution, disqualifies certain theoretically crafted distributions. One such example is the Cantor distribution, which consists of infinitely many disconnected points each of zero probability (Teschl, 2014). In summary, the purpose of the conditions in Theorem 2 is to guarantee that samples from $\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)}$ can be arbitrarily close to the diagonal $\mathcal{D}(\Delta^d)$ and that this closeness indicates how well h matches h^* on $\mathcal{D}(\Delta^d)$.

Remark. *Alternatively to our approach, one might also use a loss for finding a probabilistic model $\hat{g}(p) \approx \mathbb{P}_{Y|f(X)=p}$ and then use $\hat{h}(p, p') := \langle p - \hat{g}(p), p' - \hat{g}(p') \rangle \approx \langle p - \mathbb{P}_{Y|f(X)=p}, p' - \mathbb{P}_{Y|f(X)=p'} \rangle$ as a solution. However, learning a predictive space Δ^d becomes increasingly more difficult for higher dimensions d than a regression problem in \mathbb{R} . For example, our approach is invariant to any orthogonal matrix M since $\langle Mx, My \rangle = \langle x, y \rangle$.*

3.2 Estimating Calibration for finite data

In this section, we propose a novel calibration evaluation pipeline for the finite data regime according to our theory. First, we give an unbiased and consistent estimator of the risk in form of an U-statistic. Then, we mimic the training-validation-testing pipeline of a conventional machine learning model for a debiased calibration estimate. This procedure allows to select between different calibration estimators and optimize their hyperparameters.

3.2.1 Risk Estimator

Note that the risk as formulated in Equation (17), is an expectation of two i.i.d. tuples of random variables (X, Y) and (X', Y') . Consequently, given an i.i.d. dataset $(X_1, Y_1), \dots, (X_n, Y_n) \sim \mathbb{P}_{XY}$, we can construct an U-statistic estimator (Shao, 2003) via

$$\hat{\mathcal{R}}_{\text{CE}}(h) := \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n \left(\langle f(X_i) - e_{Y_i}, f(X_j) - e_{Y_j} \rangle - h(f(X_i), f(X_j)) \right)^2. \quad (20)$$

It holds that $\mathbb{E}[\hat{\mathcal{R}}_{\text{CE}}(h)] = \mathcal{R}_{\text{CE}}(h)$, and $\hat{\mathcal{R}}_{\text{CE}}(h) \rightarrow \mathcal{R}_{\text{CE}}(h)$ in distribution if $n \rightarrow \infty$. The estimator has quadratic complexity in n . However, an estimator with linear complexity can be constructed as well by excluding certain index combinations.

3.2.2 Calibration-Evaluation Pipeline

In general, we cannot expect to find h^* . However, the empirical risk allows us to find a parametrized $h_{\theta, \eta}$ close to h^* , where θ denote its parameters and η its hyperparameters. Consequently, similar to how traditional machine learning works, we need a training-validation-test split to achieve unbiased estimation of the calibration error once we optimized $h_{\theta, \eta}$. Specifically, we use a training set to find θ_{tr} and a validation set to find η_{val} . The final calibration estimate is then computed by

$$\widehat{\text{CE}}_2(f) := \frac{1}{n_{\text{te}}} \sum_{X \in \mathcal{D}_{\text{te}}} h_{\theta_{\text{tr}}, \eta_{\text{val}}}(f(X), f(X)), \quad (21)$$

where \mathcal{D}_{te} is the test set of size n_{te} , and $\widehat{\text{CE}}_2$ is representative for different notions of calibration. We might also use multiple training, validation and test sets via (nested) cross validation. In that case, we average the results of the calibration estimation functions fitted in each fold.

Remark. We encourage to never compare the test set risk of different calibration estimation functions, since this would put a bias on the final calibration estimation. Neglecting this is equivalent to selecting an optimal classifier based on the test accuracy in a classification task.

4 Calibration Estimation Functions

In this section, we first formulate the calibration estimation functions implicitly used in the literature. Then, we introduce two novel calibration estimators based on kernel ridge regression, which minimize regularized versions of the empirical risk in Equation (20). All calibration estimation functions can be seen as preliminary to future research, since we may find function classes with lower validation risk by expanding the search space and computational resources. All missing proofs are located in Appendix C.

4.1 Binning and Kernel Density

In this section, we show that the binning estimator $\text{TCE}_2^{\text{bin}}$ of Equation (10) and the kernel density ratio estimator $\text{CCE}_2^{\text{kde}}$ of Equation (11) are the mean prediction of an implicit calibration estimator function of the form

$$\frac{1}{n} \sum_{i=1}^n h(f(X_i), f(X_i)) \quad (22)$$

for some $h: \Delta^d \times \Delta^d \rightarrow \mathbb{R}$ and an i.i.d. dataset $(X_1, Y_1), \dots, (X_n, Y_n) \sim \mathbb{P}_{XY}$. For the binning based estimator, define

$$h_{\text{bin}}(p, p') := \left(\sum_{m=1}^M (\text{conf}(B_m) - \text{acc}(B_m)) \mathbf{1}_{p \in I_m} \right) \left(\sum_{m=1}^M (\text{conf}(B_m) - \text{acc}(B_m)) \mathbf{1}_{p' \in I_m} \right), \quad (23)$$

where $I_1 \dots I_M$, $\text{acc}(B_m)$, and $\text{conf}(B_m)$ are defined as above in Equation (10). It holds that

$$\left(\text{TCE}_2^{\text{bin}} \right)^2 = \frac{1}{n} \sum_{i=1}^n h_{\text{bin}}(f(X_i), f(X_i)). \quad (24)$$

Similarly, one may formulate the debiased (but not unbiased) estimator of Kumar et al. (2019).

Further, we can also put the estimator of Popordanoska et al. (2022b) in the form of Equation (22) by defining

$$h_{\text{kde}}(p, p') := \left\langle p - \frac{\sum_{i=1}^n e_{Y_i} k_{\text{dir}}(f(X_i), p)}{\sum_{i=1}^n k_{\text{dir}}(f(X_i), p)}, p' - \frac{\sum_{i=1}^n e_{Y_i} k_{\text{dir}}(f(X_i), p')}{\sum_{i=1}^n k_{\text{dir}}(f(X_i), p')} \right\rangle, \quad (25)$$

where k_{dir} is the Dirichlet kernel as defined above in Equation (11). Again, it holds that

$$\left(\text{CCE}_2^{\text{kde}} \right)^2 = \frac{1}{n} \sum_{i=1}^n h_{\text{kde}}(f(X_i), f(X_i)). \quad (26)$$

We will also use an analogous estimator $\text{TCE}_2^{\text{kde}}$ for estimating TCE_2 in the experiment sections. Extending the binning based estimator to CCE_2 is in general not possible. Note that the runtime complexity of $\text{CCE}_2^{\text{kde}}$ is in $O(n^2 d)$, since we compare every evaluation instance with every training instance for every class. In the following we introduce novel estimators, which are runtime invariant to the number of classes. They are based on kernel ridge regression and directly minimize the empirical risk under kernel ridge regression assumptions.

4.2 Kernel Ridge Regression

Here, we propose two novel calibration estimators, which are derived as closed-form solutions under the typical kernel ridge regression assumptions, see, e.g., (Schölkopf & Smola, 2002; Bach, 2024). The following approach is based on the notion of ordinary Kronecker kernel ridge regression (Stock et al., 2018). Specifically, we require a reproducing kernel Hilbert space (RKHS) \mathcal{H} with an associated feature map $\phi: \Delta^d \rightarrow \mathcal{H}$, kernel

$k_{\mathcal{H}}$, inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$, and norm $\|\cdot\|_{\mathcal{H}}$. Denote with $\mathcal{H} \otimes \mathcal{H}$ the tensor product of the Hilbert space with itself, with respective feature map $(\phi \otimes \phi): \Delta^d \times \Delta^d \rightarrow \mathcal{H} \otimes \mathcal{H}$. Next, assume that $\langle f(X) - e_Y, f(X') - e_{Y'} \rangle = \langle g^*, (\phi \otimes \phi)(p, p') \rangle_{\mathcal{H} \otimes \mathcal{H}} + \epsilon$ for some $g^* \in \mathcal{H} \otimes \mathcal{H}$ and zero-mean noise term ϵ . Define the kernel ridge objective for a $g \in \mathcal{H} \otimes \mathcal{H}$ via

$$\hat{\mathcal{R}}_{\text{CE}, \lambda}(g) := \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left(\langle f(X_i) - e_{Y_i}, f(X_j) - e_{Y_j} \rangle - h(f(X_j), f(X_j)) \right)^2 + \lambda \|g\|_{\mathcal{H} \otimes \mathcal{H}}^2, \quad (27)$$

where $h(p, p') = \langle g, (\phi \otimes \phi)(p, p') \rangle_{\mathcal{H} \otimes \mathcal{H}}$. Then, a closed-form minimizer can be found, which results in the predictor

$$h_{\text{kkkr}}(p, p') := \text{vec}^\top \left(\Delta_{Y f(X)}^\top \Delta_{Y f(X)} \right) \left(\mathbf{K}_{f(X)} \otimes \mathbf{K}_{f(X)} + \lambda n^2 I \right)^{-1} \left(\mathbf{k}_{f(X)}(p) \otimes \mathbf{k}_{f(X)}(p') \right), \quad (28)$$

where \otimes becomes the Kronecker product, $\Delta_{Y f(X)} := (f(X_1) - e_{Y_1} \ \cdots \ f(X_n) - e_{Y_n}) \in \mathbb{R}^{d \times n}$, $\mathbf{K}_{f(X)} := (k_{\mathcal{H}}(f(X_i), f(X_j)))_{i,j \in \{1, \dots, n\}} \in \mathbb{R}^{n \times n}$, and $\mathbf{k}_{f(X)}(p) := (k_{\mathcal{H}}(f(X_i), p))_{i \in \{1, \dots, n\}} \in \mathbb{R}^n$. Without further modifications, computing Equation (28) has runtime complexity $O(n^6)$ due to the matrix inverse, which is practically infeasible. To reduce the complexity, we classically for Kronecker products make use of the eigenvalue decomposition $\mathbf{K}_{f(X)} = Q_{f(X)} \text{diag}(\lambda_1, \dots, \lambda_n) Q_{f(X)}^\top$, which is in $O(n^3)$. Define $\tilde{\Lambda}_{f(X)} \in \mathbb{R}^{n \times n}$ with $[\tilde{\Lambda}_{f(X)}]_{ij} = \frac{1}{\lambda_i \lambda_j + \lambda n^2}$ and denote with \odot the Hadamard product. It holds that

$$h_{\text{kkkr}}(p, p') = \mathbf{k}_{f(X)}^\top(p) Q_{f(X)} \left(\tilde{\Lambda}_{f(X)} \odot Q_{f(X)}^\top \Delta_{Y f(X)}^\top \Delta_{Y f(X)} Q_{f(X)} \right) Q_{f(X)}^\top \mathbf{k}_{f(X)}(p'). \quad (29)$$

This representation consists of multiplications of $n \times n$ matrices and in consequence is in $O(n^3)$. A more general approach uses the Schur decomposition to achieve such a reduction in complexity (Moravitz Martin & Van Loan, 2007). Naively computing the empirical generalization error $\hat{\mathcal{R}}_{\text{CE}}(h_{\text{kkkr}})$ on an evaluation set $(X'_1, Y'_1), \dots, (X'_{n'}, Y'_{n'})$ for some $n' \propto n$ has complexity $O(n^5)$, which is again prohibitive in practice. However, we can also reduce this complexity to $O(n^3)$ since it holds

$$\hat{\mathcal{R}}_{\text{CE}}(h_{\text{kkkr}}) = \frac{1}{n'(n'-1)} \sum_{i=1}^{n'} \sum_{\substack{j=1 \\ j \neq i}}^{n'} \left(\langle f(X_i) - e_{Y_i}, f(X_j) - e_{Y_j} \rangle - H_{ij}^{\text{kkkr}} \right)^2, \quad (30)$$

with

$$H^{\text{kkkr}} := \mathbf{K}_{f(X)f(X')}^\top Q_{f(X)} \left(\tilde{\Lambda}_{f(X)} \odot Q_{f(X)}^\top \Delta_{Y f(X)}^\top \Delta_{Y f(X)} Q_{f(X)} \right) Q_{f(X)}^\top \mathbf{K}_{f(X)f(X')} \in \mathbb{R}^{n' \times n'} \quad (31)$$

and $[\mathbf{K}_{f(X)f(X')}]_{ij} := k_{\mathcal{H}}(f(X_i), f(X'_j))$.

Alternatively, one may also fit a kernel ridge regressor for the problem assumption $f(X) - e_Y = \tilde{g}^* \phi(X) + \epsilon$ with $\tilde{g}^* \in \{\tilde{g}: \mathcal{H} \rightarrow \Delta^d\}$, and then use the result as plug-in for the inner product. This is referred to as two-step kernel ridge regression (Stock et al., 2018) or U-statistic regression (Park et al., 2021). The model then becomes

$$h_{\text{ukkr}}(p, p') := \mathbf{k}_{f(X)}^\top(p) \left(\mathbf{K}_{f(X)} + \lambda n I \right)^{-1} \Delta_{Y f(X)}^\top \Delta_{Y f(X)} \left(\mathbf{K}_{f(X)} + \lambda n I \right)^{-1} \mathbf{k}_{f(X)}(p'). \quad (32)$$

It holds that $h_{\text{ukkr}} = h_{\text{kkkr}}$ if the kernel used for h_{kkkr} incorporates the regularisation constant λ or when $\lambda = 0$ (Stock et al., 2018).

In the next section, we perform top-label confidence and canonical calibration evaluations with the discussed and proposed estimators. Specifically, we use the proposed calibration estimation risk in Equation (20) for comparison and the proposed calibration-evaluation pipeline of Section 3.2.2 for calibration estimation.

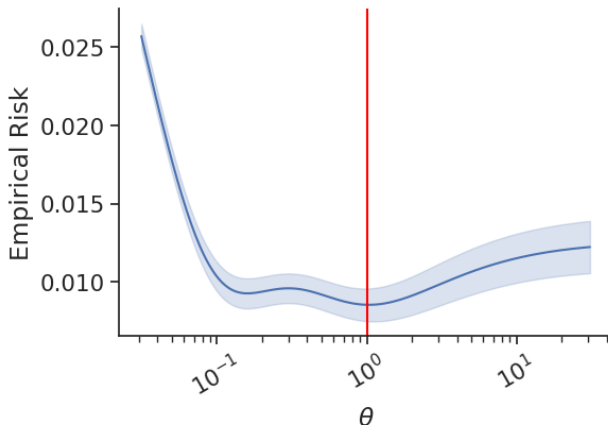


Figure 1: Simulated experiment for estimating CCE_2 in a task with 5 classes and 500 instances. The empirical risk correctly identifies the ideal calibration estimator with $\theta = 1$ indicated by the red line. Standard deviations across multiple seeds indicate the empirical risk stability.

Table 1: Validation set square root risk $\sqrt{\hat{\mathcal{R}}_{\text{CE}}} \times 100$ of TCE_2 estimators for CIFAR10 models with optimized hyperparameters. Lower is better. The estimator $\text{TCE}_2^{\text{kde}}$ performs worse and $\text{TCE}_2^{\text{ukkr}}$ better or equal than the other estimators. The large risk of $\text{TCE}_2^{\text{kde}}$ translates to an outlier calibration estimation in Figure 2a.

Model Estimator	LeNet-5	Densenet-40	ResNetWide-32	Resnet-110	Resnet-110 SD
$\text{TCE}_2^{15\text{-bins}}$	14.96 ± 0.31	6.14 ± 0.12	5.03 ± 0.2	5.4 ± 0.24	4.73 ± 0.25
$\text{TCE}_2^{\text{bins}}$	14.96 ± 0.31	6.12 ± 0.12	5.03 ± 0.2	5.39 ± 0.23	4.72 ± 0.25
$\text{TCE}_2^{\text{kde}}$	14.98 ± 0.31	13.7 ± 0.27	12.31 ± 0.65	7.46 ± 0.35	10.65 ± 0.58
$\text{TCE}_2^{\text{kkkr}}$	14.96 ± 0.31	6.13 ± 0.12	5.03 ± 0.2	5.39 ± 0.24	4.72 ± 0.25
$\text{TCE}_2^{\text{ukkr}}$	14.96 ± 0.31	6.11 ± 0.12	5.03 ± 0.2	5.38 ± 0.24	4.72 ± 0.25

5 Experiments

In this section, we demonstrate how to use our proposed risk framework in practice. We first run a simulation with known ground truth. Then, we evaluate the risk of the different calibration estimation function defined in Section 4 and the respective estimated calibration error across a variety of image classification datasets and models. We focus on top-label confidence, since it is the most prominent, and canonical calibration, which is the most general. The source code is publicly available at <https://github.com/waiting-for-acceptance>.

5.1 Simulation

We construct a simulation experiment with known ground truth to demonstrate how our proposed risk identifies the solution. For this, we draw i.i.d. samples $P_1, \dots, P_{500} \sim \text{Dir}(\alpha_1, \dots, \alpha_5)$ from a Dirichlet distribution with concentration parameters $\alpha_1 = \dots = \alpha_5 = 0.04$. These samples represent the (in practice unknown) ground truth probability vectors of a classification task with 500 instances and 5 classes. Then, we sample $Y_i \sim P_i$ for $i = 1 \dots 500$ as target labels. We set the hypothetical model predictions to $f(X_i) = \text{softmax}\left(\frac{3}{10} \log P_i\right)$ for $i = 1 \dots 500$, which represent miscalibrated predictions. The concentration parameters were set such that the model would have an accuracy of $\approx 90\%$. We then define a calibration estimation function $h_{\text{sim}}(p, p') := \left\langle p - \text{softmax}\left(\frac{10}{3} \theta \log p\right), p' - \text{softmax}\left(\frac{10}{3} \theta \log p'\right) \right\rangle$, which has a single learnable parameter $\theta \in \mathbb{R}$ referred to as temperature. The estimation function was chosen such that it matches the ground truth

Table 2: Validation set risk $\sqrt{\hat{\mathcal{R}}_{\text{CE}}} \times 100$ of TCE_2 for Cifar100 models. Lower is better. Similar as before, the estimator $\text{TCE}_2^{\text{kde}}$ performs worse than the other estimators except for LeNet-5. The best performing are $\text{TCE}_2^{\text{ukkr}}$ and $\text{TCE}_2^{\text{bins}}$.

Model	LeNet-5	Densenet-40	ResNetWide-32	Resnet-110	Resnet-110 SD
$\text{TCE}_2^{15\text{-bins}}$	18.51 ± 0.16	20.66 ± 0.40	18.40 ± 0.19	18.67 ± 0.43	17.18 ± 0.20
$\text{TCE}_2^{\text{bins}}$	18.50 ± 0.16	20.43 ± 0.38	18.24 ± 0.17	18.61 ± 0.42	17.11 ± 0.20
$\text{TCE}_2^{\text{kde}}$	18.50 ± 0.16	23.74 ± 0.41	23.28 ± 0.18	19.69 ± 0.47	18.21 ± 0.19
$\text{TCE}_2^{\text{kkkr}}$	18.50 ± 0.16	20.52 ± 0.39	18.29 ± 0.18	18.59 ± 0.42	17.11 ± 0.20
$\text{TCE}_2^{\text{ukkr}}$	18.50 ± 0.16	20.51 ± 0.40	18.28 ± 0.18	18.61 ± 0.43	17.12 ± 0.21

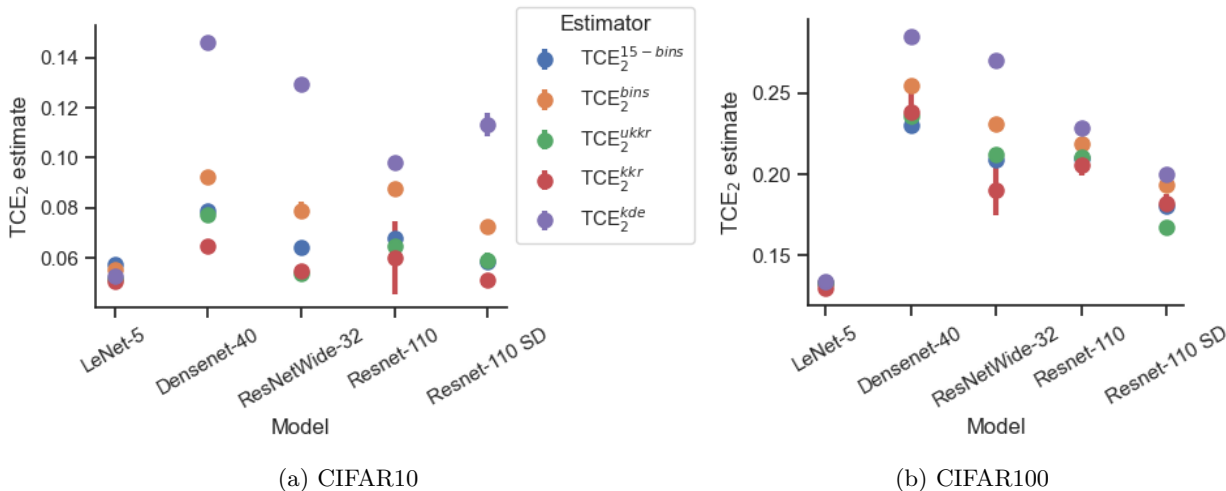


Figure 2: Different TCE_2 estimates of different models. Most calibration estimates approximately agree with each other. Only $\text{TCE}_2^{\text{kde}}$ is an outlier for Densenet-40, ResNetWide-32, and Resnet-110 SD. However, it also shows an increased calibration estimation risk in these cases (c.f. Table 1).

$h_{\text{sim}} = h^*$ if and only if $\theta = 1$. In Figure 1, we plot the mean results with standard deviations of the empirical risk according to 100 repetitions of the experiment. As can be seen, the empirical risk correctly identifies the correct calibration estimation function with $\theta = 1$.

5.2 Real World Settings

In the following, we evaluate and compare estimators discussed in Section 4 according to our novel calibration-evaluation pipeline introduced in Section 3.2.2 on standard image classifier setups, like CIFAR and ImageNet.

Technical Setup

The experiments are conducted across several model-dataset combinations, whose logit sets are openly accessible (Kull et al., 2019; Rahimi et al., 2020; Gruber & Buettner, 2022).¹ The image classification datasets in use are CIFAR10 with 10 classes, CIFAR100 with 100 classes (Krizhevsky, 2009), and ImageNet with 1,000 classes (Deng et al., 2009). Since we restrict ourselves to evaluating the calibration error estimate of the models, we only use the test set of each dataset (CIFAR: $n = 10,000$, ImageNet: $n = 25,000$). Modifying or selecting models based on the calibration estimate would require using the validation set instead. The

¹https://github.com/ML0-lab/better_uncertainty_calibration

Table 3: Validation set square root risk $\sqrt{\hat{\mathcal{R}}_{\text{CE}}} \times 100$ of CCE_2 estimators for CIFAR100 models. Again, lower is better. Contrary to previous results, the estimator $\text{CCE}_2^{\text{kde}}$ manages to outperform the kernel ridge regression based estimators in some scenarios.

Model	LeNet-5	Densenet-40	ResNetWide-32	Resnet-110	Resnet-110 SD
$\text{CCE}_2^{\text{kde}}$	8.38 ± 0.06	5.61 ± 0.09	4.98 ± 0.05	5.14 ± 0.1	4.79 ± 0.03
$\text{CCE}_2^{\text{kkkr}}$	8.36 ± 0.06	5.56 ± 0.09	5.00 ± 0.05	5.16 ± 0.1	4.80 ± 0.03
$\text{CCE}_2^{\text{ukkr}}$	8.35 ± 0.06	5.56 ± 0.09	5.01 ± 0.05	5.16 ± 0.1	4.80 ± 0.02

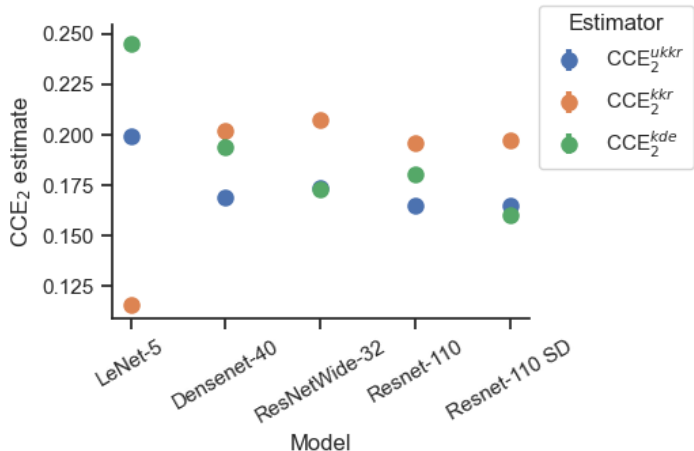


Figure 3: Different CCE_2 estimates for CIFAR100 models. The risk values of Table 3 do not relate to the calibration estimate but only indicate which estimator to trust more (here: $\text{CCE}_2^{\text{kde}}$).

included models are LeNet-5 (LeCun et al., 1998), ResNet-110, ResNet-110 SD, ResNet-152 (He et al., 2016), Wide ResNet-32 (Zagoruyko & Komodakis, 2016), DenseNet-40, DenseNet-161 (Huang et al., 2017), and PNASNet-5 Large (Liu et al., 2018). We did not conduct model training ourselves and refer to Kull et al. (2019) and Rahimi et al. (2020) for further details. We evaluate top-label confidence calibration and canonical calibration estimators for these classifiers.

We run the calibration-evaluation pipeline proposed in Section 3.2.2 with a random split of the original test set, using 80% for tuning the calibration estimator function via cross-validation and 20% for the calibration test set \mathcal{D}_{te} , which computes the mean in Equation (21). In all experiments, we use 5-fold cross-validation to optimize the hyperparameters of a calibration estimator function. For the best performing hyperparameter, the models across all folds are used as an ensemble predictor for the final calibration estimation on the set \mathcal{D}_{te} . This approach also allows us to include error bars according to the cross validation folds.

As calibration estimator functions, we consider h_{bin} , h_{kde} , h_{kkkr} , and h_{ukkr} for top-label confidence, as well as h_{kde} , h_{kkkr} , and h_{ukkr} canonical calibration. For both kernel ridge regression models, we use the RKHS of the RBF kernel $k_{\text{rbf}}(x, y) = \exp(-\gamma \|x - y\|^2)$. We set $\gamma = \frac{1}{2}$ based on preliminary evaluations. We also evaluate h_{bin} with 15 bins without hyperparameter optimization, which we refer to as $\text{TCE}_2^{15\text{-bins}}$. This corresponds to a common default choice in current practice (Guo et al., 2017; Detlefsen et al., 2022). More details on the hyperparameter search spaces are given in Appendix B.

Results

We now discuss the experimental results. All reported risks are with respect to the holdout sets in the cross-validation folds. The reported calibration estimations are with respect to the calibration test set (20% of the original test set). The error bars for the risk and calibration estimations are the standard errors according to the cross-validation folds. In Table 1 we show the performance across different models of CIFAR10, and in Table 2 across models of CIFAR100. Here, we compare the calibration estimation functions for top-label confidence calibration. As can be seen, no calibration estimation function dominates all others. Specifically, $\text{TCE}_2^{\text{kde}}$ performs worst for all models, even when we consider the error bars. The estimator $\text{TCE}_2^{\text{ukkr}}$ outperforms the other estimators, however, the difference is too marginal with respect to the error bars to come to a confident conclusion. For the CIFAR100 models, $\text{TCE}_2^{\text{kde}}$ performs more similar to the other estimators. Further, the optimized binning estimator $\text{TCE}_2^{\text{bins}}$ shows the strongest performance and not $\text{TCE}_2^{\text{ukkr}}$ anymore. However, again, the error bars are too large to designate a definitive ranking. Further, for the LeNet-5 model in CIFAR10 and CIFAR100, our risk is not sufficiently sensitive to rank the estimators.

In Figure 2 we depict the corresponding TCE_2 estimations. As can be seen, only $\text{TCE}_2^{\text{kde}}$ is occasionally an outlier relative to the other estimators, which is expected based on the reported risk values of Table 1 and Table 2. Even though, we cannot spot a direct connection between all risk values and the estimated calibration values in Figure 2, the large risk of $\text{TCE}_2^{\text{kde}}$ is indicative of a worse estimation. However, it is not surprising that differences in the risk do not always translate to differences in the estimated values, since the loss does not measure in which direction a calibration estimator function gives wrong predictions.

In Table 3, we report the risk values of the canonical calibration estimator functions for CIFAR100 classifiers. As can be seen, $\text{CCE}_2^{\text{kde}}$ performs better than in previous results, outperforming the other approaches in some cases. We can also see that $\text{CCE}_2^{\text{ukkr}}$ performs better than $\text{CCE}_2^{\text{kkkr}}$, which is a continuous trend across all results. However, the error bars dominate the performance difference and no clear cut conclusion can be made. In Figure 3, we show the respective calibration estimates.

In Appendix B, we offer additional results regarding top-label confidence calibration for ImageNet classifiers and canonical calibration for CIFAR10 classifiers. In summary, no calibration estimator outperforms the other approaches across all settings. Additionally, risk performance is often indicative of outlying calibration estimates. This underlines the requirement of a risk to assess which estimator to use for evaluating the calibration of a new model in practice. We may expect to find better estimators by extending the search space (e.g., by considering different kernels), or by including other model classes, like boosted trees or neural networks (Bishop & Nasrabadi, 2006). However, the proposed risk may not be sufficiently sensitive to rank the estimators according to their performance. Future research may involve exploring alternative loss functions for more sensitive results.

6 Conclusion

In this work, we introduced a mean-squared error based risk to compare different calibration estimators. This is the first approach in the literature to compare different calibration estimators on real-world datasets. We offer measure theoretic conditions for when the risk identifies an ideal estimator. We also derive novel calibration estimators as closed-form minimizers of the empirical risk based on kernel ridge regression assumptions. Further, using an empirical risk enables to perform hyperparameter optimization and estimator selection via a training-validation-testing pipeline, similar to conventional machine learning. In the experiments, we optimize the hyperparameters of common calibration estimators in the literature on popular real-world benchmarks, and compare the risks of different optimized estimators. No dominating calibration estimator was found, which emphasises the requirement of using our risk to detect an appropriate estimator for new settings in practice.

References

Francis Bach. *Learning Theory from First Principles*. MIT Press, 2024.

- Herman J. Bierens. *Topics in Advanced Econometrics: Estimation, Testing, and Specification of Cross-section and Time Series Models*. Cambridge University Press, 1996.
- Christopher M. Bishop and Nasser M. Nasrabadi. *Pattern Recognition and Machine Learning*, volume 4. Springer, 2006.
- Glenn W. Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1):1–3, 1950.
- Marek Capiński and Peter Ekkehard Kopp. *Measure, Integral and Probability*, volume 14. Springer, 2004.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45, 2024.
- Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, and Ibrahim Alabdulmohsin. Scaling vision transformers to 22 billion parameters. In *International Conference on Machine Learning*, pp. 7480–7512, 2023.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009.
- Nicki Skaftø Detlefsen, Jiri Borovec, Justus Schock, Ananya Harsh Jha, Teddy Koker, Luca Di Liello, Daniel Stancl, Changsheng Quan, Maxim Grechkin, and William Falcon. Torchmetrics - measuring reproducibility in pytorch. *Journal of Open Source Software*, 7(70):4101, 2022.
- Bradley Efron. *An Introduction to the Bootstrap*. CRC press, 1994.
- Hongxiang Fan, Martin Ferianc, Zhiqiang Que, Xinyu Niu, Miguel L. Rodrigues, and Wayne Luk. Accelerating Bayesian neural networks via algorithmic and hardware optimizations. *IEEE Transactions on Parallel and Distributed Systems*, 33(12):3387–3399, 2022.
- Di Feng, Christian Haase-Schütz, Lars Rosenbaum, Heinz Hertlein, Claudius Glaeser, Fabian Timm, Werner Wiesbeck, and Klaus Dietmayer. Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges. *IEEE Transactions on Intelligent Transportation Systems*, 22(3):1341–1360, 2020.
- Halina Frydman, Edward I. Altman, and Duen-Li Kao. Introducing recursive partitioning for financial classification: the case of financial distress. *The Journal of Finance*, 40(1):269–291, 1985.
- Tilmann Gneiting and Adrian E. Raftery. Weather forecasting with ensemble methods. *Science*, 310:248–249, 2005.
- 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.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- Sebastian G. Gruber and Florian Buettner. Better uncertainty calibration via proper scores for classification and beyond. In *Advances in Neural Information Processing Systems*, 2022.
- Sebastian G. Gruber, Teodora Popordanoska, Aleksei Tiulpin, Florian Buettner, and Matthew B. Blaschko. Consistent and asymptotically unbiased estimation of proper calibration errors. In *International Conference on Artificial Intelligence and Statistics*, pp. 3466–3474, 2024.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, pp. 1321–1330, 2017.
- Ashim Gupta, Giorgi Kvernadze, and Vivek Srikumar. Bert & family eat word salad: Experiments with text understanding. *ArXiv*, abs/2101.03453, 2021.

- Chirag Gupta and Aaditya Ramdas. Top-label calibration and multiclass-to-binary reductions. In *International Conference on Learning Representations*, 2022.
- Sarah Haggemüller, Roman C. Maron, Achim Hekler, Jochen S. Utikal, Catarina Barata, Raymond L. Barnhill, Helmut Beltraminelli, Carola Berking, Brigid Betz-Stablein, Andreas Blum, Stephan A. Braun, Richard Carr, Marc Combalia, Maria-Teresa Fernandez-Figueras, Gerardo Ferrara, Sylvie Fraitag, Lars E. French, Frank F. Gellrich, Kamran Ghoreshi, Matthias Goebeler, Pascale Guitera, Holger A. Haenssle, Sebastian Haferkamp, Lucie Heinzerling, Markus V. Heppt, Franz J. Hilke, Sarah Hobelsberger, Dieter Krahl, Heinz Kutzner, Aimilios Lallas, Konstantinos Liopyris, Mar Llamas-Velasco, Josep Malvehy, Friedegund Meier, Cornelia S.L. Müller, Alexander A. Navarini, Cristián Navarrete-Dechent, Antonio Perasole, Gabriela Poch, Sebastian Podlipnik, Luis Requena, Veronica M. Rotemberg, Andrea Saggini, Omar P. Sanguenza, Carlos Santonja, Dirk Schadendorf, Bastian Schilling, Max Schlaak, Justin G. Schlager, Mildred Sergon, Wiebke Sondermann, H. Peter Soyer, Hans Starz, Wilhelm Stolz, Esmeralda Vale, Wolfgang Weyers, Alexander Zink, Eva Krieghoff-Henning, Jakob N. Kather, Christof von Kalle, Daniel B. Lipka, Stefan Fröhling, Axel Hauschild, Harald Kittler, and Titus J. Brinker. Skin cancer classification via convolutional neural networks: systematic review of studies involving human experts. *European Journal of Cancer*, 156: 202–216, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Achim Hekler, Titus J. Brinker, and Florian Buettner. Test time augmentation meets post-hoc calibration: uncertainty quantification under real-world conditions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 14856–14864, 2023.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *Conference on Computer Vision and Pattern Recognition*, pp. 4700–4708, 2017.
- Ashraf Islam, Chun-Fu Chen, Rameswar Panda, Leonid Karlinsky, Richard J. Radke, and Rogério Schmidt Feris. A broad study on the transferability of visual representations with contrastive learning. *International Conference on Computer Vision (ICCV)*, pp. 8825–8835, 2021.
- Taejong Joo, Uijung Chung, and Minji Seo. Being Bayesian about categorical probability. In *ICML*, 2020.
- Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being Bayesian, even just a bit, fixes overconfidence in relu networks. In *ICML*, 2020.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Master’s thesis, University of Toronto, 2009.
- Volodymyr Kuleshov and Shachi Deshpande. Calibrated and sharp uncertainties in deep learning via density estimation. In *International Conference on Machine Learning*, pp. 11683–11693, 2022.
- Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, 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*, 32:12316–12326, 2019.
- Ananya Kumar, Percy Liang, and Tengyu Ma. Verified uncertainty calibration. In *Advances on Neural Information Processing Systems*, pp. 3792–3803, 2019.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 19–34, 2018.
- Charlie Marx, Sofian Zalouk, and Stefano Ermon. Calibration by distribution matching: Trainable kernel calibration metrics. *Advances in Neural Information Processing Systems*, 36, 2024.

- Aditya Krishna Menon, Ankit Singh Rawat, Sashank J. Reddi, Seungyeon Kim, and Sanjiv Kumar. A statistical perspective on distillation. In *ICML*, 2021.
- Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- Pablo Morales-Álvarez, Daniel Hernández-Lobato, Rafael Molina, and José Miguel Hernández-Lobato. Activation-level uncertainty in deep neural networks. In *ICLR*, 2021.
- Carla D. Moravitz Martin and Charles F. Van Loan. Shifted kronecker product systems. *SIAM Journal on Matrix Analysis and Applications*, 29(1):184–198, 2007.
- Allan H. Murphy. A new vector partition of the probability score. *Journal of Applied Meteorology and Climatology*, 12(4):595 – 600, 1973.
- Allan H. Murphy and Robert L. Winkler. Reliability of subjective probability forecasts of precipitation and temperature. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 26(1):41–47, 1977.
- Kevin P. Murphy. *Probabilistic Machine Learning: An Introduction*. MIT Press, 2022.
- Mahdi Pakdaman Naeni, Gregory F. Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using Bayesian binning. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pp. 2901–2907, 2015.
- Jeremy Nixon, Michael W. Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran. Measuring calibration in deep learning. In *CVPR Workshops*, volume 2, 2019.
- Frédéric Ouimet and Raimon Tolosana-Delgado. Asymptotic properties of dirichlet kernel density estimators. *Journal of Multivariate Analysis*, 187:104832, 2022.
- Junhyung Park, Uri Shalit, Bernhard Schölkopf, and Krikamol Muandet. Conditional distributional treatment effect with kernel conditional mean embeddings and u-statistic regression. In *International Conference on Machine Learning*, pp. 8401–8412, 2021.
- Kanil Patel, William H. Beluch, Bin Yang, Michael Pfeiffer, and Dan Zhang. Multi-class uncertainty calibration via mutual information maximization-based binning. In *International Conference on Learning Representations*, 2021.
- Teodora Popordanoska, Raphael Sayer, and Matthew B. Blaschko. A consistent and differentiable L_p canonical calibration error estimator. In *Advances in Neural Information Processing Systems*, 2022a.
- Teodora Popordanoska, Raphael Sayer, and Matthew B. Blaschko. A consistent and differentiable L_p canonical calibration error estimator. In *Advances in Neural Information Processing Systems*, 2022b.
- Amir Rahimi, Amirreza Shaban, Ching-An Cheng, Richard Hartley, and Byron Boots. Intra order-preserving functions for calibration of multi-class neural networks. *Advances in Neural Information Processing Systems*, 33:13456–13467, 2020.
- Rebecca Roelofs, Nicholas Cain, Jonathon Shlens, and Michael C. Mozer. Mitigating bias in calibration error estimation. In *International Conference on Artificial Intelligence and Statistics*, pp. 4036–4054, 2022.
- Bernhard Schölkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2002.
- Jun Shao. *Mathematical statistics*. Springer Science & Business Media, 2003.
- Michiel Stock, Tapio Pahikkala, Antti Airola, Bernard De Baets, and Willem Waegeman. A comparative study of pairwise learning methods based on kernel ridge regression. *Neural Computation*, 30(8):2245–2283, 2018.

- Zeyu Sun, Dogyoon Song, and Alfred Hero. Minimum-risk recalibration of classifiers. *Advances in Neural Information Processing Systems*, 36, 2024.
- Gerald Teschl. *Mathematical Methods in Quantum Mechanics*, volume 157. American Mathematical Soc., 2014.
- Junjiao Tian, Dylan Yung, Yen-Chang Hsu, and Zsolt Kira. A geometric perspective towards neural calibration via sensitivity decomposition. In *NeurIPS*, 2021.
- Christian Tomani, Sebastian G. Gruber, Muhammed Ebrar Erdem, Daniel Cremers, and Florian Buettner. Post-hoc uncertainty calibration for domain drift scenarios. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10124–10132, June 2021.
- 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*, pp. 3459–3467, 2019.
- Xiao Wang, Hongrui Liu, Chuan Shi, and Cheng Yang. Be confident! towards trustworthy graph neural networks via confidence calibration. In *NeurIPS*, 2021.
- David Widmann, Fredrik Lindsten, and Dave Zachariah. Calibration tests in multi-class classification: A unifying framework. *Advances in Neural Information Processing Systems*, 32:12257–12267, 2019.
- David Widmann, Fredrik Lindsten, and Dave Zachariah. Calibration tests beyond classification. In *International Conference on Learning Representations*, 2021.
- Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass probability estimates. In *International Conference on Knowledge Discovery and Data Mining*, pp. 694–699. Association for Computing Machinery, 2002.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision Conference*, 2016.

Table 4: Validation set square root risk $\sqrt{\hat{\mathcal{R}}_{\text{CE}}} \times 100$ of TCE_2 for different ImageNet models. Lower is better. The estimator $\text{TCE}_2^{\text{kde}}$ performs worse than the other estimators for DenseNet-161 and ResNet-152. For Pnasnet-5, all estimators perform similar.

Model Estimator	DenseNet-161	Resnet-152	Pnasnet-5
$\text{TCE}_2^{15\text{-bins}}$	12.17 ± 0.16	12.58 ± 0.14	10.63 ± 0.16
$\text{TCE}_2^{\text{bins}}$	12.17 ± 0.16	12.57 ± 0.14	10.63 ± 0.16
$\text{TCE}_2^{\text{kde}}$	12.31 ± 0.17	12.77 ± 0.14	10.63 ± 0.16
$\text{TCE}_2^{\text{kkkr}}$	12.17 ± 0.16	12.57 ± 0.14	10.63 ± 0.16
$\text{TCE}_2^{\text{ukkr}}$	12.17 ± 0.16	12.57 ± 0.14	10.63 ± 0.16

A Overview

Here, we first give additional experimental results in Appendix B. The missing proofs are located in Appendix C.

B Extended Experiments

Additional details

We use the implementation of the calibration estimator function h_{kde} given by the original authors (Popordanoska et al., 2022b). For a small fraction of inputs, this implementation returns NaN as prediction. We remove these instances from the risk and calibration estimation calculation of h_{kde} , which has a neglectable effect.

As hyperparameter search spaces for the TCE experiments, we consider $\{5i \mid i = 1, \dots, 20\}$ for the number of bins in h_{bin} , a bandwidth in $\{10^{-5(i-1)/14-(1-(i-1)/14)} \mid i = 1, \dots, 15\} \cup \{0.2i \mid i = 1, \dots, 5\}$ for the Dirichlet kernel of h_{kde} according to Popordanoska et al. (2022a), a regularization constant $\lambda \in \{n^{0.5}10^{-2i+1} \mid i = 1, \dots, 9\}$ for h_{kkkr} , and $\lambda \in \{n^{0.5}10^{-i} \mid i = 1, \dots, 9\}$ for h_{ukkr} . For the CCE experiments, we consider the same set of bandwidths for the Dirichlet kernel of h_{kde} , a regularization constant $\lambda \in \{n^{0.5}10^{-i+9} \mid i = 1, \dots, 18\}$ for h_{kkkr} , and $\lambda \in \{n^{0.5}10^{-0.5i+4.5} \mid i = 1, \dots, 18\}$ for h_{ukkr} .

Additional results

In the following, we discuss the risks and calibration estimations of some left-out cases from the main paper.

In Table 4 we show the risk of the top-label confidence calibration estimators for ImageNet with various models. All calibration estimation functions show similar risk except $\text{TCE}_2^{\text{kde}}$, which is worse for DenseNet-161 and Resnet-152. This is in agreement with Figure 2b, where the estimated calibration values are also mostly similar. The risks in Table 5 for canonical calibration estimators in the case of CIFAR10 show slightly different results: Here, the risk fails to distinguish the performance between the different estimators. Only $\text{CCE}_2^{\text{kde}}$ outperforms the other approaches for Resnet-110.

In summary, the results mimic the ones in the main paper and it is not apparent which estimator to use in practice without considering our proposed risk. However, the risk may be insensitive regarding the various estimator performances.

Table 5: Validation set risk $\sqrt{\hat{\mathcal{R}}_{\text{CE}}} \times 100$ of CCE_2 for CIFAR10 models. Lower is better. All estimators perform fairly similar.

Model	LeNet-5	Densenet-40	ResNetWide-32	Resnet-110	Resnet-110 SD
CCE_2^{kde}	13.53 ± 0.27	5.13 ± 0.13	4.22 ± 0.16	4.46 ± 0.14	3.89 ± 0.2
CCE_2^{kkr}	13.53 ± 0.27	5.13 ± 0.13	4.22 ± 0.16	4.47 ± 0.14	3.89 ± 0.2
CCE_2^{ukkr}	13.53 ± 0.27	5.13 ± 0.13	4.22 ± 0.16	4.47 ± 0.14	3.89 ± 0.2

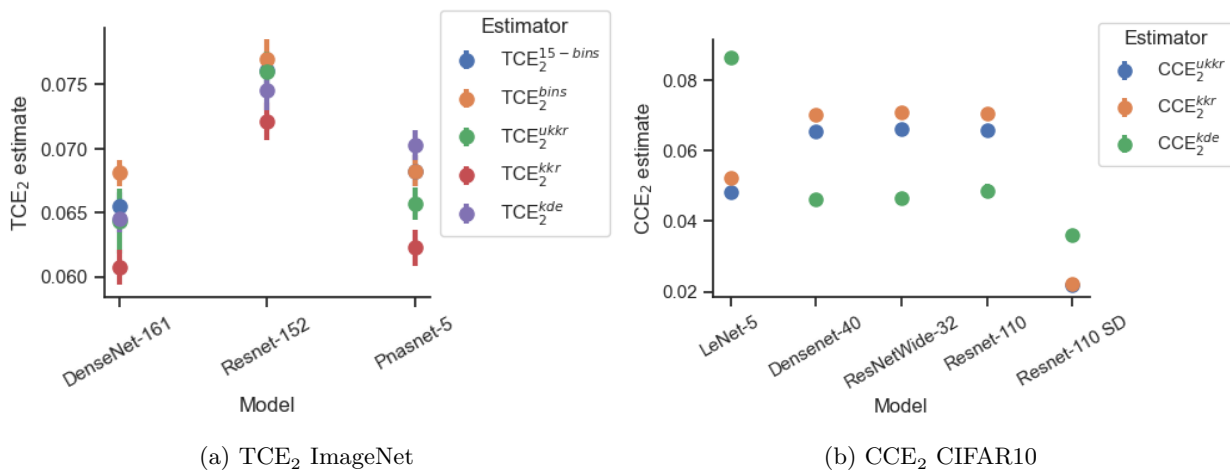


Figure 4: Different calibration estimates of different models. Most calibration estimates approximately agree with each other. This is in agreement with the similar risk values for each estimator in Table 4 and Table 5.

C Missing Proofs

Here, we present the missing proofs of the main part. Specifically, we prove Theorem 1 in Section C.1, Theorem 2 in Section C.2, and various statements of Section 3.2 in Section C.3.

C.1 Proof for Theorem 1

We show that $\mathcal{R}_{\text{CE}}(h) > \mathcal{R}_{\text{CE}}(h^*)$.

For this, we require that $\langle p - \mathbb{P}_Y, p' - \mathbb{P}_V \rangle$ is the unique minimizer of $L_{\text{CE}}(\cdot, \mathbb{P}_Y \otimes \mathbb{P}_V; p, p')$, which holds since

$$\begin{aligned}
 & \frac{\partial}{\partial c} L_{\text{CE}}(c, \mathbb{P}_Y \otimes \mathbb{P}_V; p, p') \\
 &= \frac{\partial}{\partial c} \mathbb{E}_{Y,V} [(c - \langle p - e_Y, p' - e_V \rangle)^2] \\
 &= 2 \mathbb{E}_{Y,V} [(c - \langle p - e_Y, p' - e_V \rangle)] \\
 &= 2(c - \langle p - \mathbb{P}_Y, p' - \mathbb{P}_V \rangle),
 \end{aligned} \tag{33}$$

and $\frac{\partial^2}{\partial c^2} L_{\text{CE}}(c, \mathbb{P}_Y \otimes \mathbb{P}_V; p, p') > 0$.

Based on the assumption that $\exists A \in \mathcal{F}_{f(X)}$ with $\mathbb{P}_{f(X)}(A) > 0$ we have

$$\begin{aligned}
& \forall p, p' \in A: h(p, p') \neq h^*(p, p') \\
& \iff \forall p, p' \in A: L_{\text{CE}}(h(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') > L_{\text{CE}}(h^*(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') \\
& \implies \int_{A \times A} L_{\text{CE}}(h(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') d(\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)})(p, p') \\
& \quad > \int_{A \times A} L_{\text{CE}}(h^*(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') d(\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)})(p, p'),
\end{aligned} \tag{34}$$

where the inequality follows since $h^*(p, p') = \langle p - \mathbb{P}_{Y|f(X)=p}, p' - \mathbb{P}_{Y|f(X)=p'} \rangle$ is the unique minimizer of $L_{\text{CE}}(\cdot, \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p')$.

From the unique minimizer property also follows that for all $B \in \mathcal{F}_{f(X)} \otimes \mathcal{F}_{f(X)}$ holds

$$\begin{aligned}
& \int_B L_{\text{CE}}(h(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') d(\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)})(p, p') \\
& \quad \geq \int_B L_{\text{CE}}(h^*(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') d(\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)})(p, p').
\end{aligned} \tag{35}$$

Since $(\Delta^d \times \Delta^d) \setminus (A \times A) \in \mathcal{F}_{f(X)} \otimes \mathcal{F}_{f(X)}$, it holds

$$\begin{aligned}
& \mathcal{R}_{\text{CE}}(h) \\
& = \mathbb{E}_{X, X'} [L_{\text{CE}}(h(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p')] \\
& = \int_{(\Delta^d \times \Delta^d) \setminus (A \times A)} L_{\text{CE}}(h(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') d(\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)})(p, p') \\
& \quad + \int_{A \times A} L_{\text{CE}}(h(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') d(\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)})(p, p') \\
& > \int_{(\Delta^d \times \Delta^d) \setminus (A \times A)} L_{\text{CE}}(h^*(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') d(\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)})(p, p') \\
& \quad + \int_{A \times A} L_{\text{CE}}(h^*(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p') d(\mathbb{P}_{f(X)} \otimes \mathbb{P}_{f(X)})(p, p') \\
& = \mathbb{E}_{X, X'} [L_{\text{CE}}(h^*(p, p'), \mathbb{P}_{Y|f(X)=p} \otimes \mathbb{P}_{Y|f(X)=p'}; p, p')] \\
& = \mathcal{R}_{\text{CE}}(h^*).
\end{aligned} \tag{36}$$

C.2 Proof for Theorem 2

A sketch of the necessity of Theorem 2 is given in Figure 5, which also illustrates the proof.

Proof. We use $P := f(X)$ and $f(p, p') := \begin{cases} (h(p, p') - h^*(p, p'))^2, & p, p' \in \mathcal{P}_Y \\ 0, & \text{else,} \end{cases}$ for simplicity. It is continuous

at every point in which h and h^* are continuous. We denote with D the \mathbb{P}_P -null set of p 's for which h and h^* are not continuous at point (p, p) . Then,

$$\begin{aligned}
& \mathbb{E}[f(P, P)] \\
& = \int_{\mathbb{R}^d} f(p, p) d\mathbb{P}_P(p) \\
& = \int_{\text{supp}(\mathbb{P}_P)} f(p, p) d\mathbb{P}_P(p) \\
& = \int_{\text{int supp}(\mathbb{P}_P)} f(p, p) d\mathbb{P}_P(p) + \int_{\text{bd supp}(\mathbb{P}_P)} f(p, p) d\mathbb{P}_P(p) \\
& = \int_{\text{int supp}(\mathbb{P}_P) \setminus D} f(p, p) d\mathbb{P}_P(p) + \int_{\text{bd supp}(\mathbb{P}_P)} f(p, p) d\mathbb{P}_P(p).
\end{aligned} \tag{37}$$

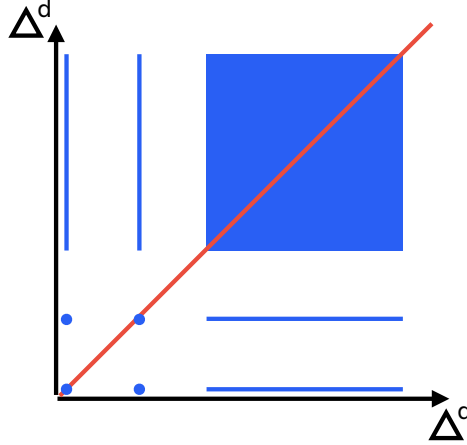


Figure 5: Blue indicates a possible support set $\text{supp}(\mathbb{P}_P \otimes \mathbb{P}_P) \subseteq \Delta^d \times \Delta^d$, and red line indicates $\{(p, p) \mid p \in \Delta^d\}$. During training we optimize all blue dots and areas, but during testing we only evaluate their intersection with the red line (a possible null set).

The last line holds since the support of a measure is closed, and, consequently, splitting it up into interior and boundary does not add new 'mass' (Teschl, 2014). For the following, denote with $B(p, \epsilon) := \{x \in \mathbb{R}^d \mid \|x - p\|_2 < \epsilon\}$ the open (euclidean) ball with center p and radius ϵ . Further, since h and h^* are by assumption continuous in the points $\{(p, p) \mid p \in \mathcal{P}_Y \setminus D\}$, it follows that f is also continuous at these points, and, consequently, also lower semicontinuous. We first deal with the interior term by using this lower-semicontinuous property giving

$$\begin{aligned}
& \int_{\text{int supp}(\mathbb{P}_P) \setminus D} f(p, p) d\mathbb{P}_P(p) \\
&= \int_{\text{int supp}(\mathbb{P}_P) \setminus D} \liminf_{(p_1, p_2) \rightarrow (p, p)} f(p_1, p_2) d\mathbb{P}_P(p) \\
&= \lim_{\epsilon \rightarrow 0} \int_{\text{int supp}(\mathbb{P}_P) \setminus D} \inf \{f(p_1, p_2) \mid (p_1, p_2) \in B((p, p), \epsilon) \setminus \{(p, p)\}\} d\mathbb{P}_P(p) \\
&\leq \lim_{\epsilon \rightarrow 0} \int_{\text{int supp}(\mathbb{P}_P) \setminus D} \text{ess inf} \{f(p_1, p_2) \mid (p_1, p_2) \in B((p, p), \epsilon) \setminus \{(p, p)\}\} d\mathbb{P}_P(p).
\end{aligned} \tag{38}$$

Since $p \in \text{int supp}(\mathbb{P}_P)$, it holds $\mathbb{P}_P(B(p, \epsilon/2) \setminus \{p\}) > 0$ for all $\epsilon > 0$ following the definition of int and supp (Teschl, 2014). Let $\text{Sq}(p, \epsilon) = \{x \in \mathbb{R}^d \mid \|p - x\|_\infty < \epsilon\}$ be the open hypercube with center p and side length 2ϵ . It holds $B(p, \epsilon) \subset \text{Sq}(p, \epsilon)$ and $\text{Sq}(p, \sqrt{d}\epsilon) \subset B(p, \epsilon)$. Then

$$\begin{aligned}
0 &< \mathbb{P}_P(B(p, \sqrt{2d}\epsilon) \setminus \{p\}) \\
&\leq \mathbb{P}_P(\text{Sq}(p, \sqrt{2d}\epsilon) \setminus \{p\}) \\
&= \sqrt{\mathbb{P}_P \otimes \mathbb{P}_P((\text{Sq}(p, \sqrt{2d}\epsilon) \setminus \{p\}) \times (\text{Sq}(p, \sqrt{2d}\epsilon) \setminus \{p\}))} \\
&\leq \sqrt{\mathbb{P}_P \otimes \mathbb{P}_P((\text{Sq}(p, \sqrt{2d}\epsilon) \times \text{Sq}(p, \sqrt{2d}\epsilon)) \setminus \{(p, p)\})} \\
&= \sqrt{\mathbb{P}_P \otimes \mathbb{P}_P(\text{Sq}((p, p), \sqrt{2d}\epsilon) \setminus \{(p, p)\})} \\
&\leq \sqrt{\mathbb{P}_P \otimes \mathbb{P}_P(B((p, p), \epsilon) \setminus \{(p, p)\})}.
\end{aligned} \tag{39}$$

Consequently, $B((p, p), \epsilon) \setminus \{(p, p)\}$ is not a null set (w.r.t. $\mathbb{P}_P \otimes \mathbb{P}_P$), and, thus,

$$\begin{aligned}
& \text{ess inf } \{f(p_1, p_2) \mid (p_1, p_2) \in B((p, p), \epsilon) \setminus \{(p, p)\}\} \\
& \leq \int_{B((p, p), \epsilon) \setminus \{(p, p)\}} f(p_1, p_2) \, d(\mathbb{P}_P \otimes \mathbb{P}_P)(p_1, p_2) \\
& \leq \int_{B((p, p), \epsilon)} f(p_1, p_2) \, d(\mathbb{P}_P \otimes \mathbb{P}_P)(p_1, p_2) \\
& \leq \mathbb{E}[f(P, P')] = 0,
\end{aligned} \tag{40}$$

where we used the given assumption $h(P, P') \stackrel{a.s.}{=} h^*(P, P')$ in the last line. Continuing Equation (38) it follows

$$\lim_{\epsilon \rightarrow 0} \int_{\text{int supp}(\mathbb{P}_P) \setminus D} \underbrace{\text{ess inf } \{f(p_1, p_2) \mid (p_1, p_2) \in B((p, p), \epsilon) \setminus \{(p, p)\}\}}_{=0} \, d\mathbb{P}_P(p) = 0. \tag{41}$$

Next, we deal with the boundary of the support. Since we assume that it consists of at most countably infinite elements with probability mass (which we denote as $\{p_1, \dots\}$), it holds

$$\begin{aligned}
& \int_{\text{bd supp}(\mathbb{P}_P)} f(p, p) \, d\mathbb{P}_P(p) \\
& = \int_{\{p_1, \dots\}} f(p, p) \, d\mathbb{P}_P(p) \\
& = \sum_{p \in \{p_1, \dots\}} f(p, p) \mathbb{P}_P(p) \\
& = \sum_{p \in \{p_1, \dots\}} \sum_{p' \in \{p_1, \dots\}} \mathbf{1}_{p=p'} f(p, p') \sqrt{\mathbb{P}_P(p)} \sqrt{\mathbb{P}_P(p')} \\
& = \int_{\{p_1, \dots\} \times \{p_1, \dots\}} \mathbf{1}_{p=p'} f(p, p') \, d(\sqrt{\mathbb{P}_P} \otimes \sqrt{\mathbb{P}_P})(p, p') \\
& = \int_{\{(p, p) \mid p \in \{p_1, \dots\}\}} f(p, p') \, d(\sqrt{\mathbb{P}_P} \otimes \sqrt{\mathbb{P}_P})(p, p') \\
& \leq \int_{\mathbb{R}^d} f(p, p') \, d(\sqrt{\mathbb{P}_P} \otimes \sqrt{\mathbb{P}_P})(p, p') \\
& = 0,
\end{aligned} \tag{42}$$

where the last line holds since for all $A \in \mathcal{F}_{P \times P}$ we have $\sqrt{\mathbb{P}_P} \otimes \sqrt{\mathbb{P}_P}(A) = 0 \iff \mathbb{P}_P \otimes \mathbb{P}_P(A) = 0$.

From Equation (41) and Equation (42) follows that $f(P, P') \stackrel{a.s.}{=} 0 \implies f(P, P) \stackrel{a.s.}{=} 0$. Consequently, we have $h(P, P') \stackrel{a.s.}{=} h^*(P, P') \implies h(P, P) \stackrel{a.s.}{=} h^*(P, P)$. \square

Remark. Note that any random variable with outcomes restricted to $\Delta^d = \{(p_1, \dots, p_d)^\top \in [0, 1]^d \mid \sum_{i=1}^d p_i = 1\} \subset \mathbb{R}^d$ has a singular distribution, since Δ^d is a null set with respect to the d dimensional Lebesgue measure λ^d . This would then fall outside of the conditions stated in Theorem 2. However, we can circumvent this simply by transforming (bijectively) the outcome space to $\Delta_r^d = \{(p_1, \dots, p_{d-1})^\top \in [0, 1]^{d-1} \mid 0 \leq \sum_{i=1}^{d-1} p_i \leq 1\} \subset \mathbb{R}^{d-1}$, which has non-zero mass according to the $d-1$ dimensional Lebesgue measure λ^{d-1} .

C.3 Proofs for Section 3.2

We give proofs of various statements of Section 3.2.

Proof for Equation (27)

Instead of proving the whole kernel ridge regression approach end-to-end, we bring Equation (27) into the form of an ordinary kernel ridge regression objective and then show that our solution in Equation (28) matches the ordinary solution.

For this, define $\tilde{Y}_i := \text{vec}_i \left(\Delta_{Yf(X)}^\top \Delta_{Yf(X)} \right) = \langle f(X_{i \bmod n+1}) - e_{Y_{i \bmod n+1}}, f(X_{\lceil i/n \rceil}) - e_{Y_{\lceil i/n \rceil}} \rangle \in \mathbb{R}$ and $\tilde{F}_i := (f(X_{i \bmod n+1}), f(X_{\lceil i/n \rceil})) \in \Delta^d \times \Delta^d$ with $i = 1 \dots n^2$, and $\tilde{h}(\tilde{p}) := h(\tilde{p}_1, \tilde{p}_2)$ for $\tilde{p} \in \Delta^d \times \Delta^d$, as well as $\tilde{\mathcal{H}} := \mathcal{H} \otimes \mathcal{H}$ with $\tilde{\phi}(\tilde{p}) := (\phi \otimes \phi)(\tilde{p}_1, \tilde{p}_2)$.

Then, we can write Equation (27) as

$$\begin{aligned} \hat{\mathcal{R}}_{\text{CE}, \lambda}(g) &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left(\langle f(X_i) - e_{Y_i}, f(X_j) - e_{Y_j} \rangle - h(f(X_j), f(X_j)) \right)^2 + \lambda \|g\|_{\mathcal{H} \otimes \mathcal{H}}^2 \\ &= \frac{1}{n^2} \sum_{i=1}^{n^2} \left(\tilde{Y}_i - \langle \tilde{g}, \tilde{\phi}(\tilde{F}_i) \rangle_{\tilde{\mathcal{H}}} \right)^2 + \lambda \|\tilde{g}\|_{\tilde{\mathcal{H}}}^2, \end{aligned} \quad (43)$$

which is ordinary kernel ridge regression in the last line (Bach, 2024). Bach (2024) shows its unique minimum is reached under certain assumptions if

$$\tilde{g} = (\tilde{Y}_1, \dots, \tilde{Y}_{n^2}) \left(\tilde{K}_{f(X)} + \lambda n^2 I \right)^{-1} \left(\tilde{\phi}(\tilde{F}_1), \dots, \tilde{\phi}(\tilde{F}_{n^2}) \right)^\top \quad (44)$$

with $[\tilde{K}_{f(X)}]_{ij} = \langle \tilde{\phi}(\tilde{F}_i), \tilde{\phi}(\tilde{F}_j) \rangle_{\tilde{\mathcal{H}}}$. Now, to reach our solution, note that it holds $\langle \tilde{\phi}(\tilde{F}_i), \tilde{\phi}(\tilde{p}) \rangle_{\tilde{\mathcal{H}}} = k(f(X_{i \bmod n+1}), \tilde{p}_1) k(f(X_{\lceil i/n \rceil}), \tilde{p}_2)$ and $\tilde{K}_{f(X)} = \mathbf{K}_{f(X)} \otimes \mathbf{K}_{f(X)}$, which gives

$$\begin{aligned} &\langle \tilde{g}, \tilde{\phi}(\tilde{p}) \rangle_{\tilde{\mathcal{H}}} \\ &= (\tilde{Y}_1, \dots, \tilde{Y}_{n^2}) \left(\tilde{K}_{f(X)} + \lambda n^2 I \right)^{-1} \left(k(f(X_{1 \bmod n+1}), \tilde{p}_1) k(f(X_{\lceil 1/n \rceil}), \tilde{p}_2), \dots, k(f(X_{n^2 \bmod n+1}), \tilde{p}_1) k(f(X_{\lceil n^2/n \rceil}), \tilde{p}_2) \right)^\top \\ &= \text{vec} \left(\Delta_{Yf(X)}^\top \Delta_{Yf(X)} \right) \left(\mathbf{K}_{f(X)} \otimes \mathbf{K}_{f(X)} + \lambda n^2 I \right)^{-1} \left(\mathbf{k}_{f(X)}(\tilde{p}_1) \otimes \mathbf{k}_{f(X)}(\tilde{p}_2) \right)^\top. \end{aligned} \quad (45)$$

The last line is the predictor we stated in Equation (28).

Proof for Equation (29)

By definition of h_{kr} and by using the eigenvalue decomposition $\mathbf{K}_{f(X)} = Q_{f(X)} \Lambda_{f(X)} Q_{f(X)}^\top$, we have

$$\begin{aligned} &h_{\text{kr}}(p, p') \\ &:= \text{vec}^\top \left(\Delta_{Yf(X)}^\top \Delta_{Yf(X)} \right) \left(\mathbf{K}_{f(X)} \otimes \mathbf{K}_{f(X)} + \lambda n^2 I \right)^{-1} \left(\mathbf{k}_{f(X)}(p) \otimes \mathbf{k}_{f(X)}(p') \right) \\ &= \text{vec}^\top \left(\Delta_{Yf(X)}^\top \Delta_{Yf(X)} \right) \left(Q_{f(X)} \otimes Q_{f(X)} \right) \left(\Lambda_{f(X)} \otimes \Lambda_{f(X)} + \lambda n^2 I \right)^{-1} \left(Q_{f(X)}^\top \otimes Q_{f(X)}^\top \right) \left(\mathbf{k}_{f(X)}(p) \otimes \mathbf{k}_{f(X)}(p') \right) \\ &= \left(\left(Q_{f(X)}^\top \otimes Q_{f(X)}^\top \right) \text{vec} \left(\Delta_{Yf(X)}^\top \Delta_{Yf(X)} \right) \right)^\top \left(\Lambda_{f(X)} \otimes \Lambda_{f(X)} + \lambda n^2 I \right)^{-1} \left(Q_{f(X)}^\top \otimes Q_{f(X)}^\top \right) \left(\mathbf{k}_{f(X)}(p) \otimes \mathbf{k}_{f(X)}(p') \right). \end{aligned} \quad (46)$$

Note it holds that $(A \otimes B) \text{vec}(C) = \text{vec}(BCA^\top)$ for matrices A, B, C and $\text{vec}^\top(A) \text{vec}(B) = \text{tr}(A^\top B)$.

Then, with the Hadamard product \odot and $\tilde{\Lambda}_X \in \mathbb{R}^{n \times n}$ with $[\tilde{\Lambda}_X]_{ij} := \frac{1}{(\Lambda_{f(X)})_{ii} (\Lambda_{f(X)})_{jj} + \lambda n^2}$, we have

$$\begin{aligned}
& \left(\left(Q_{f(X)}^\top \otimes Q_{f(X)}^\top \right) \text{vec} \left(\Delta_{Yf(X)}^\top \Delta_{Yf(X)} \right) \right)^\top \left(\Lambda_{f(X)} \otimes \Lambda_{f(X)} + \lambda n^2 I \right)^{-1} \left(Q_{f(X)}^\top \otimes Q_{f(X)}^\top \right) \left(\mathbf{k}_{f(X)}(p) \otimes \mathbf{k}_{f(X)}(p') \right) \\
&= \sum_{i=1}^{n^2} \text{vec}_i \left(Q_{f(X)}^\top \Delta_{Yf(X)}^\top \Delta_{Yf(X)} Q_{f(X)} \right) \text{vec}_i \left(Q_{f(X)}^\top \mathbf{k}_{f(X)}(p) \mathbf{k}_{f(X)}^\top(p') Q_{f(X)} \right) [\tilde{\Lambda}_{f(X)}]_{ij} \\
&= \text{tr} \left(Q_{f(X)}^\top \mathbf{k}_{f(X)}(p) \mathbf{k}_{f(X)}^\top(p') Q_{f(X)} \left(\tilde{\Lambda}_{f(X)} \odot Q_{f(X)}^\top \Delta_{Yf(X)}^\top \Delta_{Yf(X)} Q_{f(X)} \right) \right) \\
&= \mathbf{k}_{f(X)}^\top(p') Q_{f(X)} \left(\tilde{\Lambda}_{f(X)} \odot Q_{f(X)}^\top \Delta_{Yf(X)}^\top \Delta_{Yf(X)} Q_{f(X)} \right) Q_{f(X)}^\top \mathbf{k}_{f(X)}(p),
\end{aligned} \tag{47}$$

which shows Equation (29).

Proof for Equation (30)

Given the definition of $H^{\text{kk}r}$ we have

$$\begin{aligned}
H_{ij}^{\text{kk}r} &= \left[\mathbf{K}_{f(X)f(X')}^\top Q_{f(X)} \left(\tilde{\Lambda}_{f(X)} \odot Q_{f(X)}^\top \Delta_{Yf(X)}^\top \Delta_{Yf(X)} Q_{f(X)} \right) Q_{f(X)}^\top \mathbf{K}_{f(X)f(X')} \in \mathbb{R}^{n' \times n'} \right]_{ij} \\
&= \mathbf{k}_{f(X)}^\top(f(X'_i)) Q_{f(X)} \left(\tilde{\Lambda}_{f(X)} \odot Q_{f(X)}^\top \Delta_{Yf(X)}^\top \Delta_{Yf(X)} Q_{f(X)} \right) Q_{f(X)}^\top \mathbf{k}_{f(X)}(f(X'_j)) \in \mathbb{R}^{n' \times n'},
\end{aligned} \tag{48}$$

where the last line follows since $[\mathbf{k}_{f(X)}(f(X'_j))]_i = k(f(X_i), f(X'_j)) = [\mathbf{K}_{f(X)f(X')}]_{ij}$.