
Post-Calibration Techniques: Balancing Calibration and Score Distribution Alignment

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A binary scoring classifier can appear well-calibrated according to standard cal-
2 ibration metrics, even when the distribution of scores does not align with the
3 distribution of the true events. In this paper, we investigate the impact of post-
4 rocessing calibration on the score distribution (sometimes named “recalibration”).
5 Using simulated data, where the true probability is known, followed by real-world
6 datasets with prior knowledge on event distributions, we compare the performance
7 of an XGBoost model before and after applying calibration techniques. The results
8 show that while applying methods such as Platt scaling or isotonic regression can
9 improve the model’s calibration, they may also lead to an increase in the divergence
10 between the score distribution and the underlying event probability distribution.

11 1 Introduction

12 When estimating a probabilistic scoring classifier, the model must not only discriminate between
13 observations according to their class but also return scores that can be interpreted as probabilities. The
14 distribution of scores produced by the classifier should align with the underlying event distribution.
15 To assess whether classifiers return probabilistic scores, one must evaluate the model’s calibration
16 [4, 23, 8]. While some models, such as logistic regression (LR) when correctly specified, are known
17 to be well-calibrated [21], others, including ensemble methods like Random Forests (RF) [13, 3] and
18 XGBoost [10], are not inherently calibrated [27]. To assess a model’s ability to provide probabilistic
19 scores, the literature recommends evaluating its calibration using metrics like the Brier Score (BS, [4])
20 or the Integrated Calibration Index (ICI, [1]). When a model is not well-calibrated, post-processing
21 calibration methods including Platt scaling [28], or isotonic regression [36] are often applied to adjust
22 the scores [20, 11, 17]. After applying calibration techniques, these metrics generally indicate an
23 improvement in the model’s calibration relative to its initial state.

24 Since the true underlying probability distribution of the data is typically unobserved in practice,
25 calibration metrics are assessed solely on the classifier’s output range. [9] demonstrated with
26 simulated data that methods such as RF and XGBoost, can appear well-calibrated according to
27 standard calibration metrics and exhibit strong discrimination based on performance metrics, yet
28 still fail to align the score distribution with the true event distribution. This discrepancy can arise
29 when predicted scores from those algorithms lack the heterogeneity present in the underlying data
30 distribution. To address this, they suggest selecting model hyperparameters based on Kullback-Leibler
31 (KL) divergence, knowing the true event distribution in the case of simulated data. For real data,
32 where the true distribution is unknown, the approach involves prior information about the underlying
33 data distribution to better align it with the predicted score distribution. Their analysis only considers
34 model evaluation to accurately interpret predicted scores as probabilities, typically through calibration
35 metrics. However, many practitioners employ calibration techniques to ensure that output scores

36 represent probabilistic estimates. In this paper, we examine how calibration techniques affect the
 37 variability of score distribution in XGBoost binary classifiers, comparing it to the true underlying
 38 data distribution using KL divergence. We find that these post-processing methods often reduce score
 39 heterogeneity. Additionally, tree-based models optimized for KL divergence, rather than calibration or
 40 performance metrics, tend to align more closely with the underlying data distribution after calibration.
 41 Using simulated data, where the true probability is known, followed by real-world datasets with prior
 42 knowledge of event distributions, we assess the accuracy of XGBoost predicted scores as probabilistic
 43 estimates before and after applying calibration techniques.

44 2 Calibration

45 We focus on the context of a binary scoring classifier. Let $Y \in \mathcal{Y} = \{0, 1\}$ be a binary response
 46 variable, and let $\mathbf{X} \in \mathcal{X} = \mathbb{R}^d$ denote features. The goal is to predict $s(\mathbf{X}) = \mathbb{P}(Y = 1 | \mathbf{X})$, using
 47 a sample of n i.i.d. observations $(\mathbf{x}_i, y_i)_{i=1}^n$. We estimate this probability $\hat{s}(\mathbf{x}_i) \in [0, 1]$ using an
 48 XGBoost classifier, which produces a distribution of estimated scores $\hat{s}(\mathbf{X})$. If the score distribution is
 49 poorly calibrated, these scores cannot be interpreted as the “true underlying probabilities” [32, 18, 16].
 50 A model \hat{s} is well-calibrated for a binary variable Y when [30]:

$$\mathbb{P}(Y = 1 | \hat{s}(\mathbf{X})) = \mathbb{E}[Y | \hat{s}(\mathbf{X})] = \hat{s}(\mathbf{X}) \quad \text{a.s.}, \quad (1)$$

51 i.e., equivalently, $\mathbb{E}[Y | \hat{s}(\mathbf{X}) = p] = p, \forall p \in [0, 1]$.

52 2.1 Calibration Metrics

53 To measure calibration, the literature suggests various metrics. Here, we focus on two of them:
 54 BS and ICI. The former [12, 17, 28, 29] is often used to assess a model’s calibration. It writes:
 55 $\text{BS} = n^{-1} \sum_{i=1}^n (\hat{s}(\mathbf{x}_i) - y_i)^2$ [4]. More recently, [1] introduced the ICI, a metric that relies on the
 56 calibration curve. In the binary case, the calibration curve writes $g : [0, 1] \rightarrow [0, 1], \quad p \mapsto g(p) :=$
 57 $\mathbb{E}[Y | \hat{s}(\mathbf{X}) = p]$. For a well-calibrated model, the g function is the identity function $g(p) = p$.
 58 While the calibration curve is usually estimated using bins [34, 19, 26], the ICI relies on a smoother
 59 version, based on splines. The empirical version writes $\text{ICI} = n^{-1} \sum_{i=1}^n |\hat{s}(\mathbf{x}_i) - \hat{g}(\hat{s}(\mathbf{x}_i))|$, which
 60 corresponds to computing the average of the absolute difference between the calibration curve and
 61 the 45-degree diagonal line, the latter representing perfect calibration.

62 2.2 Calibration Methods

63 When using scores generated by a model estimating the probability of a binary event, the literature
 64 advocates calibrating the model by applying the calibration curve g —which serves as a transformation
 65 function—on the scores [28, 36, 17, 19]. In this paper, we focus on two calibration methods: Platt
 66 Scaling, and isotonic regression.

67 **Platt Scaling** This methodology was initially introduced to map SVM outputs to well-calibrated
 68 posterior probabilities [28]. This parametric approach consists of fitting a LR to the binary response
 69 variable using predicted scores of a binary classifier as the unique feature. Platt scaling learns
 70 parameters, μ and s ($s > 0$ for a non-decreasing calibration map g) on a calibration set. The obtained
 71 calibrated probabilities are $g(\hat{s}(\mathbf{x})) = (1 + \exp\{-\frac{1}{s}(\hat{s}(\mathbf{x}) - \mu)\})^{-1}$. As pointed out in [17], Platt
 72 scaling is unable to learn the identity function g if the predicted scores are already calibrated.

73 **Isotonic Regression** This solution arises from a constrained optimization problem [36], solved
 74 using the Pool-Adjacent-Violators Algorithm, ensuring that corrected predicted scores remain
 75 monotonic: $\min_{\beta_1, \dots, \beta_n} \sum_{i=1}^n (y_{(i)} - \beta_i)^2$, s.t. $\beta_1 \leq \dots \leq \beta_n$. $y_{(i)}$ corresponds to the value
 76 in $\{y_1, \dots, y_n\}$ associated with the i -th largest predicted score $\{\hat{s}(\mathbf{x}_1), \dots, \hat{s}(\mathbf{x}_n)\}$. Isotonic regres-
 77 sion, like Platt scaling, assumes the initial model’s predicted scores $\hat{s}(\mathbf{x})$ are well-ordered, limiting
 78 its ability to correct non-monotonic probability distortions.

79 3 Score Heterogeneity

80 To accurately interpret predicted scores from a binary classifier as probabilistic estimates, since the
 81 true underlying probability $s(\mathbf{X})$ is unobservable, calibration metrics rely solely on the predicted

82 score range. When a binary classification model is well-calibrated, the distribution of its scores $\hat{s}(\mathbf{X})$,
 83 as defined by Eq. 1, should align with the actual probability of the event in the vicinity of score values.
 84 Therefore, calibration metrics cannot fully capture discrepancies between the score distribution and
 85 the true probability distribution of the response variable Y when the predicted score variability does
 86 not accurately reflect the latter.

87 **Kullback-Leibler divergence** [9] demonstrated through simulated data that scores from ensemble
 88 methods may exhibit less variability compared to the true underlying probabilities when selecting
 89 hyperparameters based on calibration (ICI) or performance (AUC) metrics. This reduced hetero-
 90 geneity makes calibration metrics less reliable for interpreting output scores as probabilities of
 91 event occurrence. Instead of evaluating discrepancies solely on predicted score values, the authors
 92 suggest evaluating the model’s probabilistic estimates using KL divergence between the overall score
 93 distribution, $\hat{s}(\mathbf{X})$, and the available information on the “true” distribution, $s(\mathbf{X})$. Additionally, the
 94 flexibility of tree-based methods such as XGBoost allows for the selection of model hyperparameters
 95 based on KL divergence rather than traditional performance metrics, without incurring significant
 96 losses in performance or calibration and ensuring that the predicted score distribution more closely
 97 aligns with prior knowledge.

98 **Bayesian Framework** When working with simulated data, the distribution $s(\mathbf{X})$ is fully known,
 99 allowing for the direct computation of KL divergence with $\hat{s}(\mathbf{X})$. However, with real data, the
 100 KL divergence can only be computed by relying on a prior belief about the distribution of $s(\mathbf{X})$,
 101 potentially informed by expert opinion, and thus assuming a prior distribution \mathcal{B} . In the following,
 102 as in [9], we take $s(\mathbf{X}) \sim \mathcal{B} = \text{Beta}(\alpha, \beta)$ as the assumed prior distribution where each true
 103 probability p_i of the i -th observation is a sample from \mathcal{B} . We observe a sequence of n independent
 104 (as the features \mathbf{X}_i are considered n i.i.d. random variables) but non-identically distributed binary
 105 random variables Y_i where $Y_i | s(\mathbf{X}_i) = p_i \sim \text{Bernoulli}(p_i)$. In this case, instead of selecting
 106 the model with hyperparameters that minimize the empirical mean of the KL divergence across
 107 individual distributions, given by $\frac{1}{n} \sum_{i=1}^n \hat{s}(\mathbf{x}_i) \log \frac{\hat{s}(\mathbf{x}_i)}{p_i} + (1 - \hat{s}(\mathbf{x}_i)) \log \frac{1 - \hat{s}(\mathbf{x}_i)}{1 - p_i}$, we adopt the
 108 approach proposed by [9] and minimize the distance between the prior distribution \mathcal{B} and the overall
 109 distribution of $\hat{s}(\mathbf{X})$.

110 **Calibration techniques** We extend the work of [9] by investigating how score heterogeneity
 111 predicted by certain XGBoost algorithms is affected after applying post-calibration techniques such
 112 as Platt scaling or isotonic regression. These methods can potentially reduce score heterogeneity;
 113 for instance, isotonic regression applies a stepwise function g . Additionally, with Platt scaling, the
 114 range of calibrated predicted scores is always narrower than the range of the initial scores when
 115 the parameter $s \geq \frac{1}{4}$ (see Appendix A.1). And, due to the concavity of the sigmoid function over
 116 $[0, +\infty]$, this post-calibration method tends to reduce the range of predicted scores more significantly
 117 when the initial scores are highly concentrated.

118 4 Numerical Experiments

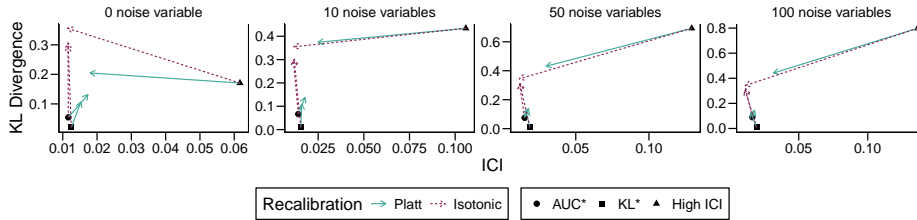
119 4.1 Simulated Data

120 We use the simulated data from [9]. We consider four data-generating processes (DGPs), all of which
 121 use a logistic link function. The first three are from [25], the fourth adds interaction terms. More
 122 details can be found in Appendix B.1. For each DGP, we generate data that include more or less
 123 noise variables: 0, 10, 50 or 100. We split the data into four samples: the train and validation samples
 124 used to train an XGBoost model and select the set of hyperparameters, the calibration sample to
 125 train a calibrator using the selected model, and lastly, a test sample to assess the performance of
 126 models on unseen data. We select the model’s hyperparameters (number of boosting iterations and
 127 maximum tree depth) to optimize three different criteria on the validation set: maximizing AUC
 128 (AUC*), minimizing KL divergence (KL*), or, for illustrative purposes, producing a model that
 129 is poorly calibrated based on the ICI metric (High ICI). Once the hyperparameters are selected, a
 130 calibration technique is applied to the scores. This allows for a comparison of models on the test
 131 set, both before and after calibration, according to the chosen optimization criterion. We run the
 132 simulations on 100 replications for each configuration. The results for DGP 1 are shown in Fig. 1

133 (see Fig. C18 for full results and Table C3 for numerical values). The shapes represent models before
 134 calibration, and arrows indicate model performance after calibration.

135 The results highlight two key findings. First, when hyperparameters are optimized for AUC, isotonic
 136 regression enhances calibration if the model is initially well-calibrated, while Platt scaling often fails
 137 due to the logistic family lacking the identity function. Second, calibration techniques generally
 138 increase the KL divergence between scores and true probabilities, except when the model is poorly
 139 calibrated (High ICI) with multiple noise variables. As shown in the histograms from Figs. C2 to C17,
 140 after applying calibration techniques, the model becomes better calibrated but the score distribution
 141 moves further from the true distribution.

142 We conclude that optimizing hyperparameters based on KL divergence rather than AUC leads
 143 to a closer alignment between the score distribution and the true distribution, even after post-
 144 processing calibration. Additionally, calibration techniques can even increase the distance from the
 145 true probability distribution by reducing score variability.



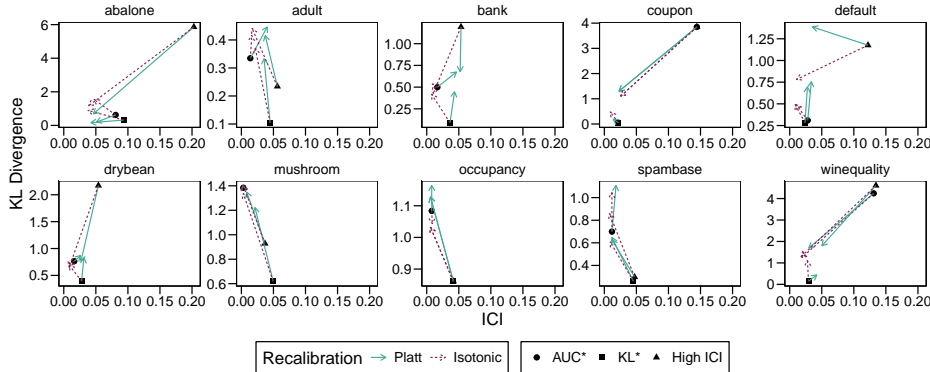
Notes: AUC*, KL*, High ICI: models selected by optimizing AUC, KL divergence, or by selecting a high ICI.

Figure 1: Average KL divergence and ICI before and after recalibration, for DGP 1.

146 4.2 Real Data

147 The 10 datasets from the UCI ML Repository used in [9] are used here (see details in Appendix B.2).
 148 For each dataset, we apply the method outlined in Section 4.1, this time calculating the KL divergence
 149 between the predicted score distribution and the prior distribution described in Section 3.¹

150 The results across the 10 datasets are shown in Fig. 2, with detailed metric values in Table C4. The
 151 findings are consistent with Section 4.1: isotonic regression preserves calibration when the model is
 152 already well-calibrated, while Platt scaling often worsens it. After calibration, KL* models align best
 153 with the prior distribution and, in some cases, show better ICI than AUC* or High ICI models. For
 154 the adult, mushroom, occupancy, and spambase datasets, calibration techniques generally increase
 155 KL divergence, reducing score distribution heterogeneity.



Notes: AUC*, KL*, High ICI: models selected by optimizing AUC, KL divergence, or by selecting a high ICI.

Figure 2: Average KL divergence and ICI before and after recalibration.

¹For illustration purposes, the parameters of the prior distribution \mathcal{B} are estimated via maximum likelihood using scores from a GAMSEL model [6], where the event is regressed on the variables generating the data.

References

- 156
157 [1] Austin, P. C. and Steyerberg, E. W. (2019). The integrated calibration index (ici) and related metrics for
158 quantifying the calibration of logistic regression models. *Statistics in Medicine*, 38(21):4051–4065.
- 159 [2] Becker, B. and Kohavi, R. (1996). Adult. UCI Machine Learning Repository. doi:10.24432/C5XW20.
- 160 [3] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1):5–32.
- 161 [4] Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly Weather Review*,
162 78(1):1–3.
- 163 [5] Candanedo, L. (2016). Occupancy Detection . UCI Machine Learning Repository. doi:10.24432/C5X01N.
- 164 [6] Chouldechova, A. and Hastie, T. (2015). Generalized additive model selection. *arXiv:1506.03850*.
- 165 [7] Cortez, P., Cerdeira, A., Almeida, F., Matos, T., and Reis, J. (2009). Wine Quality. UCI Machine Learning
166 Repository. doi:10.24432/C56S3T.
- 167 [8] Dawid, A. P. (1982). The well-calibrated bayesian. *Journal of the American Statistical Association*,
168 77(379):605–610.
- 169 [9] Fernandes Machado, A., Charpentier, A., Flachaire, E., Gallic, E., and Hu, F. (2024). Probabilistic scores of
170 classifiers, calibration is not enough. *arXiv:2408.03421*.
- 171 [10] Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *The Annals of
172 Statistics*, 29(5):1189 – 1232.
- 173 [11] Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. (2017). On calibration of modern neural networks. In
174 Precup, D. and Teh, Y. W., editors, *Proceedings of the 34th International Conference on Machine Learning*,
175 volume 70 of *Proceedings of Machine Learning Research*, pages 1321–1330. PMLR.
- 176 [12] Gupta, K., Rahimi, A., Ajanthan, T., Sminchisescu, C., Mensink, T., and Hartley, R. I. (2021). Calibration
177 of neural networks using splines. In *International Conference on Learning Representations (ICLR)*.
- 178 [13] Ho, T. K. (1995). Random decision forests. In *Proceedings of 3rd International Conference on Document
179 Analysis and Recognition*, volume 1, pages 278–282 vol.1.
- 180 [14] Hopkins, M., Reeber, E., Forman, G., and Suermondt, J. (1999). Spambase. UCI Machine Learning
181 Repository. doi:10.24432/C53G6X.
- 182 [15] Koklu, M. and Ali Ozkan, I. (2020). Dry Bean. UCI Machine Learning Repository. doi:10.24432/C50S4B.
- 183 [16] Konek, J. (2016). Probabilistic knowledge and cognitive ability. *Philosophical Review*, 125(4):509–587.
- 184 [17] Kull, M., Filho, T. M. S., and Flach, P. (2017). Beyond sigmoids: How to obtain well-calibrated probabilities
185 from binary classifiers with beta calibration. *Electronic Journal of Statistics*, 11(2):5052 – 5080.
- 186 [18] Kull, M. and Flach, P. A. (2014). Reliability maps: a tool to enhance probability estimates and improve
187 classification accuracy. In *Machine Learning and Knowledge Discovery in Databases: European Conference,
188 ECML PKDD 2014, Nancy, France, September 15-19, 2014. Proceedings, Part II 14*, pages 18–33. Springer.
- 189 [19] Kumar, A., Liang, P. S., and Ma, T. (2019). Verified uncertainty calibration. In Wallach, H., Larochelle,
190 H., Beygelzimer, A., d’Alché-Buc, F., Fox, E., and Garnett, R., editors, *Advances in Neural Information
191 Processing Systems*, volume 32. Curran Associates, Inc.
- 192 [20] Leathart, T., Frank, E., Pfahringer, B., and Holmes, G. (2019). On calibration of nested dichotomies. In
193 *Springer-Verlag*, page 69–80.
- 194 [21] Mildenhall, S. J. (1999). A systematic relationship between minimum bias and generalized linear models.
195 In *Journal Proceedings of the Casualty Actuarial Society*, volume 86, pages 393–487.
- 196 [22] Moro, S., Rita, P., and Cortez, P. (2012). Bank Marketing. UCI Machine Learning Repository. doi:
197 <https://doi.org/10.24432/C5K306>.
- 198 [23] Murphy, A. H. (1972). Scalar and vector partitions of the probability score: Part i. two-state situation.
199 *Journal of Applied Meteorology and Climatology*, 11(2):273–282.
- 200 [24] Nash, W., Sellers, T., Talbot, S., Cawthorn, A., and Ford, W. (1995). Abalone. UCI Machine Learning
201 Repository. doi:10.24432/C55C7W.

- 202 [25] Ojeda, F. M., Jansen, M. L., Thiéry, A., Blankenberg, S., Weimar, C., Schmid, M., and Ziegler, A. (2023).
 203 Calibrating machine learning approaches for probability estimation: A comprehensive comparison. *Statistics*
 204 *in Medicine*, 42(29):5451–5478.
- 205 [26] Pakdaman Naeini, M., Cooper, G., and Hauskrecht, M. (2015). Obtaining well calibrated probabilities
 206 using bayesian binning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 29(1):2901–2907.
- 207 [27] Park, Y. and Ho, J. C. (2020). Califorest: Calibrated random forest for health data. *Proceedings of the*
 208 *ACM Conference on Health, Inference, and Learning 2020*, pages 40–50.
- 209 [28] Platt, J. (1999). Probabilistic outputs for support vector machines and comparisons to regularized likelihood
 210 methods. *Advances in large margin classifiers*, 10(3):61–74.
- 211 [29] Rahimi, A., Shaban, A., Cheng, C.-A., Hartley, R., and Boots, B. (2020). Intra order-preserving functions
 212 for calibration of multi-class neural networks. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and
 213 Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 13456–13467. Curran
 214 Associates, Inc.
- 215 [30] Schervish, M. J. (1989). A General Method for Comparing Probability Assessors. *The Annals of Statistics*,
 216 17(4):1856–1879.
- 217 [31] Schlimmer, J. (1987). Mushroom. UCI Machine Learning Repository. doi:10.24432/C5959T.
- 218 [32] Van Fraassen, B. C. (1995). Fine-grained opinion, probability, and the logic of full belief. *Journal of*
 219 *Philosophical logic*, 24(4):349–377.
- 220 [33] Wang, T., Rudin, C., Doshi-Velez, F., Liu, Y., Klampfl, E., and MacNeille, P. (2020). In-Vehicle Coupon
 221 Recommendation. UCI Machine Learning Repository. doi:10.24432/C5GS4P.
- 222 [34] Wilks, D. S. (1990). On the combination of forecast probabilities for consecutive precipitation periods.
 223 *Weather and Forecasting*, 5(4):640–650.
- 224 [35] Yeh, I.-C. (2016). Default of Credit Card Clients. UCI Machine Learning Repository.
 225 doi:10.24432/C55S3H.
- 226 [36] Zadrozny, B. and Elkan, C. (2002). Transforming classifier scores into accurate multiclass probability
 227 estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and*
 228 *data mining*, pages 694–699.

229 A Platt scaling

230 A.1 Reduction in Score Range

231 Platt scaling learns parameters, μ and s ($s > 0$ for a non-decreasing calibration map g) on a calibration
 232 set. The obtained calibrated probabilities are:

$$g(\hat{s}(\mathbf{x})) = \frac{1}{1 + \exp\left\{-\frac{1}{s}(\hat{s}(\mathbf{x}) - \mu)\right\}}.$$

233 With Platt scaling, the range of calibrated predicted scores is always narrower than the range of
 234 the initial scores when the parameter $s \geq \frac{1}{4}$. Indeed, since $\rho = \frac{1}{4}$ is the minimum value for which
 235 $\sigma(x) = \frac{1}{1 + \exp -x}$ remains ρ -Lipschitz on \mathbb{R} , for $x_1 < x_2 \in \mathbb{R}$, we have:

$$\left| \sigma\left(\frac{x_2 - \mu}{s}\right) - \sigma\left(\frac{x_1 - \mu}{s}\right) \right| \leq \frac{1}{4} \left| \frac{x_2 - \mu}{s} - \frac{x_1 - \mu}{s} \right| \leq \frac{1}{4s} |x_2 - x_1| \text{ with } s > 0.$$

236 As a result, if $s \geq \frac{1}{4}$, the range of $(\hat{s}(\mathbf{x}_i))_{i=1}^n$ is larger than the range of the calibrated scores with
 237 Platt scaling $(g(\hat{s}(\mathbf{x}_i)))_{i=1}^n$. Let \hat{s}_m (resp. \hat{s}_M) denote the minimum (resp. the maximum) value of
 238 $(\hat{s}(\mathbf{x}_i))_{i=1}^n$. If $s \geq \frac{1}{4}$, we have:

$$\left| \frac{g(\hat{s}_M) - g(\hat{s}_m)}{\hat{s}_M - \hat{s}_m} \right| = \left| \frac{\sigma\left(\frac{\hat{s}_M - \mu}{s}\right) - \sigma\left(\frac{\hat{s}_m - \mu}{s}\right)}{\hat{s}_M - \hat{s}_m} \right| \leq 1.$$

239 And, due to the concavity of the sigmoid function over $[0, +\infty]$, this post-calibration method tends to
 240 reduce the range of predicted scores more significantly when the initial scores are highly concentrated.

241 **B Data**

242 **B.1 Simulated Data**

243 To simulate data, we consider the DGPs from [9]. The first three are from Ojeda et al. [25]. In the
 244 fourth, an interaction term between two predictors is added. Each scenario uses a logistic model to
 245 generate the outcome. Let Y_i be a binary variable following a Bernoulli distribution: $Y_i \sim B(p_i)$,
 246 where p_i is the probability of observing $Y_i = 1$. The probability p_i is defined by:

$$p_i = \mathbb{P}(Y = 1 \mid \mathbf{x}_i) = [1 + \exp(-\eta_i)]^{-1}. \quad (2)$$

247 For the second DGP, to introduce non-linearities, p^3 is used as true probabilities instead of p .

248 For all DGPs, $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$, where \mathbf{x}_i is a vector of covariates and $\boldsymbol{\beta}$ is a vector of arbitrary scalars.
 249 The covariate vector includes two continuous predictors for DGPs 1 and 2. For DGP 3, it includes
 250 five continuous and five categorical predictors. For DGP 4, it contains three continuous variables,
 251 the square of the first variable, and an interaction term between the second and third variables.
 252 Specifically, $\eta_i = \beta_1 x_{1,i} + \beta_2 x_{2,i} + \beta_3 x_{3,i} + \beta_4 x_{1,i}^2 + \beta_5 x_{2,i} \times x_{3,i}$. Continuous predictors are
 253 drawn from $\mathcal{N}(0, 1)$. Categorical predictors consist of two variables with two categories, one with
 254 three categories, and one with five categories, all uniformly distributed. The values of coefficients $\boldsymbol{\beta}$
 255 are reported in Table B1.

256 For each DGP, we generate data considering four scenarios with varying numbers of noise variables:
 257 0, 10, 50, or 100 variables drawn from $\mathcal{N}(0, 1)$.

258 For the fourth DGP, to achieve a similar probability distribution to DGP 1, we perform resampling
 259 using a rejection algorithm (the algorithm is detailed in [9]).

DGP	No. Cont.	No. Cat.	No. Noise	$\boldsymbol{\beta}$	Type η
1	2	0	{0, 10, 50, 100}	(.5, 1)	Linear terms
2			Same as DGP 2, but with probabilities p^3		
3	5	5	{0, 10, 50, 100}	(.1, .2, .3, .4, .5, .01, .02, .03, .04, .05)	Linear terms
4	3	0	{0, 10, 50, 100}	(.5, 1, .3)	Non-linear terms

Notes: No. Cont., No. Cat. and No. Noise correspond to the number of continuous, categorical and noise variables, respectively.

Table B1: Parameters of the different scenarios.

260 The datasets are split into four parts: a training sample, a validation sample, a calibration sample,
 261 and a test sample, each containing 10,000 observations. The empirical distribution of samples of
 262 from each DGP are shown in Fig. B1.

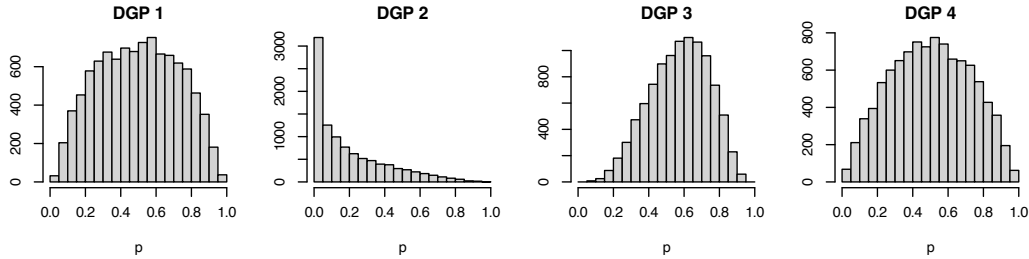


Figure B1: Distribution of the underlying probabilities in the different categories of scenarios.

263 **B.2 Real Data**

264 The main characteristics of the datasets are summarized in Table B2.

265 Most of the datasets used are associated with classification tasks. If not, they contain a binary
 266 variable suitable for classification or a variable that can be converted into a binary variable. The target
 267 variables for each dataset are as follows:

Table B2: Key characteristics of the datasets

Dataset	n	No. predictors	No. num. predictors	Prop. target = 1	Reference	License
abalone	4,177	8	8	0.37	Nash et al. [24]	CC BY 4.0
adult	32,561	14	6	0.24	Becker and Kohavi [2]	CC BY 4.0
bank	45,211	16	7	0.12	Moro et al. [22]	CC BY 4.0
default	30,000	23	14	0.22	Yeh [35]	CC BY 4.0
drybean	13,611	16	16	0.26	Koklu and Ali Ozkan [15]	CC BY 4.0
coupon	12,079	22	0	0.57	Wang et al. [33]	CC BY 4.0
mushroom	8,124	21	0	0.52	Schlimmer [31]	CC BY 4.0
occupancy	20,560	5	5	0.23	Candanedo [5]	CC BY 4.0
winequality	6,495	12	11	0.63	Cortez et al. [7]	CC BY 4.0
spambase	4,601	57	57	0.39	Hopkins et al. [14]	CC BY 4.0

Notes: n represents the number of observations, 'No. predictors' the total number of predictors, 'No. num. predictors' the number of numeric predictors, and 'Prop. target = 1' the proportion of positive observed events.

- 268 • abalone: gender of abalones (1 for male, 0 for female); originally used to predict the size
- 269 of abalones.
- 270 • adult: high income (1 if income \geq 50k per year).
- 271 • bank: subscription to a term deposit (1 if yes, 0 otherwise).
- 272 • default: default payment (1 if default, 0 otherwise).
- 273 • drybean: type of dry bean (1 if dermason, 0 otherwise); originally a multi-class variable.
- 274 • coupon: acceptance of a recommended coupon in different driving scenarios (1 if accepted,
- 275 0 otherwise).
- 276 • mushroom: mushroom classification (1 if edible, 0 otherwise).
- 277 • occupancy: prediction of room occupancy (1 if occupied, 0 otherwise); originally aimed at
- 278 predicting the age of occupancy from physical measurements.
- 279 • winequality: quality of wine (1 if quality \geq 6, 0 otherwise); originally a scale from 0 to
- 280 10, with 0 being bad quality and 10 being good quality.
- 281 • spambase: email classification (1 if spam, 0 otherwise).

282 C Numerical Experiments

283 C.1 Simulated Data

284 For each of the four DGPs (see Section B.1) and each configuration of the number of noise variables
 285 (0, 10, 50, or 100), we generate 100 sample replications. For each sample, we train an XGBoost model
 286 on 10,000 observation using the `xgb.train` function from the R package `xgboost`. The learning
 287 rate is set to 0.3. The tree depth (argument `max_depth`) varies according to the following values: 2,
 288 4, 6. The number of boosting iterations (argument `nrounds`) ranges from 1 to 400. All variables
 289 (predictors and, if applicable, noise variables) are included in the model without transformation.

290 For each model configuration, we select the hyperparameters based on different criteria using the
 291 validation set results. Specifically, we make three model choices:

- 292 • AUC*: hyperparameters are selected to maximize the AUC.
- 293 • KL*: hyperparameters are chosen to minimize the Kullback-Leibler divergence between
 294 the scores on the validation set and the true probability distribution (observable here in the
 295 context of simulated data).
- 296 • High ICI: hyperparameters are selected to produce relatively poor calibration, as measured
 297 by the ICI. Specifically, we select the model with the smallest ICI among those with an ICI
 298 at least one standard deviation above the mean ICI obtained during grid search.

299 Once the hyperparameters are selected, we apply a recalibration method on an independent calibration
 300 set: either Platt scaling or isotonic regression.

301 The model performance is then evaluated on a test set, allowing for comparison based on: (i) the
 302 metric used to select the hyperparameters, and (ii) whether or not calibration techniques were applied
 303 to the scores.

304 Figs C2 to C17 display the empirical distribution of scores for a single replication (the first one) in
 305 each of the 4×4 configurations (4 DGPs and 4 different values for the number of noise variables
 306 introduced in the training data). In each figure, the first row shows the distribution of test set scores
 307 without applying any calibration technique to the selected model. The second row, in green, shows
 308 the score distributions after applying Platt scaling for calibration. The third row, in purple, displays
 309 the score distributions after applying isotonic regression. The columns correspond to the criteria used
 310 to select the hyperparameters based on the validation set results: AUC, Brier, ICI, KL, or High ICI.

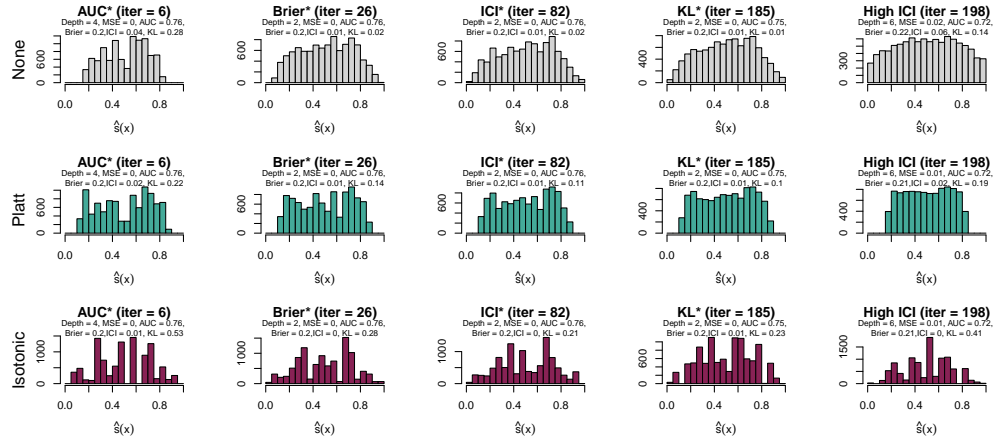


Figure C2: Distribution of estimated scores for XGB: **DGP 1, 0 noise variable**, single replication.

Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

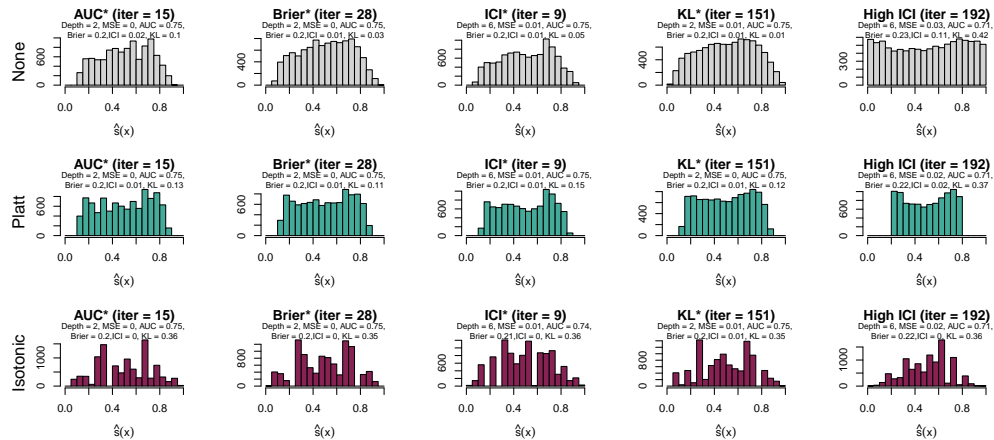


Figure C3: Distribution of estimated scores for XGB: **DGP 1, 10 noise variables**, single replication.

Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

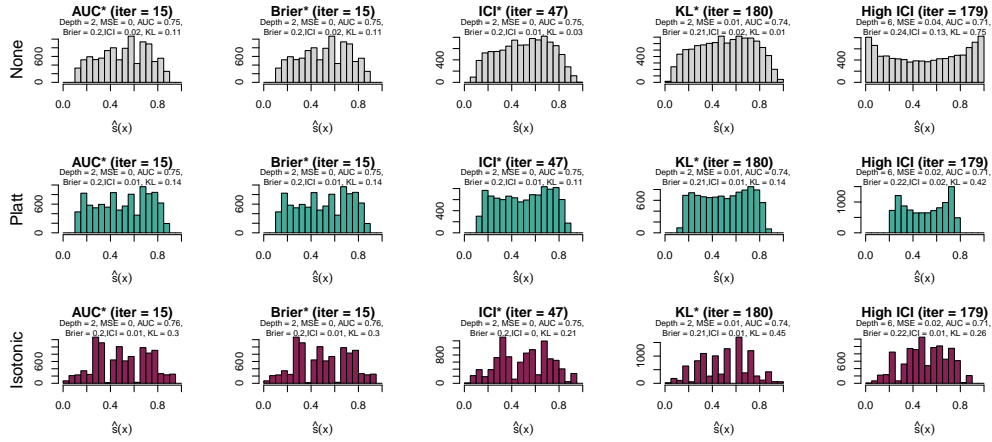


Figure C4: Distribution of estimated scores for XGB: **DGP 1, 50 noise variables**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

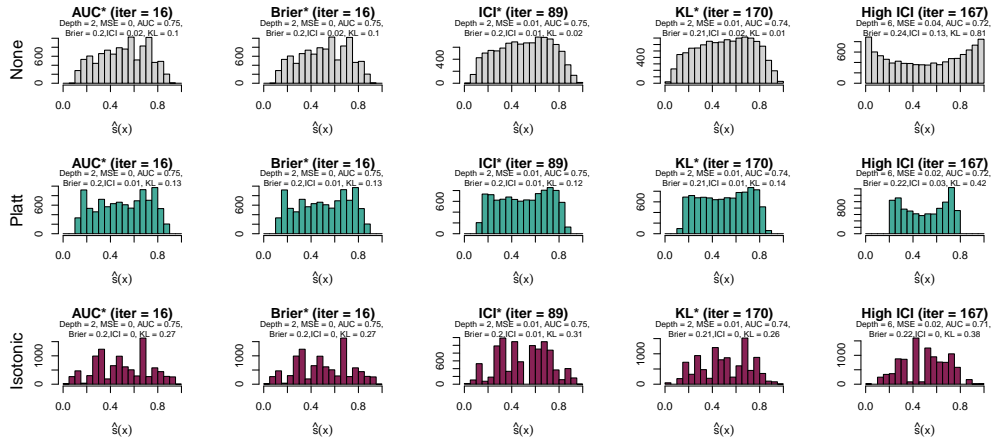


Figure C5: Distribution of estimated scores for XGB: **DGP 1, 100 noise variables**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

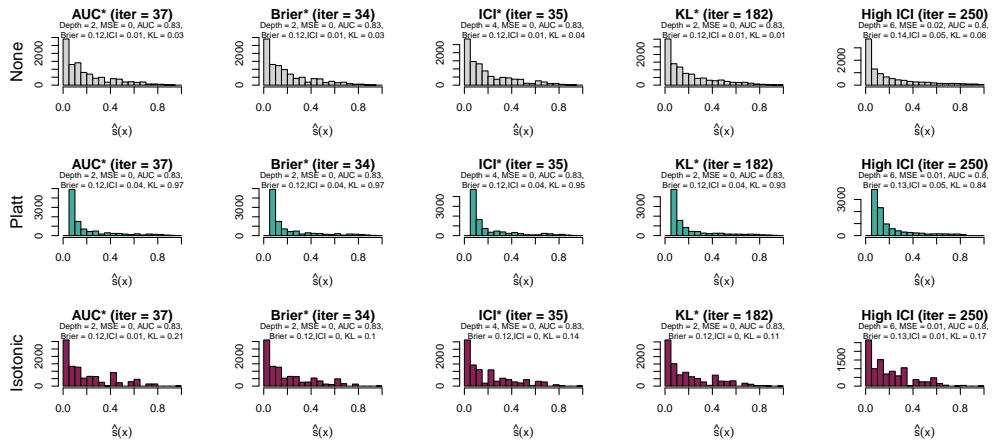


Figure C6: Distribution of estimated scores for XGB: **DGP 2, 0 noise variable**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

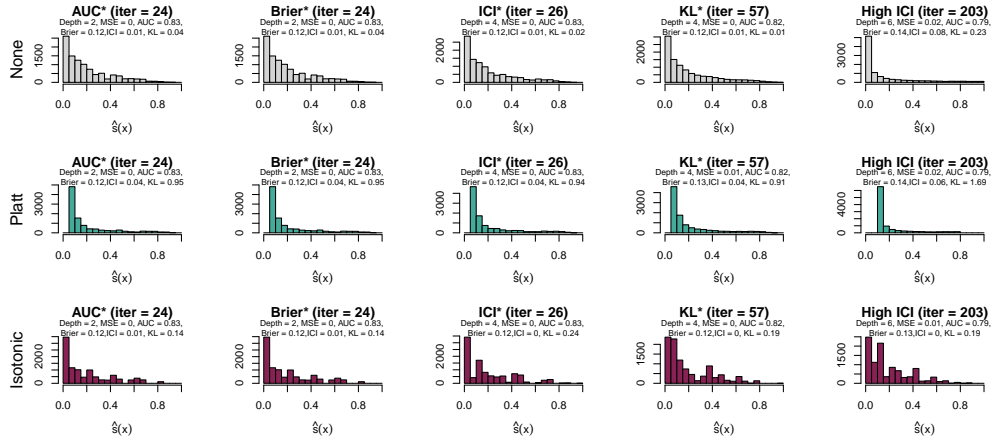


Figure C7: Distribution of estimated scores for XGB: **DGP 2, 10 noise variables**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

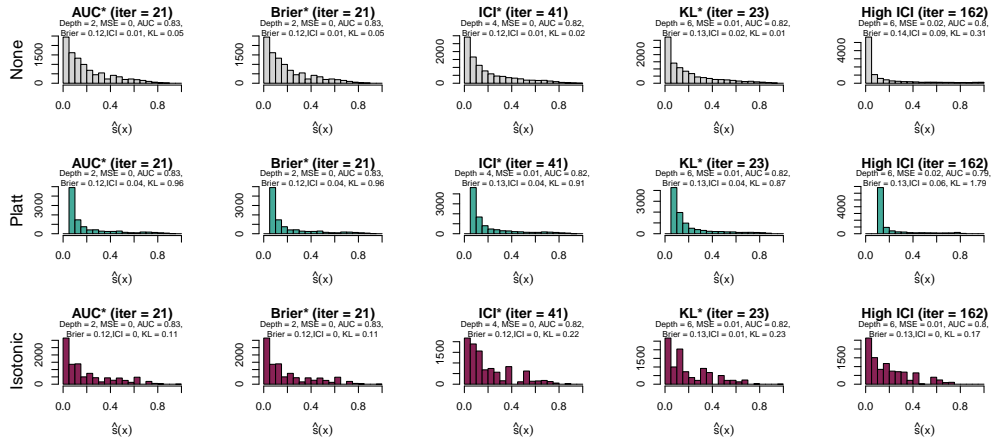


Figure C8: Distribution of estimated scores for XGB: **DGP 2, 50 noise variables**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

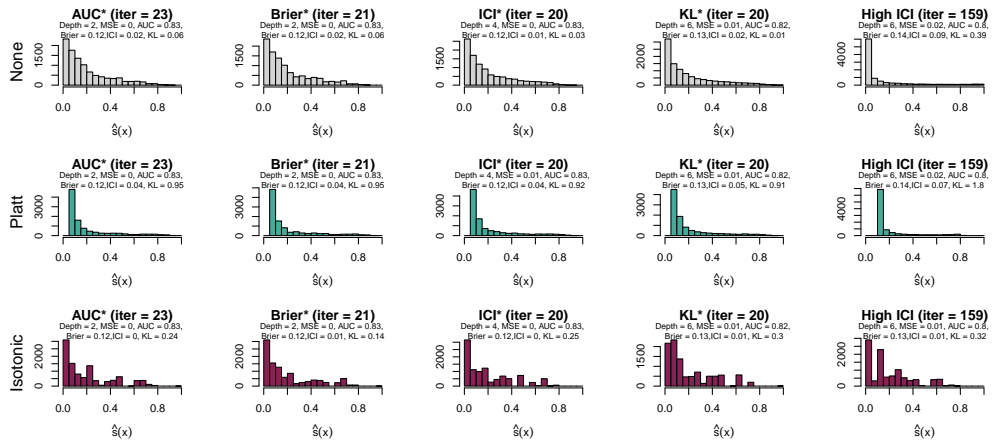


Figure C9: Distribution of estimated scores for XGB: **DGP 2, 100 noise variables**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

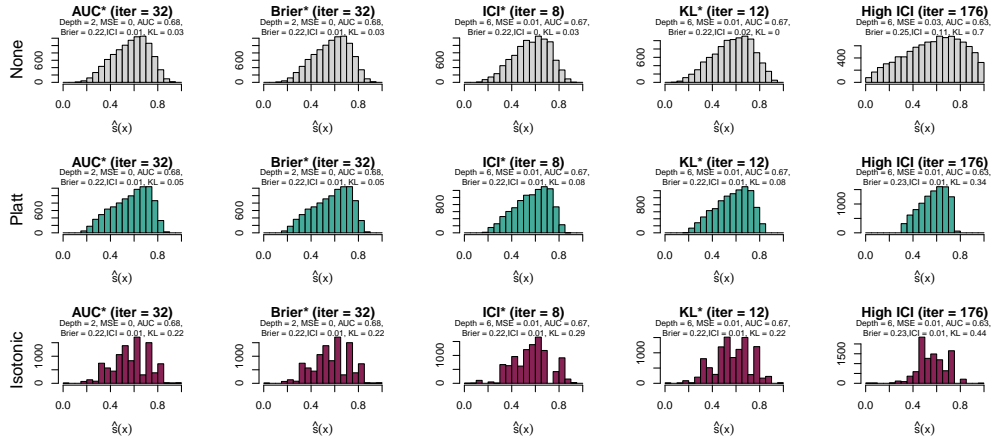


Figure C10: Distribution of estimated scores for XGB: **DGP 3, 0 noise variable**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

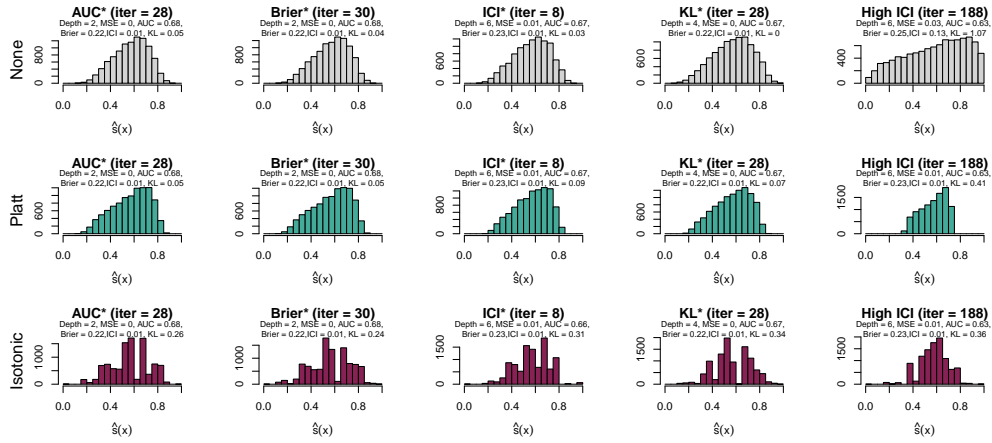


Figure C11: Distribution of estimated scores for XGB: **DGP 3, 10 noise variables**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

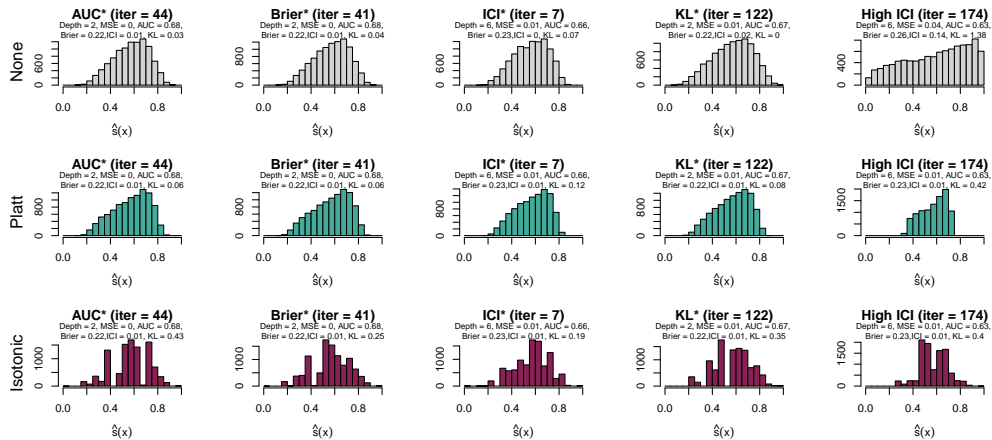


Figure C12: Distribution of estimated scores for XGB: **DGP 3, 50 noise variables**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

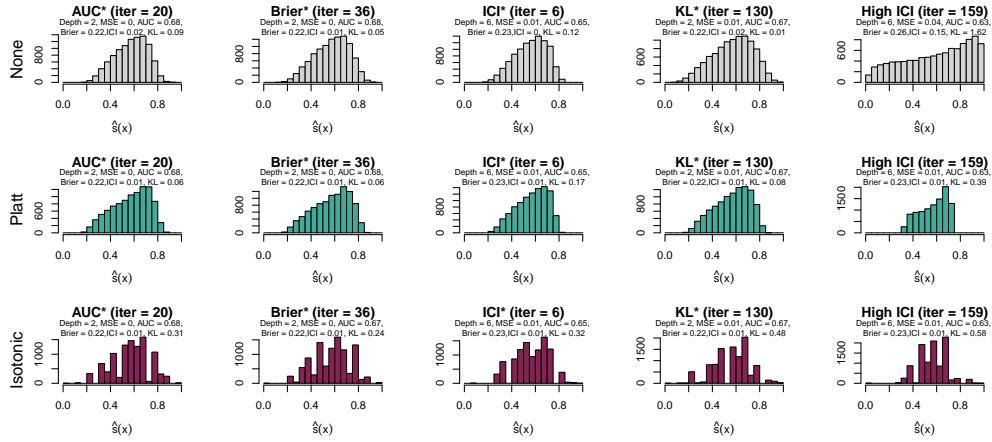


Figure C13: Distribution of estimated scores for XGB: **DGP 3, 100 noise variables**, single replication.

Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

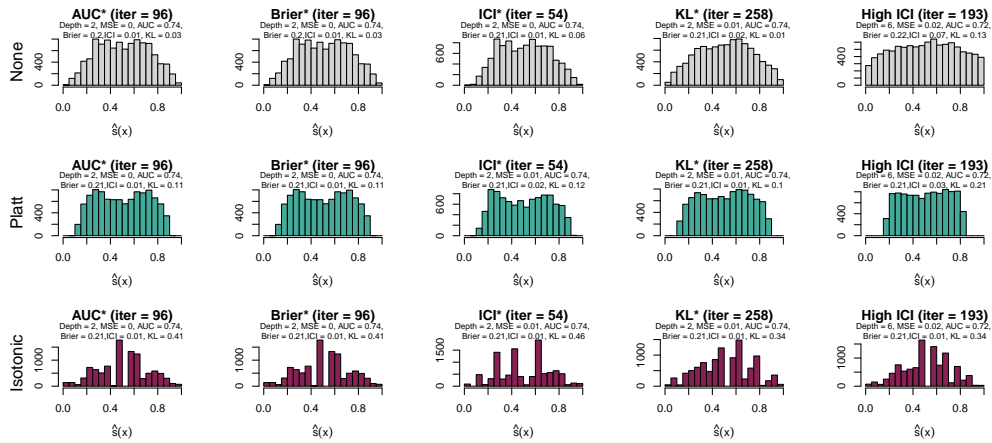


Figure C14: Distribution of estimated scores for XGB: **DGP 4, 0 noise variable**, single replication.

Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

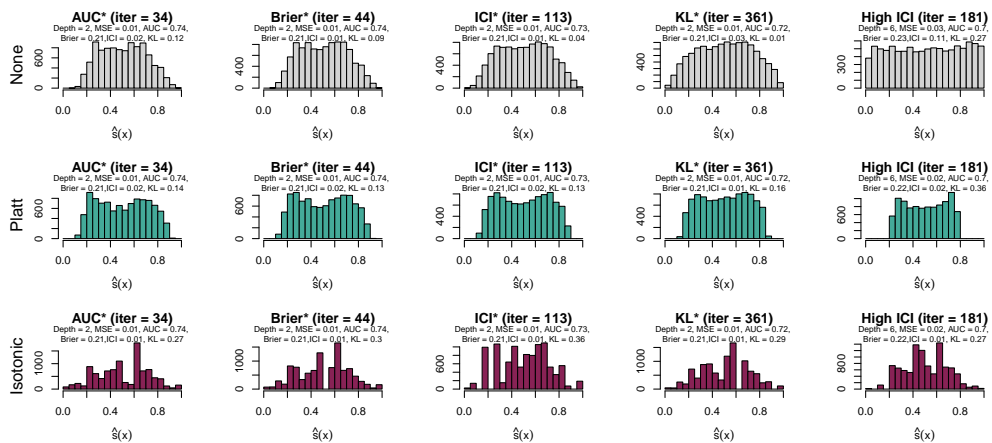


Figure C15: Distribution of estimated scores for XGB: **DGP 4, 10 noise variables**, single replication.

Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

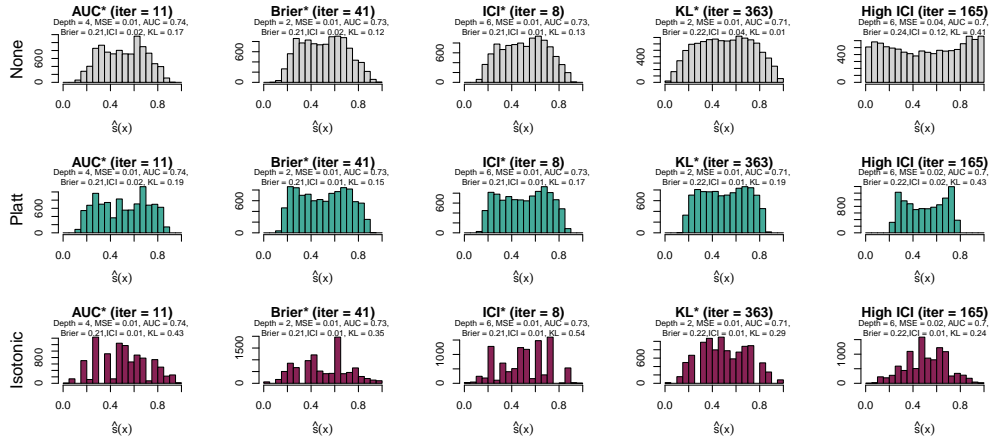


Figure C16: Distribution of estimated scores for XGB: **DGP 4, 50 noise variables**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

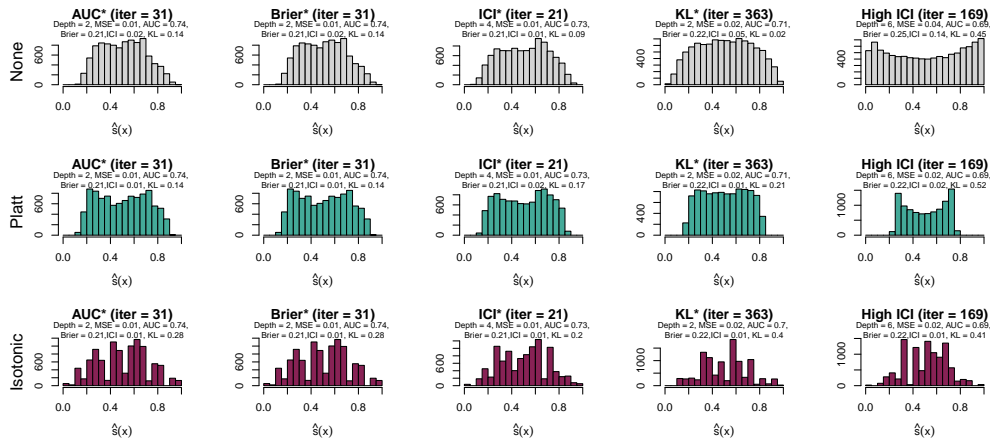


Figure C17: Distribution of estimated scores for XGB: **DGP 4, 100 noise variables**, single replication.
 Notes: AUC*, Brier*, ICI*, and KL*: models selected based on optimizing AUC, Brier score, ICI, and Kullback-Leibler divergence, resp.

311 Table C3 reports the average values of metrics calculated on the test set over 100 replications for
 312 each DGP and each number of noise variables in the training data. The values are presented for
 313 models selected by optimizing, on the validation set, either AUC (AUC*) or KL divergence (KL*),
 314 as well as for a model with poor calibration (High ICI). The metrics are calculated before applying
 315 any calibration method (column “None”), after applying Platt Scaling, and after applying isotonic
 316 regression calibration.

317 Fig. C18 shows the performance of the models, measured by the KL divergence between the test set
 318 score distribution (x -axis) and the true probability distribution (y -axis), before and after applying
 319 calibration methods. The values represent the average of these two metrics over 100 replications for
 320 each DGP (rows), based on the number of noise variables in the training set (columns). The point
 321 corresponds to the model whose hyperparameters (number of boosting iterations and tree depth)
 322 are selected to maximize AUC on the validation set. The square represents the model selected by
 323 minimizing the Kullback-Leibler divergence between the score distribution on the validation set and
 324 the true probability distribution. The triangle denotes a model with poor calibration on the test set.
 325 Solid green arrows illustrate the change in metrics after Platt scaling calibration, while dashed purple
 326 arrows indicate changes after isotonic regression calibration.

Table C3: Comparison of metrics computed on the validation set for models selected based on AUC, KL divergence, or ICI across 100 replications. Standard errors are provided in parentheses.

DGP	Noise	Optim.	None			Platt Scaling			Isotonic		
			BS	ICI	KL	BS	ICI	KL	BS	ICI	KL
1	0	AUC*	0.201 (0.002)	0.012 (0.005)	0.054 (0.038)	0.201 (0.002)	0.017 (0.005)	0.131 (0.029)	0.201 (0.002)	0.012 (0.004)	0.295 (0.099)
		KL*	0.202 (0.002)	0.012 (0.005)	0.021 (0.006)	0.202 (0.002)	0.015 (0.005)	0.107 (0.016)	0.202 (0.002)	0.012 (0.004)	0.307 (0.101)
		High ICI	0.217 (0.003)	0.062 (0.006)	0.171 (0.101)	0.212 (0.002)	0.018 (0.005)	0.205 (0.065)	0.212 (0.002)	0.011 (0.005)	0.356 (0.122)
	10	AUC*	0.201 (0.002)	0.014 (0.005)	0.066 (0.033)	0.201 (0.002)	0.018 (0.005)	0.138 (0.028)	0.201 (0.002)	0.011 (0.004)	0.298 (0.106)
		KL*	0.204 (0.002)	0.015 (0.005)	0.01 (0.004)	0.204 (0.002)	0.016 (0.005)	0.109 (0.015)	0.204 (0.002)	0.012 (0.004)	0.302 (0.11)
		High ICI	0.229 (0.003)	0.106 (0.006)	0.434 (0.179)	0.216 (0.002)	0.025 (0.005)	0.374 (0.032)	0.216 (0.002)	0.012 (0.004)	0.355 (0.117)
	50	AUC*	0.201 (0.002)	0.015 (0.005)	0.075 (0.028)	0.201 (0.002)	0.018 (0.005)	0.139 (0.028)	0.201 (0.002)	0.012 (0.004)	0.307 (0.1)
		KL*	0.205 (0.002)	0.019 (0.005)	0.009 (0.003)	0.205 (0.002)	0.016 (0.004)	0.12 (0.02)	0.205 (0.002)	0.011 (0.004)	0.312 (0.11)
		High ICI	0.235 (0.003)	0.129 (0.006)	0.693 (0.089)	0.216 (0.002)	0.03 (0.005)	0.432 (0.036)	0.215 (0.002)	0.011 (0.004)	0.344 (0.115)
	100	AUC*	0.201 (0.002)	0.016 (0.005)	0.09 (0.031)	0.201 (0.002)	0.018 (0.005)	0.145 (0.026)	0.201 (0.002)	0.012 (0.004)	0.317 (0.111)
		KL*	0.206 (0.002)	0.02 (0.005)	0.009 (0.004)	0.206 (0.002)	0.015 (0.004)	0.125 (0.022)	0.206 (0.002)	0.011 (0.004)	0.304 (0.105)
		High ICI	0.236 (0.003)	0.136 (0.006)	0.798 (0.063)	0.216 (0.002)	0.032 (0.005)	0.442 (0.028)	0.215 (0.002)	0.011 (0.004)	0.342 (0.114)
2	0	AUC*	0.118 (0.002)	0.01 (0.004)	0.031 (0.027)	0.12 (0.002)	0.038 (0.004)	0.778 (0.221)	0.118 (0.002)	0.009 (0.004)	0.213 (0.073)
		KL*	0.12 (0.002)	0.01 (0.004)	0.013 (0.005)	0.121 (0.002)	0.038 (0.004)	0.863 (0.141)	0.12 (0.002)	0.009 (0.004)	0.214 (0.072)
		High ICI	0.131 (0.003)	0.048 (0.004)	0.131 (0.131)	0.13 (0.003)	0.05 (0.005)	0.839 (0.081)	0.127 (0.003)	0.009 (0.003)	0.241 (0.073)
	10	AUC*	0.119 (0.002)	0.012 (0.005)	0.036 (0.067)	0.12 (0.002)	0.038 (0.004)	0.758 (0.226)	0.119 (0.002)	0.009 (0.004)	0.209 (0.069)
		KL*	0.12 (0.002)	0.011 (0.003)	0.007 (0.003)	0.122 (0.002)	0.04 (0.004)	0.887 (0.076)	0.121 (0.002)	0.01 (0.004)	0.205 (0.074)
		High ICI	0.137 (0.003)	0.075 (0.004)	0.294 (0.136)	0.132 (0.003)	0.058 (0.005)	1.235 (0.363)	0.128 (0.002)	0.009 (0.003)	0.263 (0.093)
	50	AUC*	0.119 (0.002)	0.013 (0.004)	0.036 (0.016)	0.12 (0.002)	0.038 (0.004)	0.727 (0.226)	0.119 (0.002)	0.009 (0.004)	0.204 (0.073)
		KL*	0.121 (0.002)	0.011 (0.003)	0.006 (0.003)	0.123 (0.002)	0.041 (0.004)	0.894 (0.045)	0.121 (0.002)	0.009 (0.003)	0.215 (0.068)
		High ICI	0.139 (0.003)	0.089 (0.004)	0.434 (0.14)	0.133 (0.003)	0.064 (0.006)	1.793 (0.165)	0.127 (0.002)	0.009 (0.003)	0.232 (0.077)
	100	AUC*	0.119 (0.002)	0.014 (0.004)	0.041 (0.02)	0.121 (0.002)	0.038 (0.004)	0.727 (0.22)	0.119 (0.002)	0.009 (0.004)	0.212 (0.068)
		KL*	0.122 (0.002)	0.012 (0.004)	0.006 (0.003)	0.124 (0.002)	0.041 (0.003)	0.89 (0.032)	0.122 (0.002)	0.009 (0.003)	0.214 (0.074)
		High ICI	0.14 (0.003)	0.094 (0.004)	0.476 (0.082)	0.133 (0.003)	0.066 (0.004)	1.852 (0.112)	0.127 (0.002)	0.009 (0.004)	0.234 (0.067)
3	0	AUC*	0.22 (0.002)	0.01 (0.004)	0.012 (0.008)	0.221 (0.002)	0.012 (0.004)	0.042 (0.013)	0.221 (0.002)	0.011 (0.004)	0.285 (0.113)
		KL*	0.221 (0.002)	0.012 (0.004)	0.005 (0.002)	0.221 (0.002)	0.011 (0.004)	0.047 (0.014)	0.222 (0.002)	0.011 (0.004)	0.284 (0.115)
		High ICI	0.246 (0.002)	0.105 (0.005)	0.63 (0.086)	0.231 (0.001)	0.014 (0.004)	0.268 (0.047)	0.231 (0.001)	0.011 (0.004)	0.375 (0.106)
	10	AUC*	0.221 (0.001)	0.01 (0.004)	0.027 (0.021)	0.221 (0.002)	0.012 (0.004)	0.046 (0.014)	0.221 (0.002)	0.011 (0.004)	0.287 (0.097)
		KL*	0.222 (0.002)	0.014 (0.005)	0.004 (0.002)	0.222 (0.002)	0.012 (0.004)	0.056 (0.019)	0.222 (0.002)	0.011 (0.004)	0.284 (0.103)
		High ICI	0.253 (0.003)	0.127 (0.005)	0.933 (0.114)	0.232 (0.001)	0.015 (0.004)	0.366 (0.035)	0.232 (0.001)	0.011 (0.004)	0.382 (0.113)
	50	AUC*	0.221 (0.001)	0.013 (0.006)	0.054 (0.032)	0.221 (0.002)	0.012 (0.004)	0.049 (0.015)	0.222 (0.002)	0.011 (0.004)	0.27 (0.122)
		KL*	0.224 (0.002)	0.018 (0.006)	0.004 (0.002)	0.224 (0.002)	0.011 (0.004)	0.075 (0.023)	0.224 (0.002)	0.011 (0.004)	0.284 (0.104)
		High ICI	0.259 (0.003)	0.145 (0.006)	1.286 (0.167)	0.233 (0.001)	0.017 (0.005)	0.403 (0.028)	0.232 (0.001)	0.011 (0.004)	0.417 (0.113)
	100	AUC*	0.222 (0.001)	0.015 (0.006)	0.067 (0.031)	0.221 (0.002)	0.012 (0.004)	0.052 (0.016)	0.222 (0.002)	0.011 (0.004)	0.291 (0.111)
		KL*	0.225 (0.002)	0.019 (0.005)	0.004 (0.002)	0.224 (0.002)	0.011 (0.004)	0.08 (0.021)	0.225 (0.002)	0.011 (0.004)	0.301 (0.122)
		High ICI	0.261 (0.003)	0.152 (0.005)	1.454 (0.181)	0.233 (0.001)	0.017 (0.004)	0.418 (0.036)	0.233 (0.001)	0.011 (0.005)	0.426 (0.111)
4	0	AUC*	0.204 (0.002)	0.011 (0.004)	0.041 (0.021)	0.205 (0.002)	0.017 (0.004)	0.132 (0.02)	0.205 (0.002)	0.011 (0.005)	0.305 (0.095)
		KL*	0.206 (0.002)	0.018 (0.005)	0.011 (0.004)	0.206 (0.002)	0.015 (0.004)	0.115 (0.012)	0.206 (0.002)	0.011 (0.004)	0.288 (0.105)
		High ICI	0.222 (0.003)	0.074 (0.005)	0.176 (0.196)	0.216 (0.002)	0.02 (0.005)	0.243 (0.128)	0.215 (0.002)	0.011 (0.004)	0.359 (0.191)
	10	AUC*	0.206 (0.002)	0.014 (0.005)	0.089 (0.026)	0.206 (0.002)	0.016 (0.005)	0.142 (0.023)	0.206 (0.002)	0.012 (0.005)	0.302 (0.108)
		KL*	0.211 (0.002)	0.028 (0.005)	0.015 (0.005)	0.21 (0.002)	0.015 (0.004)	0.156 (0.02)	0.21 (0.002)	0.012 (0.005)	0.306 (0.089)
		High ICI	0.232 (0.002)	0.106 (0.005)	0.324 (0.279)	0.219 (0.002)	0.021 (0.005)	0.396 (0.125)	0.219 (0.002)	0.011 (0.005)	0.435 (0.206)
	50	AUC*	0.207 (0.002)	0.019 (0.006)	0.127 (0.031)	0.207 (0.002)	0.016 (0.005)	0.145 (0.025)	0.207 (0.002)	0.012 (0.005)	0.292 (0.105)
		KL*	0.215 (0.002)	0.034 (0.006)	0.017 (0.004)	0.213 (0.002)	0.014 (0.004)	0.19 (0.021)	0.214 (0.002)	0.012 (0.004)	0.347 (0.1)
		High ICI	0.238 (0.003)	0.126 (0.006)	0.422 (0.048)	0.221 (0.002)	0.024 (0.005)	0.424 (0.039)	0.22 (0.002)	0.012 (0.005)	0.382 (0.108)
	100	AUC*	0.208 (0.002)	0.021 (0.006)	0.144 (0.034)	0.207 (0.002)	0.015 (0.005)	0.147 (0.024)	0.207 (0.002)	0.012 (0.005)	0.326 (0.104)
		KL*	0.216 (0.002)	0.037 (0.006)	0.017 (0.004)	0.215 (0.002)	0.014 (0.005)	0.202 (0.022)	0.215 (0.002)	0.012 (0.005)	0.368 (0.111)
		High ICI	0.241 (0.003)	0.133 (0.005)	0.486 (0.046)	0.221 (0.002)	0.025 (0.004)	0.471 (0.06)	0.22 (0.002)	0.011 (0.005)	0.399 (0.098)

Notes: AUC*, KL*, High ICI: models selected by optimizing AUC, KL divergence, or by selecting a high ICI.

327 C.2 Real Data

328 We train XGBoost models on the 10 datasets presented in Section B.2. Unlike Section C.1, the true
 329 probabilities underlying the binary events are not observable. Here, we assume that we have prior
 330 knowledge about the probability distribution, which can be considered as expert opinion. To simulate
 331 this expert opinion, we assume that the true probabilities follow a Beta distribution. The parameters
 332 of this distribution, specific to each dataset, are estimated via MLE using the scores from a GAMSEL
 333 model [6].

334 Using these prior distributions, it is possible to replicate the estimation procedure previously applied
 335 to the simulated data. Each dataset is split into two parts: 70% of the observations are used to train
 336 an XGBoost model (on a training set comprising 80% of these observations, with hyperparameters
 337 selected based on metrics calculated on the remaining 20% validation set), and the remaining 30%

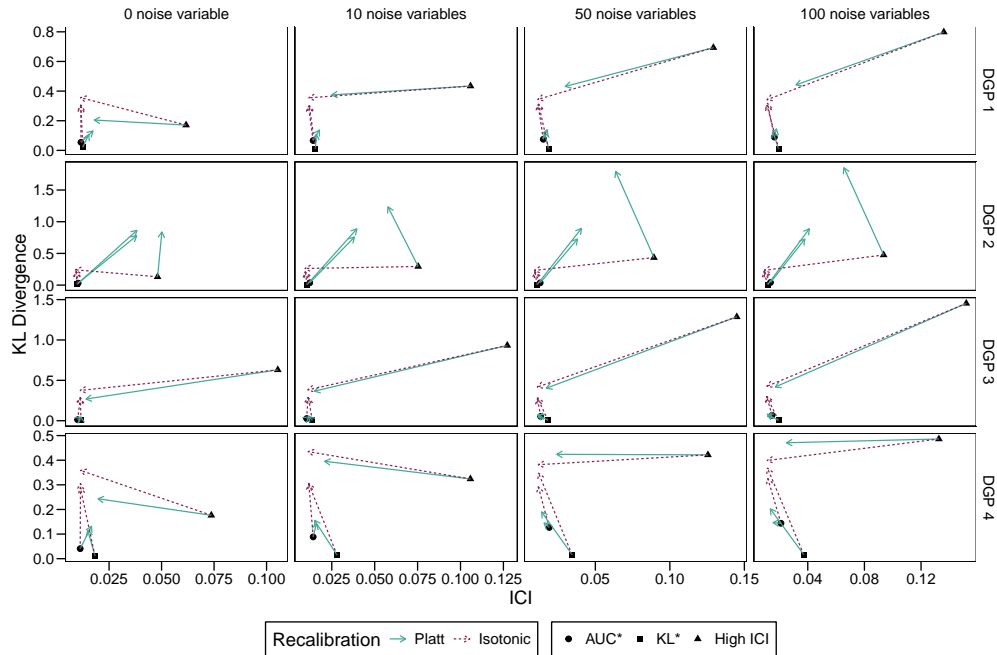


Figure C18: Average KL divergence and ICI before and after recalibration of the estimated scores.

Notes: AUC*, KL*, High ICI: models selected by optimizing AUC, KL divergence, or by selecting a high ICI.

338 are used for model calibration (with 80% of this subset forming the calibration set) and for testing
 339 model performance on unseen data (the remaining 20%).

340 Table C4 presents the metrics calculated on the test set for models selected based on their validation
 341 performance, according to AUC (AUC*), KL divergence between the score distribution and the
 342 prior distribution (KL*), or to intentionally obtain poor calibration (High ICI), both before and
 343 after applying calibration methods. For real datasets, High ICI refers to the algorithm with hyperparameters
 344 yielding the highest AUC among models with an ICI at least one standard deviation above the mean
 345 ICI observed during grid search. This table complements Fig. 2 from the main part of the article.

Table C4: Comparison of metrics computed on the validation set for models selected based on AUC, KL divergence, or ICI. Standard errors are provided in parentheses.

Dataset	Optim.	None			Platt			Isotonic		
		BS	ICI	KL	BS	ICI	KL	BS	ICI	KL
abalone	AUC*	0.218	0.081	0.623	0.208	0.051	0.197	0.208	0.047	1.541
	KL*	0.212	0.093	0.319	0.204	0.042	0.163	0.205	0.038	0.887
	High ICI	0.266	0.203	5.861	0.209	0.046	0.689	0.208	0.039	1.174
adult	AUC*	0.086	0.014	0.335	0.088	0.040	0.446	0.086	0.016	0.416
	KL*	0.097	0.045	0.105	0.096	0.035	0.333	0.096	0.022	0.370
	High ICI	0.096	0.056	0.234	0.094	0.037	0.416	0.093	0.017	0.441
bank	AUC*	0.066	0.017	0.499	0.070	0.046	0.674	0.066	0.009	0.535
	KL*	0.077	0.036	0.091	0.079	0.043	0.446	0.075	0.008	0.417
	High ICI	0.078	0.053	1.193	0.078	0.052	0.679	0.069	0.010	0.427
default	AUC*	0.132	0.028	0.309	0.133	0.034	0.756	0.132	0.010	0.490
	KL*	0.133	0.024	0.283	0.133	0.028	0.695	0.132	0.009	0.491
	High ICI	0.151	0.122	1.176	0.135	0.036	1.387	0.134	0.011	0.781
drybean	AUC*	0.034	0.016	0.764	0.037	0.026	0.865	0.033	0.008	0.721
	KL*	0.037	0.028	0.392	0.040	0.031	0.826	0.036	0.008	0.692
	High ICI	0.037	0.054	2.172	0.036	0.027	0.769	0.034	0.008	0.650
coupon	AUC*	0.202	0.144	3.849	0.180	0.023	1.300	0.178	0.026	1.106
	KL*	0.196	0.022	0.052	0.196	0.014	0.191	0.197	0.009	0.478
	High ICI	0.202	0.144	3.849	0.180	0.023	1.300	0.178	0.026	1.106
mushroom	AUC*	0.000	0.003	1.384	0.000	0.002	1.403	0.000	0.002	1.403
	KL*	0.013	0.050	0.620	0.011	0.021	1.222	0.005	0.003	1.337
	High ICI	0.006	0.037	0.928	0.001	0.008	1.348	0.000	0.003	1.403
occupancy	AUC*	0.009	0.008	1.084	0.009	0.008	1.164	0.009	0.009	1.044
	KL*	0.011	0.041	0.862	0.009	0.006	1.128	0.009	0.006	1.032
	High ICI	0.011	0.041	0.862	0.009	0.006	1.128	0.009	0.006	1.032
winequality	AUC*	0.178	0.131	4.244	0.162	0.030	1.697	0.157	0.019	1.253
	KL*	0.173	0.029	0.146	0.174	0.042	0.439	0.171	0.028	1.109
	High ICI	0.172	0.134	4.609	0.157	0.050	1.797	0.151	0.018	1.288
spambase	AUC*	0.044	0.012	0.699	0.045	0.018	1.107	0.045	0.012	1.044
	KL*	0.059	0.045	0.260	0.057	0.012	0.647	0.056	0.009	0.607
	High ICI	0.061	0.047	0.297	0.058	0.012	0.626	0.058	0.008	0.858

Notes: AUC*, KL*, High ICI: models selected by optimizing AUC, KL divergence, or by selecting a high ICI.