

# Inductive Venn–Abers and related regressors

Ivan Petej and Vladimir Vovk



практические выводы  
теории вероятностей  
могут быть обоснованы  
в качестве следствий  
гипотез о *предельной*  
при данных ограничениях  
сложности изучаемых явлений

**On-line Compression Modelling Project (New Series)**

Working Paper #47

May 7, 2026

Project web site:  
<http://alrw.net>

# Abstract

Venn–Abers predictors are probabilistic predictors that enjoy appealing properties of validity, but their major limitation is that they are applicable only to the case of binary classification, with a recent extension to bounded regression. We generalize them to the case of unbounded regression, which requires adding an element of conformal prediction. In our simulation and empirical studies we investigate the predictive efficiency of point regressors derived from Venn–Abers regressors and argue that they somewhat improve the predictive efficiency of standard regressors for larger training sets.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Comparisons with Literature</b>	<b>2</b>
<b>3</b>	<b>Inductive Venn–Abers Regressors</b>	<b>3</b>
<b>4</b>	<b>Validity</b>	<b>5</b>
<b>5</b>	<b>Algorithm</b>	<b>6</b>
<b>6</b>	<b>Merging a Regression Interval into a Single Value</b>	<b>7</b>
<b>7</b>	<b>Cross Venn–Abers Regressors</b>	<b>8</b>
<b>8</b>	<b>Experimental Results</b>	<b>8</b>
<b>9</b>	<b>Conclusion</b>	<b>16</b>
	<b>References</b>	<b>16</b>
<b>A</b>	<b>Proofs</b>	<b>18</b>
<b>B</b>	<b>Further Experimental Results</b>	<b>20</b>
<b>C</b>	<b>Code Availability</b>	<b>32</b>

# 1 Introduction

This paper explores the problem of regression estimation as presented in, e.g., Vapnik (1998, Sect. 1.4). In the standard setting of statistical learning, where we observe an IID sequence of random pairs  $(X, Y)$  consisting of objects  $X$  and their labels  $Y \in \mathbb{R}$ , the (ideal) *regression estimator* (of  $Y$  given  $X$ ) maps each object  $x$  to the expected value  $\mathbb{E}(Y \mid X = x)$  of its label  $Y$  conditional on observing  $x$ . Vapnik lists regression estimation as one of three basic statistical problems (Vapnik, 1998, Sect. 1.2). We develop ideas of conformal prediction (Angelopoulos et al., 2026; Vovk et al., 2022) and introduce an algorithm for regression estimation that satisfies a natural distribution-free notion of validity.

We are interested in the distribution-free setting, when nothing is known about the distribution generating one pair  $(X, Y)$  and we are only given a training sequence of labelled objects and an unlabelled test object. In typical cases we cannot hope to find the true regression estimate  $\mathbb{E}(Y \mid X = x)$  (even when it is well-defined), and so we lower the bar in three respects in our definition of a valid regression estimator. First, we allow estimators of the form  $\mathbb{E}(Y \mid \mathcal{F})$  for some  $\sigma$ -algebra  $\mathcal{F}$ ; we will express this by saying that our regression estimator is auto-calibrated. Ideally,  $\mathcal{F}$  should be close to the  $\sigma$ -algebra generated by the training set and test object. Second, we allow our algorithms to output intervals (ideally short “imprecise regression estimates”) containing  $\mathbb{E}(Y \mid \mathcal{F})$ . And third, we replace  $\mathbb{E}(Y \mid \mathcal{F})$  by  $\mathbb{E}(Y' \mid \mathcal{F})$ , where  $Y'$  is a “regularized” version of  $Y$ , but  $Y' = Y$  with high probability. (This is the element of conformal prediction that we mentioned in the abstract; it is required in the absence of bounds on the test label.)

As usual, we consider two principal requirements for our algorithms, validity (technically, auto-calibration) and efficiency. Validity (in this technical sense) will be guaranteed for our algorithms, but to achieve efficiency we will use existing regression algorithms that we believe to be efficient, although perhaps miscalibrated. We will develop methods for improving their calibration. These methods will be adaptations of the methods used in Vovk et al. (2015) (and presented in a much greater detail in Vovk et al. 2022, Chap. 6) in the context of binary classification, with bounded regression considered earlier by van der Laan and Alaa (2024).

The problem of regression estimation considered in this paper is very different from the regression problems in conformal prediction, where the task is to output a prediction region (typically a prediction interval) for a test label with a pre-specified coverage probability. In regression estimation the task is to cover the expected label  $\mathbb{E}(Y \mid \mathcal{F})$  rather than the label  $Y$  itself. This will allow us to produce much shorter intervals (especially that we replace  $Y$  by  $Y'$  that is not guaranteed to coincide with  $Y$ ).

The goal of the computational experiments reported in this paper is to demonstrate that application of our methods leads to an improvement in the performance of standard point regressors (we consider those implemented in `scikit-learn`). The improvement is not as significant as we had hoped and disappears for smaller datasets.

## 2 Comparisons with Literature

Results of this paper are not directly comparable with the existing literature because of a somewhat unusual notion of validity that we use. But there are several related approaches.

A strand of research that is closest to what we do in this paper is conformal regression, already mentioned in the previous section. Important advances in conformal regression include Romano et al. (2019), offering very flexible methods based on quantile regression, and Gibbs et al. (2025), establishing conditional validity results. See, e.g., Angelopoulos et al. (2026) for a recent review of conformal prediction in general and conformal regression in particular.

The goal of conformal regression is to produce provably valid, under the assumption of IID data, prediction intervals for the label of a test object. Since validity here means a guaranteed coverage probability, this is a much more ambitious problem than regression estimation dealt with in this paper, as mentioned earlier.

Conformal predictive distributions (Vovk et al., 2022, Chap. 7) are different from conformal regression in that, instead of prediction intervals, they output full predictive distributions for future labels (assumed to be real-valued, as in conformal regression). To achieve validity these predictive distributions should be imprecise in a certain sense. Our current task is easier in that we only aim to cover  $\mathbb{E}(Y' | \mathcal{F})$ , but we will achieve another property of validity, namely auto-calibration instead of calibration in probability (which is stronger in the binary case but not comparable in general). See Allen et al. (2025) for a recent analysis of achievable properties of validity for conformal predictive distributions.

The property of validity used in this paper is inherited from Venn prediction, which is, however, typically applied to produce imprecise probability forecasts in classification problems. Venn prediction (including Venn–Abers prediction) is reviewed in Vovk et al. (2022, Chap. 6) and Venn–Abers prediction is reviewed in Angelopoulos et al. (2026, Sect. 12.4).

Venn–Abers predictors were extended to the case of bounded regression by van der Laan and Alaa (2024), as mentioned earlier. Our bounded Venn–Abers regressor in the next section is just an inductive version of van der Laan and Alaa (2024, Algorithm 1). As the next step, van der Laan and Alaa develop an algorithm (van der Laan and Alaa, 2024, Algorithm 2) for “self-calibrating conformal prediction” based on their Algorithm 1. Finally, van der Laan and Alaa (2025) extend these results to general loss functions interpreting the usual regression setting as the case of squared error loss.

A key step in our algorithm is replacing the true labels  $Y$  by their regularized versions  $Y'$ , as discussed in the previous section. We need it to avoid the vacuous regression interval  $(-\infty, \infty)$  (see the end of Sect. 3). Assumptions made by van der Laan and Alaa (2024) and van der Laan and Alaa (2025) allow them to bypass this step, but we avoid making any assumptions on the data-generating distribution apart from the observations being IID.

In general, the properties of validity of our procedures are sufficiently different from the properties of validity considered in literature to make direct com-

parison between them in empirical or simulation studies difficult or impossible. Therefore, in our experimental section (Sect. 8) we concentrate on evaluating the predictive efficiency of point regressors derived from imprecise regressors implementing our methods; namely, we compare the predictive performance of those point regressors with that of the base algorithms.

### 3 Inductive Venn–Abers Regressors

In this section we ignore the computational complexity of our methods (it will be the topic of Sect. 5). Fix a measurable space  $\mathbf{X}$  (the *object space*) and set  $\mathbf{Z} := \mathbf{X} \times \mathbb{R}$  (this is the *example space*). Each example  $z = (x, y) \in \mathbf{Z}$  consists of an object  $x \in \mathbf{X}$  and its real-valued label  $y$ .

We consider two settings (leading to formally different, albeit similar, prediction algorithms). In the setting of *bounded regression* we are given a finite interval  $[C_*, C^*] \subseteq \mathbb{R}$  guaranteed to contain all labels, training and test. Otherwise, we have *unbounded regression*. In the former case, the algorithm and its property of validity are simpler, and we consider them separately even though we are not particularly interested in them per se as we would like to avoid any assumptions apart from IID observations.

Notice that an algorithm for unbounded regression can also be applied in the bounded setting, and it may well be preferable if the bounds  $C_*$  and  $C^*$  are loose. This is another reason why our main interest is in unbounded regression.

#### 3.1 Bounded regression

Fix any regression algorithm (such as a neural net) producing point predictions as the *base algorithm*. We are given a training set consisting of  $l > k$  examples (where  $k \in \mathbb{N} := \{1, 2, \dots\}$  is typically large and  $l - k \in \mathbb{N}$  is also large) and a test object  $x \in \mathbf{X}$ ;  $k$  is a parameter of our algorithm.

In bounded regression, we are given an interval  $[C_*, C^*]$  guaranteed to contain all labels. The corresponding *bounded inductive Venn–Abers regressor* (bounded IVAR) produces the *regression interval*

$$[\hat{y}_*, \hat{y}^*] := [f_*(r), f^*(r)] \tag{1}$$

for the test label, where  $f^*$ ,  $f_*$ , and  $r$  are defined as follows:

1. Randomly split the training set of size  $l > k$  into two parts, a *proper training set* of size  $l - k$  and a *calibration set*  $z_1, \dots, z_k$  of size  $k$ ; for each  $i \in \{1, \dots, k\}$ ,  $z_i = (x_i, y_i)$  consists of an object  $x_i$  and its label  $y_i$ .
2. Train the regression algorithm on the proper training set obtaining a *prediction rule*  $R : \mathbf{X} \rightarrow \mathbb{R}$  (a measurable function) mapping objects to their predicted labels.
3. Find the prediction  $r_i := R(x_i)$  (*base prediction*) for the calibration object  $x_i$ ,  $i = 1, \dots, k$ , and the base prediction  $r := R(x)$  for the test object  $x$ .

4. Fit isotonic regression to  $(r_1, y_1), \dots, (r_k, y_k), (r, C^*)$  obtaining an isotonic calibrator  $f^*$ .
5. Set  $\hat{y}^* := f^*(r)$ .
6. Fit isotonic regression to  $(r_1, y_1), \dots, (r_k, y_k), (r, C_*)$  obtaining an isotonic calibrator  $f_*$ .
7. Set  $\hat{y}_* := f_*(r)$ .

The bounded IVAR is a simple modification of the inductive Venn–Abers predictor as defined in Vovk et al. (2015). First, we relax the condition  $y_i \in \{0, 1\}$  of binary labels by allowing the labels to take intermediate values,  $y_i \in [0, 1]$ , and then we scale the interval  $[0, 1]$  to arbitrary  $[C_*, C^*]$ . As mentioned in the previous section, non-inductive Venn–Abers regressors were introduced by van der Laan and Alaa (2024).

### 3.2 Unbounded regression

We again fix a regression algorithm and are given a training set of size  $l > k$ . Another parameter of our algorithm is  $m \in \{1, \dots, \lfloor (k-1)/2 \rfloor\}$ ; we are interested in a small  $m$ , such as 1. The corresponding *inductive Venn–Abers regressor* (IVAR) produces the following regression interval of the form (1) for the test label:

1. Randomly split the training set of size  $l > k$  into two parts, a *proper training set* of size  $l - k$  and a *calibration set*  $z_1, \dots, z_k$  of size  $k$ , as before.
2. Train the regression algorithm on the proper training set obtaining a prediction rule  $R : \mathbf{X} \rightarrow \mathbb{R}$ .
3. Find the base predictions  $r_i := R(x_i)$ ,  $i = 1, \dots, k$ , and  $r := R(x)$ .
4. Replace the  $m$  smallest calibration labels  $y_i$  by the  $(m + 1)$ th smallest calibration label  $y_*$  and replace the  $m - 1$  largest calibration labels  $y_i$  by the  $m$ th largest calibration label  $y^*$ . (Notice the asymmetry in the definitions of  $y_*$  and  $y^*$ .) In other words, let the new calibration labels be

$$y'_i := \begin{cases} y_* & \text{if } y_i < y_* \\ y^* & \text{if } y_i > y^* \\ y_i & \text{otherwise.} \end{cases} \quad (2)$$

(Notice that the recipe is still unambiguous when there are ties among  $y_i$ .)

5. Fit isotonic regression to  $(r_1, y'_1), \dots, (r_k, y'_k), (r, y^*)$  obtaining an isotonic calibrator  $f^*$ .
6. Set  $\hat{y}^* := f^*(r)$ .

7. Replace the  $m$  largest calibration labels  $y_i$  by the  $(m + 1)$ th largest calibration label  $y^*$  and replace the  $m - 1$  smallest calibration labels  $y_i$  by the  $m$ th smallest calibration label  $y_*$ . (Notice that these  $y_*$  and  $y^*$  are different from those in item 4.) In other words, let the new calibration labels be defined as (2) for the new  $y_*$  and  $y^*$ .
8. Fit isotonic regression to  $(r_1, y'_1), \dots, (r_k, y'_k), (r, y_*)$  obtaining an isotonic calibrator  $f_*$ .
9. Set  $\hat{y}_* := f_*(r)$ .

When fitting isotonic regression in steps 5 and 8 (and the analogous steps, 4 and 6, in the bounded IVAR), we always use the standard pool-adjacent-violators algorithm (PAVA; see, e.g., Barlow et al. 1972, Sect. 1.2). Notice that the new steps 4 and 7 in the IVAR as compared with the bounded IVAR are really needed: if we just set  $y^* := \infty$  and  $y_* := -\infty$ , the PAVA will produce  $(-\infty, \infty)$  as the regression interval.

## 4 Validity

We often write  $Z_1, \dots, Z_k$ , where  $Z_i = (X_i, Y_i)$ , for calibration examples and  $(X, Y)$  for the test example, in order to emphasize that they are considered as random elements.

First we state a simpler property of validity, the one for bounded regression (including binary classification as special case). So we assume that the labels take values in  $[C_*, C^*]$ . Let us say that a random variable  $S$  is *auto-calibrated* as regression estimate of  $Y$  if  $S = \mathbb{E}(Y \mid S)$  a.s. (This is often used in the case of binary classification and was referred to as perfect calibration in, e.g., Vovk et al. 2022, Sect. 6.2.1, van der Laan and Alaa 2024, and van der Laan and Alaa 2025; for a general definition, see, e.g., Krüger and Ziegel 2021, Definition 3.1.)

*Remark 1.* An equivalent definition of auto-calibration is that a random variable  $S$  is said to be auto-calibrated as regression estimate of  $Y$  if  $S = \mathbb{E}(Y \mid \mathcal{F})$  a.s. for some  $\sigma$ -algebra  $\mathcal{F}$  (in which case we may also say that  $S$  is *ideal* relative to  $\mathcal{F}$ ).

A *selector* for an IVAR is a random variable that always belongs to the regression interval output by the IVAR.

**Theorem 2.** *For any bounded IVAR, there is a selector  $S$  that is auto-calibrated for the test label, i.e.,  $\mathbb{E}(Y \mid S) = S$  a.s.*

Our interpretation of Theorem 2 is that the regression interval produced by the bounded IVAR is approximately auto-calibrated provided it is narrow enough.

Now we consider the case of unbounded regression and continue to assume  $2m < k$ . A *selector* is defined as before, and the *Winsorized* test label  $Y$  is

defined by

$$Y' := \begin{cases} Y_{(m)} & \text{if } Y < Y_{(m)} \\ Y_{(k-m+1)} & \text{if } Y > Y_{(k-m+1)} \\ Y & \text{otherwise,} \end{cases} \quad (3)$$

where  $Y_{(1)} \leq \dots \leq Y_{(k)}$  is the sequence  $Y_1, \dots, Y_k$  of calibration labels sorted in the ascending order. We can regard (3) as a more feasible version of  $Y$  corrected for the possibility of  $Y$  being an outlier; we make it easier to predict by restricting it to be in the range of the calibration labels.

*Remark 3.* See Dixon (1960, Sect. 1) and Tukey (1962, Sect. 14) for the original definition of Winsorization, which we use in this paper. This definition is sometimes modified, as in, e.g., Wilcox (2012, Sect. 2.2.2). While the original definition is about moderating a dataset, the modified definition is about moderating the population.

**Theorem 4.** *We have  $Y = Y'$  with probability at least  $1 - \frac{2m}{k+1}$ . For any IVAR, there is a selector  $S$  that is auto-calibrated for the Winsorized test label  $Y'$ , i.e.,  $\mathbb{E}(Y' | S) = S$ .*

The obvious first statement of Theorem 4 says that the Winsorized label  $Y'$  is the same as the original label  $Y$  with high probability assuming  $m \ll k$ . The second statement asserts the validity of the IVAR as regression function for the Winsorized label. See Appendix A for the proofs.

An alternative parametrization of the IVAR is in terms of  $\epsilon := 2m/(k+1)$ ; we then have  $Y = Y'$  with probability at least  $1 - \epsilon$ . If a given  $\epsilon \in [2/(k+1), 1)$  is not of the form  $2m/(k+1)$  for an integer  $m$ , we decrease it as little as possible so that it takes this form.

## 5 Algorithm

In our description of the algorithm, we will just refer to Vovk et al. (2022, Chap. 6, especially Sect. 6.5.3) (which in turn follows Vovk et al. 2015). We will discuss computing  $f^*$  and  $f_*$ , which are denoted by  $f^1$  and  $f^0$ , respectively, in Vovk et al. (2022, Sect. 6.5.3). Remember that we moderate the original labels  $y_i$  to their less extreme versions  $y'_i$ ; after that, we forget the original labels and drop the primes. But it should always be remembered that  $y_i$  stand for the moderated labels.

First we see how to compute  $F^1$ , which determines  $f^1 = f^*$ . We are given the base predictions  $r_i$  and the moderated labels  $y_1, \dots, y_k$ . Define  $k', r'_1, \dots, r'_{k'}, w_1, \dots, w_{k'}$ , and  $y'_1, \dots, y'_{k'}$  as in Vovk et al. (2022, Sect. 6.5.3) (where base predictions were called scores and were denoted by  $s_i$ ). The CSD consisting of the points  $P_i$ ,  $i \in \{0, \dots, k'\}$ , is defined as before, but now it is extended by adding the point  $P_{-1} := (-1, -y^*)$ ; this corresponds to adding the test example to the calibration set assuming that the base prediction for it is smaller than the base prediction for any calibration example while its label is  $y^*$  (which is a most unusual combination). Algorithm 6.3 in Vovk et al. (2022, Sect. 6.5.3)

will compute the corners of the resulting CSD  $P_{-1}, \dots, P_{k'}$  and Algorithm 6.4 will do the rest.

Computing  $F^0$ , which determines  $f^0 = f_*$ , is analogous. Now the CSD consisting of the points  $P_i$ ,  $i \in \{0, \dots, k'\}$ , is extended by adding the point  $P_{k'+1} := P_{k'} + (1, y_*)$ ; this corresponds to adding the test example to the calibration set assuming that its base prediction is larger than any calibration base prediction while its label is  $y_*$  (another most unusual combination). Algorithm 6.5 in Vovk et al. (2022, Sect. 6.5.3) will compute the corners of the resulting CSD  $P_0, \dots, P_{k'+1}$  and Algorithm 6.6 will do the rest.

Algorithm 6.7 in Vovk et al. (2022, Sect. 6.5.3) shows how to arrange the calibration set into a convenient binary search tree, and Algorithm 6.8 shows how to use this tree for computationally efficient prediction. The computational complexity of these procedures is summarized in the following theorem.

**Theorem 5.** *The computation time of our prediction algorithm is  $O(k \log k)$  for preprocessing ( $O(k)$  apart from sorting the calibration set). Once preprocessing is completed, processing each test object can be done in time  $O(\log k)$ .*

## 6 Merging a Regression Interval into a Single Value

The IVAR introduced in the previous section achieves our goal of producing provably valid regression intervals that have a potential to be predictively efficient for efficient base algorithms. However, in order to be able to compare the predictive efficiency of our methods with traditional regression algorithms, in this and following section we will define natural modifications of the IVAR that produce point predictions. In this section we see how to replace the regression intervals output by IVARs by point predictions, and in the following section we will see how we can combine several IVARs to achieve both predictive and computational efficiency.

Given  $y_* < \hat{y}_* < \hat{y}^* < y^*$ , let us see how to replace the regression interval  $[\hat{y}_*, \hat{y}^*]$  with a single regression value  $\hat{y}$ . Following the minimax approach of Vovk et al. (2022, Sect. 6.4.3), we need to solve the equation

$$(\hat{y} - y_*)^2 - (\hat{y}_* - y_*)^2 = (y^* - \hat{y})^2 - (y^* - \hat{y}^*)^2. \quad (4)$$

It is clear that there is a unique solution  $\hat{y}$  in the interval  $[\hat{y}_*, \hat{y}^*]$ , since, over that interval, the left-hand side of (4) increases in  $\hat{y}$  from 0 to a positive number, while the right-hand side decreases from a positive number to 0.

After simplification (4) becomes a linear equation in  $\hat{y}$  with solution

$$\hat{y} = \frac{\hat{y}_*^2 - \hat{y}^{*2} + 2\hat{y}^*y^* - 2\hat{y}_*y_*}{2(y^* - y_*)}, \quad (5)$$

and there is only one solution (both in the interval  $[\hat{y}_*, \hat{y}^*]$  and overall).

We can also solve (4) approximately obtaining a much more intuitive expression. Rewriting (4) as

$$2(\hat{y} - \hat{y}_*)(\hat{y}_* - y_*) + (\hat{y} - \hat{y}_*)^2 = 2(\hat{y}^* - \hat{y})(y^* - \hat{y}^*) + (\hat{y}^* - \hat{y})^2,$$

assuming  $\hat{y}_* \approx \hat{y}^*$ , and ignoring the quadratic terms, we obtain the approximate equation

$$(\hat{y} - \hat{y}_*)(\hat{y}_* - y_*) = (\hat{y}^* - \hat{y})(y^* - \hat{y}^*).$$

Regrouping its terms as

$$\hat{y}(\hat{y}_* - y_* + y^* - \hat{y}^*) = \hat{y}_*(\hat{y}_* - y_*) + \hat{y}^*(y^* - \hat{y}^*),$$

we can write its solution as the weighted average

$$\hat{y} = \frac{\hat{y}_* - y_*}{\hat{y}_* - y_* + y^* - \hat{y}^*} \hat{y}_* + \frac{y^* - \hat{y}^*}{\hat{y}_* - y_* + y^* - \hat{y}^*} \hat{y}^* \quad (6)$$

of  $\hat{y}_*$  and  $\hat{y}^*$ .

## 7 Cross Venn–Abers Regressors

We define *cross Venn–Abers regressors* (CVARs) as in Vovk et al. (2022, Sect. 6.4.4) except that the regression interval coming from each fold is replaced by the regression estimate computed as described in Sect. 6, either as in (5) or as in (6) (in our experiments in the next section we use the latter). Namely, we divide the training set into  $K$  folds of approximately equal size, and use each fold in turn as the calibration set while the union of the remaining folds is used as the proper training set. The overall regression estimate is found as the arithmetic mean of the  $K$  regression estimates obtained by merging the regression intervals coming from the  $K$  IVARs corresponding to the  $K$  folds. In our experiments we set the parameter  $K$  to 10.

## 8 Experimental Results

As discussed at the end of Sect. 2, in our experimental results we only compare our point regressors, namely the CVARs as defined in the previous section. These point regressors no longer satisfy any formal properties of validity, but the hope is that the validity properties of IVARs will show in superior performance of CVARs as compared with standard point regressors and measured in a traditional way.

We evaluate our approach on a suite of controlled synthetic regression benchmarks designed to isolate different statistical challenges. We generate datasets with  $n = 10000$  examples under the following generative scenarios. Each dataset consists of  $n$  examples whose objects are vectors of  $d = 10$  features. As before, the objects are denoted by  $x_i$  and the corresponding labels by  $y_i$ ,  $i = 1, \dots, n$ . Across all datasets, the level of noise in the generated labels is parametrized by

$\sigma \in \{1, 3\}$ . The results in this section represent the more challenging random noise level  $\sigma = 3$ , and results for  $\sigma = 1$ , as well as for  $n = 1000$ , are included in Appendix B. Unless stated otherwise, all random draws are independent.

**Bounded Logistic dataset** Each object is generated as  $x_i \sim \mathcal{N}(0, I_{10})$ , the weight vector as  $w \sim \mathcal{N}(0, I_{10})$ , and the Gaussian noise variables as  $\xi_i \sim \mathcal{N}(0, \sigma^2)$ . (We parametrize the Gaussian distribution by its mean and variance, or its covariance matrix in the multidimensional case.) The labels are then generated using the sigmoid function:

$$y_i := \frac{10}{1 + e^{-w^\top x_i}} + \xi_i.$$

Notice that while the conditional expectation of  $y_i$  is bounded (belongs to  $(0, 10]$ ), the labels are not bounded because of the Gaussian noise. Since  $y_i$  are “almost bounded” (the Gaussian distribution having thin tails), this is the most benign case, close to the setting of Sect. 3.1, and our methods work best for it.

**Linear Gaussian dataset** We generate  $x_i \sim \mathcal{N}(0, I_{10})$ ,  $w \sim \mathcal{N}(0, I_{10})$ , and  $\xi_i \sim \mathcal{N}(0, \sigma^2)$  as before. The labels are then defined by

$$y_i := w^\top x_i + \xi_i.$$

**Nonlinear dataset** Again  $x_i \sim \mathcal{N}(0, I_{10})$  and  $w \sim \mathcal{N}(0, I_{10})$ . The label of  $x_i$  combines a linear contribution with smooth nonlinear components:

$$y_i := w^\top x_i + 2 \sin(x_{i,1}) + \frac{1}{2} x_{i,2}^2 - \cos(2x_{i,3}) + \xi_i,$$

where  $\xi_i \sim \mathcal{N}(0, \sigma^2)$  and  $x_{i,j}$  denotes the  $j$ th feature of  $x_i$ .

**Heteroscedastic noise dataset** We again draw  $x_i \sim \mathcal{N}(0, I_{10})$  and  $w \sim \mathcal{N}(0, I_{10})$ . However, the observation noise now depends on the objects. Specifically,

$$y_i := w^\top x_i + \xi_i, \quad \xi_i \sim \mathcal{N}(0, \sigma_i^2),$$

where the variance increases with the magnitude of the first feature,  $\sigma_i := 0.5\sigma + |x_{i,1}|$ . This creates object-dependent uncertainty.

**Heavy-tailed noise dataset** Here  $x_i \sim \mathcal{N}(0, I_{10})$  and  $w \sim \mathcal{N}(0, I_{10})$  as before, but the additive noise follows a Student- $t$  distribution with low degrees of freedom  $\nu$ :

$$y_i := w^\top x_i + \xi_i, \quad \xi_i \sim t_\nu,$$

where we use  $\nu := 3\sigma$ . This produces infrequent but very large deviations.

**Outlier contamination dataset** We generate  $x_i \sim \mathcal{N}(0, I_{10})$  and  $w \sim \mathcal{N}(0, I_{10})$ , and define

$$y_i := w^\top x_i + \xi_i,$$

where  $\xi_i$  follows a contamination model:

$$\xi_i \sim \begin{cases} \mathcal{N}(0, \sigma^2), & \text{with probability } 1 - p, \\ \mathcal{N}(0, \tau^2), & \text{with probability } p, \end{cases}$$

with  $p := 0.01$  representing the outlier probability and  $\tau \gg 1$  (we set  $\tau := 10\sigma$ ) meaning that the rare errors are extremely large.

**Sparse high-dimensional dataset** We again draw  $x_i \sim \mathcal{N}(0, I_{10})$ , but the true vector of coefficients is sparse. Let  $k$  be the sparsity level (in our experiments we set  $k := 2$ ); we randomly select a support set  $S \subset \{1, \dots, 10\}$  with  $|S| = k$  and define

$$w_j := \begin{cases} \mathcal{N}(0, 1), & j \in S, \\ 0, & j \notin S. \end{cases}$$

The labels follow

$$y_i := w^\top x_i + \xi_i, \quad \xi_i \sim \mathcal{N}(0, \sigma^2).$$

**Covariate shift dataset** Here the training and test distributions differ. The training objects are sampled as  $x_i^{(\text{train})} \sim \mathcal{N}(0, I_{10})$ , while the test objects are drawn from a shifted distribution  $x_i^{(\text{test})} \sim \mathcal{N}(\mu, \Sigma)$  (we set  $\mu$  to the vector of 1s and  $\Sigma := I_{10}$ ). A single  $w \sim \mathcal{N}(0, I_{10})$  generates both training and test labels:

$$y_i := w^\top x_i + \xi_i, \quad \xi_i \sim \mathcal{N}(0, \sigma^2).$$

Thus the conditional distribution  $p(y \mid x)$  is the same across splits, but the marginal object distribution  $p(x)$  differs.

**Friedman’s datasets** We also show results for Friedman’s three synthetic datasets (Friedman, 1991) as numbered by Breiman (Breiman, 1996, Sect. 3 and Appendix B) and implemented in `scikit-learn`. We use the `scikit-learn` implementation with the default values of parameters except for the size (namely, the default number of features is 10 for Friedman 1).

For each dataset we randomly split the data into 80% training and 20% testing, standardize the features based on training statistics, and repeat all experiments across multiple (namely, 100) random seeds to reduce variance. At the end of the section we also briefly discuss results for the other combinations of  $n \in \{1000, 10000\}$  and  $\sigma \in \{1, 3\}$ .

We compare our methods against seven widely used regression baselines implemented in `scikit-learn`, including linear regressors (Linear Regression, Ridge, Lasso, Elastic Net), kernel-based methods (Support Vector Regression, SVR, with the RBF kernel), and tree-based ensembles (Random Forest and

Gradient Boosting). Linear Regression in fact implements Least Squares, and Ridge (or Ridge Regression), Lasso, and Elastic Net add various elements of regularization. All models are trained with default or widely adopted hyperparameters to reflect typical practitioner usage. Performance is measured using root mean squared error (RMSE) on the held-out test set. We report mean RMSE across repeated trials for each scenario to assess their accuracy.

Tables 1–11 show the average RMSE over 100 trials for each synthetic dataset with  $\sigma = 3$  and  $n = 10000$  examples and for each base algorithm without calibration (*none*) or calibrated using our CVAR method with 10 folds and parameters  $m \in \{1, 10\}$  (denoted *CVAR1* and *CVAR10*, respectively, in the tables; using  $m = 100$  never leads to improvements as compared with the smaller  $m$  in our experiments). Apart from the average RMSE we also report the standard error of the mean (SEM) obtained by dividing the sample standard deviation of the RMSE over the 100 trials by  $\sqrt{100} = 10$ ; therefore, each cell in our tables has the form

$$\text{average RMSE} \pm \text{SEM}.$$

These intervals reassure us that most of our comparisons are not unduly affected by randomness, but in our summaries we only use the averages. The smallest average in each row of our tables is shown in boldface. For each table we also report the average of each column, which is the average of the average RMSE over the seven regression algorithms.

*Remark 6.* The results for Linear Regression and Ridge coincide in Tables 1–11, but they are sometimes slightly different in Appendix B. The reason for this closeness is that our features are constructed as independent, so multicollinearity is unlikely, and the sample size is large relative to the feature count; this makes Least Squares stable, and so the default ridge penalty has relatively little impact.

The last rows in Tables 1–11 show that both of our methods, CVAR1 and CVAR10, attain better RMSE scores on average when compared with uncalibrated algorithms, especially for the bounded logistic, linear Gaussian, nonlinear, covariate shift, Friedman 1, and Friedman 2 datasets (the results for the last dataset, however, are very irregular, and in two cases our methods significantly lower the quality of predictions). In many rows the base algorithm performs better, but in a typical row of a typical table either our algorithms perform better or they lose little. For the “almost bounded” Bounded Logistic dataset our methods improve the base predictions in all rows.

We report the results for  $\sigma \in \{1, 3\}$  and  $n = 1000$  and for  $\sigma = 1$  and  $n = 10000$  settings in Appendix B. The results for  $n = 10000$  and  $\sigma = 1$  are similar to the results in this section, while the results for  $n = 1000$  are mixed for  $\sigma = 3$  and weaker for our methods in the case of low noise,  $\sigma = 1$ . It appears that our methods work well for larger datasets.

In Appendix B we also give results for four real-life datasets. Our methods tend to improve the performance of the base algorithms unless the dataset is small, but the tendency is weak.

Table 1: Bounded Logistic dataset

	none	CVAR1	CVAR10
Elastic Net	$3.809 \pm 0.015$	<b><math>3.205 \pm 0.007</math></b>	$3.206 \pm 0.007$
Gradient Boosting	$3.239 \pm 0.010$	<b><math>3.144 \pm 0.006</math></b>	$3.145 \pm 0.006$
Lasso	$4.001 \pm 0.018$	<b><math>3.588 \pm 0.017</math></b>	$3.589 \pm 0.017$
Linear Regression	$3.257 \pm 0.012$	<b><math>3.011 \pm 0.005</math></b>	$3.012 \pm 0.005$
Random Forest	$3.231 \pm 0.008$	<b><math>3.208 \pm 0.007</math></b>	$3.209 \pm 0.007$
Ridge	$3.257 \pm 0.012$	<b><math>3.011 \pm 0.005</math></b>	$3.012 \pm 0.005$
SVR (RBF)	$3.143 \pm 0.007$	<b><math>3.070 \pm 0.005</math></b>	$3.071 \pm 0.005$
average	3.420	<b>3.177</b>	3.178

Table 2: Linear Gaussian dataset

	none	CVAR1	CVAR10
Elastic Net	$3.510 \pm 0.015$	<b><math>3.147 \pm 0.006</math></b>	$3.155 \pm 0.006$
Gradient Boosting	<b><math>3.073 \pm 0.006</math></b>	$3.083 \pm 0.006$	$3.091 \pm 0.006$
Lasso	$3.707 \pm 0.018$	<b><math>3.419 \pm 0.014</math></b>	$3.424 \pm 0.014$
Linear Regression	<b><math>3.002 \pm 0.005</math></b>	$3.019 \pm 0.005$	$3.029 \pm 0.005$
Random Forest	<b><math>3.138 \pm 0.006</math></b>	$3.150 \pm 0.007$	$3.158 \pm 0.007$
Ridge	<b><math>3.002 \pm 0.005</math></b>	$3.019 \pm 0.005$	$3.029 \pm 0.005$
SVR (RBF)	<b><math>3.087 \pm 0.005</math></b>	$3.101 \pm 0.006$	$3.108 \pm 0.006$
average	3.217	<b>3.134</b>	3.142

Table 3: Nonlinear dataset

	none	CVAR1	CVAR10
Elastic Net	$3.718 \pm 0.015$	<b><math>3.327 \pm 0.006</math></b>	$3.336 \pm 0.006$
Gradient Boosting	<b><math>3.100 \pm 0.006</math></b>	$3.112 \pm 0.006$	$3.124 \pm 0.007$
Lasso	$3.904 \pm 0.018$	<b><math>3.590 \pm 0.014</math></b>	$3.596 \pm 0.014$
Linear Regression	<b><math>3.199 \pm 0.005</math></b>	$3.211 \pm 0.006$	$3.221 \pm 0.006$
Random Forest	<b><math>3.194 \pm 0.008</math></b>	$3.210 \pm 0.008$	$3.220 \pm 0.008$
Ridge	<b><math>3.199 \pm 0.005</math></b>	$3.211 \pm 0.006$	$3.221 \pm 0.006$
SVR (RBF)	<b><math>3.139 \pm 0.006</math></b>	$3.157 \pm 0.006$	$3.167 \pm 0.006$
average	3.350	<b>3.260</b>	3.269

Table 4: Heteroscedastic noise dataset

	none	CVAR1	CVAR10
Elastic Net	$4.654 \pm 0.016$	<b>4.386</b> $\pm 0.011$	$4.392 \pm 0.011$
Gradient Boosting	$4.371 \pm 0.012$	<b>4.367</b> $\pm 0.012$	$4.372 \pm 0.012$
Lasso	$4.804 \pm 0.017$	<b>4.587</b> $\pm 0.015$	$4.590 \pm 0.015$
Linear Regression	<b>4.289</b> $\pm 0.011$	$4.295 \pm 0.011$	$4.302 \pm 0.011$
Random Forest	$4.430 \pm 0.012$	<b>4.425</b> $\pm 0.012$	$4.429 \pm 0.012$
Ridge	<b>4.289</b> $\pm 0.011$	$4.295 \pm 0.011$	$4.302 \pm 0.011$
SVR (RBF)	<b>4.358</b> $\pm 0.011$	$4.365 \pm 0.011$	$4.369 \pm 0.012$
average	4.456	<b>4.388</b>	4.394

Table 5: Heavy-tailed noise dataset

	none	CVAR1	CVAR10
Elastic Net	$5.444 \pm 0.076$	$5.218 \pm 0.077$	<b>5.210</b> $\pm 0.077$
Gradient Boosting	$5.207 \pm 0.077$	$5.212 \pm 0.077$	<b>5.204</b> $\pm 0.077$
Lasso	$5.578 \pm 0.075$	$5.407 \pm 0.075$	<b>5.391</b> $\pm 0.075$
Linear Regression	<b>5.120</b> $\pm 0.078$	$5.137 \pm 0.078$	$5.132 \pm 0.078$
Random Forest	$5.280 \pm 0.077$	$5.283 \pm 0.077$	<b>5.267</b> $\pm 0.077$
Ridge	<b>5.120</b> $\pm 0.078$	$5.137 \pm 0.078$	$5.132 \pm 0.078$
SVR (RBF)	<b>5.173</b> $\pm 0.077$	$5.196 \pm 0.077$	$5.184 \pm 0.077$
average	5.275	5.227	<b>5.217</b>

Table 6: Outlier contamination dataset

	none	CVAR1	CVAR10
Elastic Net	$4.674 \pm 0.040$	$4.401 \pm 0.042$	<b>4.400</b> $\pm 0.042$
Gradient Boosting	$4.389 \pm 0.042$	$4.397 \pm 0.042$	<b>4.388</b> $\pm 0.042$
Lasso	$4.826 \pm 0.039$	$4.617 \pm 0.040$	<b>4.604</b> $\pm 0.040$
Linear Regression	<b>4.290</b> $\pm 0.043$	$4.309 \pm 0.043$	$4.311 \pm 0.043$
Random Forest	$4.456 \pm 0.042$	$4.456 \pm 0.041$	<b>4.443</b> $\pm 0.041$
Ridge	<b>4.290</b> $\pm 0.043$	$4.309 \pm 0.043$	$4.311 \pm 0.043$
SVR (RBF)	<b>4.353</b> $\pm 0.042$	$4.378 \pm 0.042$	$4.370 \pm 0.042$
average	4.468	4.409	<b>4.404</b>

Table 7: Sparse high-dimensional signals dataset

	none	CVAR1	CVAR10
Elastic Net	$3.049 \pm 0.007$	<b><math>2.999 \pm 0.005</math></b>	$3.000 \pm 0.005$
Gradient Boosting	$3.009 \pm 0.005$	$3.008 \pm 0.005$	<b><math>3.008 \pm 0.005</math></b>
Lasso	$3.079 \pm 0.009$	<b><math>3.026 \pm 0.006</math></b>	$3.027 \pm 0.006$
Linear Regression	<b><math>2.995 \pm 0.005</math></b>	$2.998 \pm 0.005$	$2.998 \pm 0.005$
Random Forest	$3.046 \pm 0.005$	$3.027 \pm 0.005$	<b><math>3.026 \pm 0.005</math></b>
Ridge	<b><math>2.995 \pm 0.005</math></b>	$2.998 \pm 0.005$	$2.998 \pm 0.005$
SVR (RBF)	$3.038 \pm 0.005$	$3.025 \pm 0.005$	<b><math>3.024 \pm 0.005</math></b>
average	3.030	3.012	<b>3.012</b>

Table 8: Covariate shift dataset

	none	CVAR1	CVAR10
Elastic Net	$3.887 \pm 0.050$	<b><math>3.352 \pm 0.023</math></b>	$3.378 \pm 0.025$
Gradient Boosting	<b><math>3.213 \pm 0.011</math></b>	$3.275 \pm 0.019$	$3.306 \pm 0.023$
Lasso	$4.191 \pm 0.065$	<b><math>3.809 \pm 0.044</math></b>	$3.825 \pm 0.045$
Linear Regression	<b><math>3.007 \pm 0.004</math></b>	$3.109 \pm 0.018$	$3.143 \pm 0.021$
Random Forest	<b><math>3.371 \pm 0.019</math></b>	$3.426 \pm 0.024$	$3.453 \pm 0.027$
Ridge	<b><math>3.007 \pm 0.004</math></b>	$3.109 \pm 0.018$	$3.143 \pm 0.021$
SVR (RBF)	<b><math>3.600 \pm 0.043</math></b>	$3.658 \pm 0.046$	$3.675 \pm 0.048$
average	3.468	<b>3.391</b>	3.417

Table 9: Friedman 1 dataset

	none	CVAR1	CVAR10
Elastic Net	$4.383 \pm 0.008$	<b><math>3.871 \pm 0.007</math></b>	$3.879 \pm 0.007$
Gradient Boosting	<b><math>3.139 \pm 0.006</math></b>	$3.173 \pm 0.006$	$3.190 \pm 0.006$
Lasso	$4.353 \pm 0.008$	<b><math>3.960 \pm 0.007</math></b>	$3.967 \pm 0.007$
Linear Regression	$3.864 \pm 0.007$	<b><math>3.858 \pm 0.007</math></b>	$3.865 \pm 0.007$
Random Forest	<b><math>3.276 \pm 0.006</math></b>	$3.312 \pm 0.006$	$3.328 \pm 0.006$
Ridge	$3.864 \pm 0.007$	<b><math>3.858 \pm 0.007</math></b>	$3.865 \pm 0.007$
SVR (RBF)	<b><math>3.226 \pm 0.006</math></b>	$3.258 \pm 0.006$	$3.273 \pm 0.006$
average	3.729	<b>3.613</b>	3.624

Table 10: Friedman 2 dataset

	none	CVAR1	CVAR10
Elastic Net	$181.783 \pm 0.340$	<b><math>82.684 \pm 0.170</math></b>	$84.122 \pm 0.168$
Gradient Boosting	<b><math>15.432 \pm 0.057</math></b>	$48.841 \pm 0.119$	$51.088 \pm 0.149$
Lasso	$137.838 \pm 0.233$	<b><math>82.678 \pm 0.171</math></b>	$84.126 \pm 0.169$
Linear Regression	$137.820 \pm 0.233$	<b><math>82.720 \pm 0.170</math></b>	$84.167 \pm 0.169$
Random Forest	<b><math>6.550 \pm 0.019</math></b>	$47.315 \pm 0.119$	$49.651 \pm 0.151$
Ridge	$137.820 \pm 0.233$	<b><math>82.720 \pm 0.170</math></b>	$84.167 \pm 0.169$
SVR (RBF)	$140.878 \pm 0.536$	<b><math>105.137 \pm 0.376</math></b>	$105.298 \pm 0.397$
average	108.303	<b>76.013</b>	77.517

Table 11: Friedman 3 dataset

	none	CVAR1	CVAR10
Elastic Net	$3.022 \pm 0.005$	$3.022 \pm 0.005$	<b><math>3.022 \pm 0.005</math></b>
Gradient Boosting	$3.023 \pm 0.005$	$3.020 \pm 0.005$	<b><math>3.018 \pm 0.005</math></b>
Lasso	$3.022 \pm 0.005$	$3.022 \pm 0.005$	<b><math>3.022 \pm 0.005</math></b>
Linear Regression	<b><math>3.013 \pm 0.005</math></b>	$3.016 \pm 0.005$	$3.014 \pm 0.005$
Random Forest	$3.111 \pm 0.005$	$3.026 \pm 0.005$	<b><math>3.024 \pm 0.005</math></b>
Ridge	<b><math>3.013 \pm 0.005</math></b>	$3.016 \pm 0.005$	$3.014 \pm 0.005$
SVR (RBF)	$3.023 \pm 0.005$	$3.020 \pm 0.005$	<b><math>3.018 \pm 0.005</math></b>
average	3.032	3.020	<b>3.019</b>

## 9 Conclusion

This paper presents new validity guarantees for regression estimation. Computational experiments show that this leads to a limited improvement in the performance of standard point regressors on large datasets.

These are some directions of further research.

- An alternative merging procedure to the ones described in Sect. 7 is to try and merge the regression intervals coming from the  $K$  folds directly (in a minimax manner, as in Vovk et al. 2022, Sect. 6.4.5) in order to obtain an overall regression estimate, without the intermediate step of merging each regression interval. It would be interesting to compare it experimentally with the procedure used in this paper.
- In this paper we explore the predictive efficiency of the CVAR in simulation and empirical studies. Alternatively, we could try and prove theoretical results along the lines of the Burnaev–Wasserman programme, as described in Vovk et al. (2022, Sects. 2.5 and 2.9.7). As a first step, we may assume that the base predictions  $r$  coincide with the true regression function,  $r := \mathbb{E}(Y \mid X = x)$ .
- Our methods tend to improve significantly the quality of predictions made by Lasso and Elastic Net. A very common pattern of the boldface (i.e., best) entries in our tables is “CVAR1, none, CVAR1, none, none, none, none” (from top to bottom); we can observe it in three tables (Tables 2, 3, and 8) in the main paper, and in more than half of the tables (20 out of 37) in Appendix B. For this pattern our methods help only for Lasso and Elastic Net. However, it can be argued that applying our methods destroys, at least partially, a valued property of Lasso and Elastic Net, automatic feature selection. Optimizing the number of features is another interesting desideratum, alongside predictive efficiency.

## Acknowledgements

We are grateful to anonymous referees for their advice.

## References

- Sam Allen, Georgios Gavrilopoulos, Alexander Henzi, Gian-Reto Kleger, and Johanna Ziegel. In-sample calibration yields conformal calibration guarantees. Technical Report arXiv:2503.03841 [stat.ME], arXiv.org e-Print archive, March 2025.
- Anastasios N. Angelopoulos, Rina Foygel Barber, and Stephen Bates. Theoretical foundations of conformal prediction. Technical Report arXiv:2411.11824 [math.ST], arXiv.org e-Print archive, March 2026. Pre-publication version of a book to be published by Cambridge University Press.

- Richard E. Barlow, D. J. Bartholomew, J. M. Bremner, and H. Daniel Brunk. *Statistical Inference under Order Restrictions: The Theory and Application of Isotonic Regression*. Wiley, London, 1972.
- Leo Breiman. Bagging predictors. *Machine Learning*, 24:123–140, 1996.
- Thomas Brooks, D. Pope, and Michael Marcolini. Airfoil self-noise [dataset]. UCI Machine Learning Repository, 1989. DOI: <https://doi.org/10.24432/C5VW2C>.
- Dongjin Cho, Cheolhee Yoo, Jungho Im, and Dong-Hyun Cha. Bias correction of numerical prediction model temperature forecast [dataset]. UCI Machine Learning Repository, 2020. DOI: <https://doi.org/10.24432/C59K76>.
- Paulo Cortez. Student performance [dataset]. UCI Machine Learning Repository, 2008. DOI: <https://doi.org/10.24432/C5TG7T>.
- Paulo Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis. Wine quality [dataset]. UCI Machine Learning Repository, 2009. DOI: <https://doi.org/10.24432/C56S3T>.
- Wilfrid J. Dixon. Simplified estimation from censored normal samples. *Annals of Mathematical Statistics*, 31:385–391, 1960.
- Jerome H. Friedman. Multivariate adaptive regression splines (with discussion). *Annals of Statistics*, 19:1–141, 1991.
- Isaac Gibbs, John J. Cherian, and Emmanuel J. Candès. Conformal prediction with conditional guarantees. *Journal of the Royal Statistical Society B*, 87: 1100–1126, 2025.
- Fabian Krüger and Johanna F. Ziegel. Generic conditions for forecast dominance. *Journal of Business and Economic Statistics*, 39:972–983, 2021.
- Ivan Petej. Inductive Venn–Abers and related regressors benchmark experiments. URL <https://github.com/ip200/ivar-experiments> (based on the “Venn–Abers calibration” library <https://github.com/ip200/venn-abers>), May 2026.
- Yaniv Romano, Evan Patterson, and Emmanuel J. Candès. Conformalized quantile regression. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, pages 3543–3553, 2019.
- John W. Tukey. The future of data analysis. *Annals of Mathematical Statistics*, 33:1–67, 1962.
- Lars van der Laan and Ahmed M. Alaa. Self-calibrating conformal prediction. In *NeurIPS*, 2024. Available on OpenReview.

Lars van der Laan and Ahmed M. Alaa. Generalized Venn and Venn–Abers calibration with applications in conformal prediction. In *ICML*, 2025. Available on OpenReview.

Vladimir N. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.

Vladimir Vovk, Ivan Petej, and Valentina Fedorova. Large-scale probabilistic predictors with and without guarantees of validity. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 892–900. Curran Associates, 2015. Full version: arXiv:1511.00213 [cs.LG].

Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random World*. Springer, Cham, second edition, 2022.

Rand R. Wilcox. *Introduction to Robust Estimation and Hypothesis Testing*. Elsevier, Amsterdam, third edition, 2012.

## A Proofs

### A.1 Proof of Theorem 2

The selector  $S$  that we use for demonstrating the theorem is a modification of items 4–7 in the description of the bounded IVAR corresponding to using the true test label. Namely, we fit isotonic regression to  $(r_1, y_1), \dots, (r_k, y_k), (r, y)$ , where  $y$  is the true label of the test object, obtaining an isotonic calibrator  $f$ . Set  $S := f(r)$ .

First we need to check that  $S$  is indeed a selector, namely that  $f_*(r) \leq S \leq f^*(r)$ . This follows from  $C_* \leq y \leq C^*$  and the monotonicity property of PAVA given by the following lemma.

**Lemma 7.** *If  $y_i \leq y'_i$  for all  $i \in \{1, \dots, n\}$ , then  $f \leq f'$ , where  $f$  (resp.  $f'$ ) is the isotonic regression fitted to  $(r_1, y_1), \dots, (r_n, y_n)$  (resp. to  $(r_1, y'_1), \dots, (r_n, y'_n)$ ).*

*Proof.* This is an immediate corollary of the max-min formulas (as given in, e.g., Barlow et al. 1972, p. 19).  $\square$

Now let us check that  $\mathbb{E}(Y \mid S) = S$  even conditionally on a bag  $\{Z_1, \dots, Z_k, Z\}$  of  $k$  calibration examples and a test example under the uniform probability measure on all  $(k + 1)!$  orderings of the bag (see Vovk et al. 2022, Lemma A.3). The value of  $S$  determines the solution block containing the test example. Therefore, both  $S$  and  $\mathbb{E}(Y \mid S)$  will equal the arithmetic mean of the labels in that solution block.

## A.2 Proof of Theorem 4

Fix an IVAR. The selector  $S$  is produced by the following ideal picture. Let  $y$  be the true label of the test object  $x$ . Then the selector  $S$  associated with the IVAR is defined by the following recipe applied after splitting the training set and training the base regression algorithm on the proper training set (which gives us a prediction rule).

1. Winsorize the labels, namely:
  - replace the  $m$  largest labels in the *augmented calibration sequence*

$$(x_1, y_1), \dots, (x_k, y_k), (x, y)$$
 by the  $(m + 1)$ th largest label;
    - replace the  $m$  smallest labels in the augmented calibration sequence by the  $(m + 1)$ th smallest label.
2. Find the base predictions  $r_1, \dots, r_k, r$  of the objects  $x_1, \dots, x_k, x$  using the prediction rule.
3. Fit isotonic regression  $f$  to  $(r_1, y_1), \dots, (r_k, y_k), (r, y)$  (with the Winsorized labels).
4. Set  $S := f(r)$ .

Step 1 regularizes the labels moderating the  $2m$  most extreme ones. The recipe treats all elements of the augmented calibration sequence symmetrically, which is essential in the proof of  $\mathbb{E}(Y' | S) = S$ ; let us call it the *symmetric recipe*.

First let us check that  $S$  is indeed a selector for the IVAR. Consider three cases:

- If  $y$  is one of the  $m$  largest labels in the augmented calibration sequence, the definition of IVAR and the symmetric recipe modify the labels in the same way; in particular, the  $m - 1$  largest calibration labels and the test label are replaced by  $y^*$ , and the  $m$  smallest calibration labels are replaced by  $y_*$ . In this case  $S = \hat{y}^*$ .
- Analogously, if  $y$  is one of the  $m$  smallest labels in the augmented calibration sequence, we have  $S = \hat{y}_*$ .
- Otherwise, the monotonicity property of isotonic regression given in Lemma 7 above implies that we still have  $S \in [\hat{y}_*, \hat{y}^*]$ .

In all three cases,  $S \in [\hat{y}_*, \hat{y}^*]$ , which means that  $S$  is a selector.

*Remark 8.* The first two cases considered in the proof demonstrate the tightness, in some sense, of the regression interval  $[\hat{y}_*, \hat{y}^*]$ .

Table 12: Bounded Logistic dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.551 \pm 0.021$	<b>1.506</b> $\pm 0.012$	$1.508 \pm 0.012$
Gradient Boosting	$1.532 \pm 0.018$	<b>1.323</b> $\pm 0.009$	$1.326 \pm 0.009$
Lasso	$2.830 \pm 0.024$	<b>2.201</b> $\pm 0.025$	$2.203 \pm 0.025$
Linear Regression	$1.596 \pm 0.021$	<b>1.031</b> $\pm 0.002$	$1.035 \pm 0.002$
Random Forest	$1.497 \pm 0.017$	<b>1.463</b> $\pm 0.014$	$1.466 \pm 0.014$
Ridge	$1.596 \pm 0.021$	<b>1.031</b> $\pm 0.002$	$1.035 \pm 0.002$
SVR (RBF)	$1.211 \pm 0.009$	<b>1.093</b> $\pm 0.002$	$1.097 \pm 0.002$
average	1.830	<b>1.378</b>	1.382

It remains to prove  $\mathbb{E}(Y' | S) = S$ . We prove this equality conditionally on the set of all permutations (equiprobable) of a given augmented calibration sequence, where a permutation  $\pi$  of  $\{1, \dots, k+1\}$  makes  $z_{\pi^{-1}(k+1)}$  the test example. The Winsorized test label (3) over the permutations is the same random variable as the test label produced according to the symmetric recipe. It remains to remember that the isotonic regression at the base prediction  $r$  for the test object is the mean of the labels in the augmented calibration sequence in the same solution block as the test example (Barlow et al., 1972, pp. 13–15). Solution blocks are level sets of the isotonic regression, which implies  $S$  being the conditional average of the Winsorized test label given  $S$ .

## B Further Experimental Results

In Tables 12–22 we show the experimental results in the case of  $\sigma = 1$  parallel to those reported in Sect. 8 of the main paper; in particular, we still have  $n = 10000$ . Our comments in the main paper are also applicable here. In particular, it is still true that on average both CVAR1 and CVAR10 produce better results than the base algorithms in all 11 tables. The level of noise does not appear to be a crucial parameter.

Tables 23–33 give results for  $n = 1000$  examples and noise level  $\sigma = 3$ . For these smaller datasets the results are mixed, as we said in the main paper.

In Tables 34–44 we have  $n = 1000$  examples, but the noise level is smaller,  $\sigma = 1$ . The results for our methods appear even worse than for the more challenging case  $\sigma = 3$ ; the base algorithms work best on average in nine cases out of eleven.

Finally, in Tables 45–48 we show results for the following four real-life datasets: Bias correction (in full, Bias correction of numerical prediction model temperature forecast) (Cho et al., 2020) with 7750 examples and 7 features, Wine quality (Cortez et al., 2009) with 4898 examples and 12 features (including

Table 13: Linear Gaussian dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.079 \pm 0.024$	<b><math>1.395 \pm 0.009</math></b>	$1.414 \pm 0.010$
Gradient Boosting	<b><math>1.155 \pm 0.008</math></b>	$1.189 \pm 0.008$	$1.213 \pm 0.009$
Lasso	$2.397 \pm 0.026$	<b><math>1.921 \pm 0.024</math></b>	$1.933 \pm 0.024$
Linear Regression	<b><math>1.001 \pm 0.002</math></b>	$1.084 \pm 0.005$	$1.110 \pm 0.006$
Random Forest	<b><math>1.281 \pm 0.013</math></b>	$1.312 \pm 0.014$	$1.333 \pm 0.014$
Ridge	<b><math>1.001 \pm 0.002</math></b>	$1.084 \pm 0.005$	$1.110 \pm 0.006$
SVR (RBF)	<b><math>1.094 \pm 0.004</math></b>	$1.155 \pm 0.006$	$1.179 \pm 0.007$
average	1.430	<b>1.306</b>	1.327

Table 14: Nonlinear dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.417 \pm 0.021$	<b><math>1.763 \pm 0.007</math></b>	$1.781 \pm 0.008$
Gradient Boosting	<b><math>1.227 \pm 0.009</math></b>	$1.260 \pm 0.009$	$1.289 \pm 0.010$
Lasso	$2.695 \pm 0.024$	<b><math>2.212 \pm 0.020</math></b>	$2.223 \pm 0.020$
Linear Regression	<b><math>1.494 \pm 0.003</math></b>	$1.539 \pm 0.005$	$1.560 \pm 0.005$
Random Forest	<b><math>1.417 \pm 0.015</math></b>	$1.455 \pm 0.015$	$1.479 \pm 0.016$
Ridge	<b><math>1.494 \pm 0.003</math></b>	$1.539 \pm 0.005$	$1.560 \pm 0.005$
SVR (RBF)	<b><math>1.202 \pm 0.004</math></b>	$1.272 \pm 0.006$	$1.301 \pm 0.007$
average	1.707	<b>1.577</b>	1.599

Table 15: Heteroscedastic noise dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.315 \pm 0.022$	<b><math>1.722 \pm 0.008</math></b>	$1.740 \pm 0.009$
Gradient Boosting	<b><math>1.550 \pm 0.007</math></b>	$1.570 \pm 0.007$	$1.591 \pm 0.008$
Lasso	$2.604 \pm 0.025$	<b><math>2.173 \pm 0.021</math></b>	$2.184 \pm 0.021$
Linear Regression	<b><math>1.430 \pm 0.004</math></b>	$1.480 \pm 0.005$	$1.504 \pm 0.006$
Random Forest	<b><math>1.640 \pm 0.010</math></b>	$1.662 \pm 0.011$	$1.681 \pm 0.011$
Ridge	<b><math>1.430 \pm 0.004</math></b>	$1.480 \pm 0.005$	$1.504 \pm 0.006$
SVR (RBF)	<b><math>1.509 \pm 0.005</math></b>	$1.546 \pm 0.006$	$1.566 \pm 0.007$
average	1.782	<b>1.662</b>	1.681

Table 16: Heavy-tailed noise dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.512 \pm 0.029$	<b><math>1.957 \pm 0.025</math></b>	$1.981 \pm 0.025$
Gradient Boosting	<b><math>1.819 \pm 0.026</math></b>	$1.830 \pm 0.026$	$1.856 \pm 0.026$
Lasso	$2.785 \pm 0.030$	<b><math>2.381 \pm 0.029</math></b>	$2.395 \pm 0.029$
Linear Regression	<b><math>1.707 \pm 0.026</math></b>	$1.742 \pm 0.026$	$1.772 \pm 0.026$
Random Forest	<b><math>1.901 \pm 0.026</math></b>	$1.915 \pm 0.026$	$1.938 \pm 0.026$
Ridge	<b><math>1.707 \pm 0.026</math></b>	$1.742 \pm 0.026$	$1.772 \pm 0.026$
SVR (RBF)	<b><math>1.770 \pm 0.025</math></b>	$1.797 \pm 0.025$	$1.822 \pm 0.025$
average	2.028	<b>1.909</b>	1.934

Table 17: Outlier contamination dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.330 \pm 0.022$	<b><math>1.718 \pm 0.014</math></b>	$1.748 \pm 0.014$
Gradient Boosting	<b><math>1.561 \pm 0.013</math></b>	$1.573 \pm 0.014$	$1.607 \pm 0.014$
Lasso	$2.620 \pm 0.025$	<b><math>2.182 \pm 0.022</math></b>	$2.199 \pm 0.022$
Linear Regression	<b><math>1.430 \pm 0.014</math></b>	$1.473 \pm 0.014$	$1.512 \pm 0.014$
Random Forest	<b><math>1.656 \pm 0.015</math></b>	$1.671 \pm 0.015$	$1.700 \pm 0.015$
Ridge	<b><math>1.430 \pm 0.014</math></b>	$1.473 \pm 0.014$	$1.512 \pm 0.014$
SVR (RBF)	<b><math>1.500 \pm 0.014</math></b>	$1.533 \pm 0.014$	$1.567 \pm 0.014$
average	1.790	<b>1.661</b>	1.692

Table 18: Sparse high-dimensional signals dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$1.144 \pm 0.012$	<b><math>1.016 \pm 0.003</math></b>	$1.019 \pm 0.003$
Gradient Boosting	<b><math>1.004 \pm 0.002</math></b>	$1.010 \pm 0.002$	$1.012 \pm 0.002$
Lasso	$1.214 \pm 0.018$	<b><math>1.085 \pm 0.011</math></b>	$1.087 \pm 0.011$
Linear Regression	<b><math>0.998 \pm 0.002</math></b>	$1.005 \pm 0.002$	$1.008 \pm 0.002$
Random Forest	<b><math>1.020 \pm 0.002</math></b>	$1.022 \pm 0.002$	$1.025 \pm 0.003$
Ridge	<b><math>0.998 \pm 0.002</math></b>	$1.005 \pm 0.002$	$1.008 \pm 0.002$
SVR (RBF)	$1.030 \pm 0.002$	<b><math>1.029 \pm 0.003</math></b>	$1.032 \pm 0.003$
average	1.058	<b>1.025</b>	1.027

Table 19: Covariate shift dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.619 \pm 0.068$	<b><math>1.794 \pm 0.040</math></b>	$1.841 \pm 0.043$
Gradient Boosting	<b><math>1.462 \pm 0.024</math></b>	$1.590 \pm 0.037$	$1.646 \pm 0.042$
Lasso	$3.047 \pm 0.083$	<b><math>2.537 \pm 0.063</math></b>	$2.563 \pm 0.064$
Linear Regression	<b><math>1.002 \pm 0.001</math></b>	$1.312 \pm 0.036$	$1.377 \pm 0.041$
Random Forest	<b><math>1.778 \pm 0.036</math></b>	$1.863 \pm 0.042$	$1.908 \pm 0.046$
Ridge	<b><math>1.002 \pm 0.001</math></b>	$1.312 \pm 0.036$	$1.377 \pm 0.041$
SVR (RBF)	<b><math>1.888 \pm 0.052</math></b>	$2.025 \pm 0.060$	$2.058 \pm 0.062$
average	1.828	<b>1.776</b>	1.824

Table 20: Friedman 1 dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$3.344 \pm 0.005$	<b><math>2.649 \pm 0.005</math></b>	$2.658 \pm 0.005$
Gradient Boosting	<b><math>1.307 \pm 0.004</math></b>	$1.418 \pm 0.004$	$1.447 \pm 0.004$
Lasso	$3.305 \pm 0.005$	<b><math>2.775 \pm 0.005</math></b>	$2.783 \pm 0.005$
Linear Regression	$2.629 \pm 0.005$	<b><math>2.628 \pm 0.005</math></b>	$2.637 \pm 0.005$
Random Forest	<b><math>1.498 \pm 0.003</math></b>	$1.603 \pm 0.003$	$1.628 \pm 0.003$
Ridge	$2.629 \pm 0.005$	<b><math>2.628 \pm 0.005</math></b>	$2.637 \pm 0.005$
SVR (RBF)	<b><math>1.309 \pm 0.003</math></b>	$1.428 \pm 0.003$	$1.458 \pm 0.003$
average	2.289	<b>2.161</b>	2.178

Table 21: Friedman 2 dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$181.762 \pm 0.340$	<b><math>82.637 \pm 0.170</math></b>	$84.073 \pm 0.168$
Gradient Boosting	<b><math>15.208 \pm 0.061</math></b>	$48.795 \pm 0.120$	$51.037 \pm 0.150$
Lasso	$137.802 \pm 0.233$	<b><math>82.631 \pm 0.171</math></b>	$84.076 \pm 0.170$
Linear Regression	$137.783 \pm 0.233$	<b><math>82.675 \pm 0.171</math></b>	$84.118 \pm 0.170$
Random Forest	<b><math>5.791 \pm 0.021</math></b>	$47.239 \pm 0.119$	$49.570 \pm 0.152$
Ridge	$137.784 \pm 0.233$	<b><math>82.675 \pm 0.171</math></b>	$84.118 \pm 0.170$
SVR (RBF)	$140.800 \pm 0.538$	<b><math>105.034 \pm 0.376</math></b>	$105.190 \pm 0.398$
average	108.133	<b>75.955</b>	77.455

Table 22: Friedman 3 dataset ( $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$1.050 \pm 0.002$	$1.050 \pm 0.002$	<b><math>1.050 \pm 0.002</math></b>
Gradient Boosting	$1.009 \pm 0.002$	$1.009 \pm 0.002$	<b><math>1.009 \pm 0.002</math></b>
Lasso	$1.050 \pm 0.002$	$1.050 \pm 0.002$	<b><math>1.050 \pm 0.002</math></b>
Linear Regression	$1.022 \pm 0.002$	$1.014 \pm 0.002$	<b><math>1.014 \pm 0.002</math></b>
Random Forest	$1.042 \pm 0.002$	$1.021 \pm 0.002$	<b><math>1.020 \pm 0.002</math></b>
Ridge	$1.022 \pm 0.002$	$1.014 \pm 0.002$	<b><math>1.014 \pm 0.002</math></b>
SVR (RBF)	$1.012 \pm 0.002$	$1.011 \pm 0.002$	<b><math>1.010 \pm 0.002</math></b>
average	1.030	1.024	<b>1.024</b>

Table 23: Bounded Logistic dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$3.831 \pm 0.021$	<b><math>3.302 \pm 0.018</math></b>	$3.417 \pm 0.018$
Gradient Boosting	$3.420 \pm 0.019$	<b><math>3.408 \pm 0.016</math></b>	$3.514 \pm 0.016$
Lasso	$4.035 \pm 0.023$	<b><math>3.691 \pm 0.023</math></b>	$3.773 \pm 0.022$
Linear Regression	$3.288 \pm 0.019$	<b><math>3.120 \pm 0.015</math></b>	$3.247 \pm 0.015$
Random Forest	<b><math>3.449 \pm 0.018</math></b>	$3.466 \pm 0.017$	$3.564 \pm 0.017$
Ridge	$3.288 \pm 0.019$	<b><math>3.120 \pm 0.015</math></b>	$3.247 \pm 0.015$
SVR (RBF)	$3.345 \pm 0.018$	<b><math>3.242 \pm 0.016</math></b>	$3.358 \pm 0.016$
average	3.522	<b>3.336</b>	3.446

Table 24: Linear Gaussian dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$3.549 \pm 0.022$	<b><math>3.289 \pm 0.019</math></b>	$3.431 \pm 0.023$
Gradient Boosting	<b><math>3.228 \pm 0.016</math></b>	$3.342 \pm 0.019$	$3.475 \pm 0.024$
Lasso	$3.755 \pm 0.025$	<b><math>3.562 \pm 0.023</math></b>	$3.661 \pm 0.026$
Linear Regression	<b><math>3.016 \pm 0.015</math></b>	$3.173 \pm 0.018$	$3.333 \pm 0.022$
Random Forest	<b><math>3.289 \pm 0.018</math></b>	$3.375 \pm 0.022$	$3.502 \pm 0.026$
Ridge	<b><math>3.016 \pm 0.015</math></b>	$3.173 \pm 0.018$	$3.333 \pm 0.022$
SVR (RBF)	<b><math>3.256 \pm 0.019</math></b>	$3.307 \pm 0.020$	$3.438 \pm 0.025$
average	<b>3.301</b>	3.317	3.453

Table 25: Nonlinear dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$3.749 \pm 0.023$	<b><math>3.473 \pm 0.019</math></b>	$3.626 \pm 0.024$
Gradient Boosting	<b><math>3.265 \pm 0.017</math></b>	$3.418 \pm 0.020$	$3.579 \pm 0.025$
Lasso	$3.941 \pm 0.026$	<b><math>3.730 \pm 0.023</math></b>	$3.841 \pm 0.027$
Linear Regression	<b><math>3.210 \pm 0.014</math></b>	$3.366 \pm 0.018$	$3.536 \pm 0.023$
Random Forest	<b><math>3.365 \pm 0.019</math></b>	$3.483 \pm 0.022$	$3.634 \pm 0.028$
Ridge	<b><math>3.209 \pm 0.014</math></b>	$3.366 \pm 0.018$	$3.536 \pm 0.023$
SVR (RBF)	<b><math>3.365 \pm 0.019</math></b>	$3.425 \pm 0.020$	$3.580 \pm 0.025$
average	<b>3.444</b>	3.466	3.619

Table 26: Heteroscedastic noise dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$4.677 \pm 0.032$	<b><math>4.474 \pm 0.031</math></b>	$4.599 \pm 0.033$
Gradient Boosting	<b><math>4.545 \pm 0.031</math></b>	$4.580 \pm 0.031$	$4.682 \pm 0.033$
Lasso	$4.837 \pm 0.033$	<b><math>4.688 \pm 0.033</math></b>	$4.771 \pm 0.035$
Linear Regression	<b><math>4.282 \pm 0.029</math></b>	$4.382 \pm 0.030$	$4.524 \pm 0.032$
Random Forest	<b><math>4.522 \pm 0.031</math></b>	$4.572 \pm 0.032$	$4.675 \pm 0.034$
Ridge	<b><math>4.282 \pm 0.029</math></b>	$4.382 \pm 0.030$	$4.524 \pm 0.032$
SVR (RBF)	<b><math>4.479 \pm 0.031</math></b>	$4.498 \pm 0.031$	$4.608 \pm 0.033$
average	4.518	<b>4.511</b>	4.626

Table 27: Heavy-tailed noise dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$5.266 \pm 0.091$	<b><math>5.094 \pm 0.092</math></b>	$5.188 \pm 0.092$
Gradient Boosting	$5.297 \pm 0.095$	<b><math>5.208 \pm 0.093</math></b>	$5.280 \pm 0.093$
Lasso	$5.409 \pm 0.089$	<b><math>5.279 \pm 0.091</math></b>	$5.336 \pm 0.091$
Linear Regression	<b><math>4.949 \pm 0.093</math></b>	$5.018 \pm 0.093$	$5.129 \pm 0.093$
Random Forest	$5.237 \pm 0.093$	<b><math>5.207 \pm 0.094</math></b>	$5.277 \pm 0.093$
Ridge	<b><math>4.949 \pm 0.093</math></b>	$5.018 \pm 0.093$	$5.129 \pm 0.093$
SVR (RBF)	<b><math>5.084 \pm 0.093</math></b>	$5.107 \pm 0.093$	$5.191 \pm 0.093$
average	5.170	<b>5.133</b>	5.218

Table 28: Outlier contamination dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$4.561 \pm 0.102$	<b><math>4.358 \pm 0.105</math></b>	$4.469 \pm 0.103$
Gradient Boosting	$4.488 \pm 0.108$	<b><math>4.457 \pm 0.106</math></b>	$4.542 \pm 0.104$
Lasso	$4.727 \pm 0.099$	<b><math>4.591 \pm 0.103</math></b>	$4.652 \pm 0.101$
Linear Regression	<b><math>4.171 \pm 0.108</math></b>	$4.268 \pm 0.107$	$4.401 \pm 0.105$
Random Forest	$4.480 \pm 0.107$	<b><math>4.479 \pm 0.105</math></b>	$4.555 \pm 0.103$
Ridge	<b><math>4.171 \pm 0.108</math></b>	$4.269 \pm 0.107$	$4.401 \pm 0.105$
SVR (RBF)	<b><math>4.329 \pm 0.106</math></b>	$4.374 \pm 0.106$	$4.474 \pm 0.104$
average	4.418	<b>4.399</b>	4.499

Table 29: Sparse high-dimensional signals dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$3.022 \pm 0.018$	<b><math>2.982 \pm 0.016</math></b>	$2.997 \pm 0.017$
Gradient Boosting	$3.106 \pm 0.017$	<b><math>3.038 \pm 0.017</math></b>	$3.038 \pm 0.018$
Lasso	$3.042 \pm 0.018$	<b><math>3.009 \pm 0.016</math></b>	$3.024 \pm 0.017$
Linear Regression	<b><math>2.978 \pm 0.016</math></b>	$2.995 \pm 0.016$	$3.004 \pm 0.017$
Random Forest	$3.064 \pm 0.017$	<b><math>3.030 \pm 0.017</math></b>	$3.032 \pm 0.019$
Ridge	<b><math>2.978 \pm 0.016</math></b>	$2.995 \pm 0.016$	$3.004 \pm 0.017$
SVR (RBF)	$3.041 \pm 0.017$	<b><math>3.030 \pm 0.017</math></b>	$3.033 \pm 0.019$
average	3.033	<b>3.011</b>	3.019

Table 30: Covariate shift dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$3.954 \pm 0.064$	<b><math>3.693 \pm 0.067</math></b>	$3.934 \pm 0.082$
Gradient Boosting	<b><math>3.490 \pm 0.029</math></b>	$3.832 \pm 0.069$	$4.051 \pm 0.084$
Lasso	$4.273 \pm 0.080$	<b><math>4.054 \pm 0.072</math></b>	$4.234 \pm 0.085$
Linear Regression	<b><math>3.040 \pm 0.018</math></b>	$3.536 \pm 0.067$	$3.808 \pm 0.083$
Random Forest	<b><math>3.624 \pm 0.045</math></b>	$3.881 \pm 0.071$	$4.097 \pm 0.085$
Ridge	<b><math>3.040 \pm 0.018</math></b>	$3.536 \pm 0.067$	$3.808 \pm 0.083$
SVR (RBF)	<b><math>4.151 \pm 0.086</math></b>	$4.179 \pm 0.089$	$4.342 \pm 0.099$
average	<b>3.653</b>	3.816	4.039

Table 31: Friedman 1 dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$4.385 \pm 0.023$	<b>4.055</b> $\pm 0.021$	$4.269 \pm 0.022$
Gradient Boosting	<b>3.378</b> $\pm 0.018$	$3.660 \pm 0.020$	$3.934 \pm 0.022$
Lasso	$4.365 \pm 0.023$	<b>4.157</b> $\pm 0.021$	$4.359 \pm 0.023$
Linear Regression	<b>3.893</b> $\pm 0.021$	$4.043 \pm 0.021$	$4.256 \pm 0.022$
Random Forest	<b>3.568</b> $\pm 0.019$	$3.808 \pm 0.020$	$4.061 \pm 0.022$
Ridge	<b>3.893</b> $\pm 0.021$	$4.043 \pm 0.021$	$4.256 \pm 0.022$
SVR (RBF)	<b>3.746</b> $\pm 0.021$	$3.865 \pm 0.021$	$4.104 \pm 0.023$
average	<b>3.890</b>	3.947	4.177

Table 32: Friedman 2 dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$181.546 \pm 1.247$	<b>151.625</b> $\pm 1.044$	$181.004 \pm 1.604$
Gradient Boosting	<b>23.936</b> $\pm 0.181$	$142.035 \pm 0.940$	$171.828 \pm 1.657$
Lasso	<b>137.739</b> $\pm 0.784$	$151.686 \pm 1.032$	$181.067 \pm 1.600$
Linear Regression	<b>137.743</b> $\pm 0.784$	$151.713 \pm 1.033$	$181.083 \pm 1.600$
Random Forest	<b>20.166</b> $\pm 0.202$	$141.653 \pm 0.914$	$171.544 \pm 1.636$
Ridge	<b>137.746</b> $\pm 0.784$	$151.711 \pm 1.033$	$181.082 \pm 1.600$
SVR (RBF)	$345.749 \pm 2.419$	<b>170.491</b> $\pm 1.685$	$193.626 \pm 2.036$
average	<b>140.661</b>	151.559	180.176

Table 33: Friedman 3 dataset ( $n = 1000$  and  $\sigma = 3$ )

	none	CVAR1	CVAR10
Elastic Net	$3.019 \pm 0.015$	<b>3.019</b> $\pm 0.015$	$3.020 \pm 0.015$
Gradient Boosting	$3.130 \pm 0.015$	$3.039 \pm 0.015$	<b>3.024</b> $\pm 0.015$
Lasso	$3.019 \pm 0.015$	<b>3.019</b> $\pm 0.015$	$3.020 \pm 0.015$
Linear Regression	<b>3.018</b> $\pm 0.015$	$3.036 \pm 0.015$	$3.019 \pm 0.015$
Random Forest	$3.146 \pm 0.016$	$3.045 \pm 0.015$	<b>3.028</b> $\pm 0.015$
Ridge	<b>3.018</b> $\pm 0.015$	$3.036 \pm 0.015$	$3.019 \pm 0.015$
SVR (RBF)	$3.054 \pm 0.016$	$3.036 \pm 0.015$	<b>3.022</b> $\pm 0.015$
average	3.058	3.033	<b>3.022</b>

Table 34: Bounded Logistic dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.580 \pm 0.022$	<b><math>1.749 \pm 0.016</math></b>	$1.849 \pm 0.015$
Gradient Boosting	<b><math>1.708 \pm 0.021</math></b>	$1.770 \pm 0.016$	$1.869 \pm 0.015$
Lasso	$2.877 \pm 0.026$	<b><math>2.390 \pm 0.026</math></b>	$2.454 \pm 0.025$
Linear Regression	$1.630 \pm 0.023$	<b><math>1.380 \pm 0.007</math></b>	$1.506 \pm 0.007$
Random Forest	<b><math>1.870 \pm 0.023</math></b>	$1.953 \pm 0.019$	$2.042 \pm 0.018$
Ridge	$1.629 \pm 0.023$	<b><math>1.380 \pm 0.007</math></b>	$1.506 \pm 0.007$
SVR (RBF)	$1.578 \pm 0.017$	<b><math>1.504 \pm 0.008</math></b>	$1.619 \pm 0.008$
average	1.982	<b>1.732</b>	1.835

Table 35: Linear Gaussian dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.141 \pm 0.029$	<b><math>1.744 \pm 0.024</math></b>	$1.956 \pm 0.031$
Gradient Boosting	<b><math>1.326 \pm 0.015</math></b>	$1.676 \pm 0.028$	$1.899 \pm 0.035$
Lasso	$2.471 \pm 0.032$	<b><math>2.189 \pm 0.030</math></b>	$2.335 \pm 0.035$
Linear Regression	<b><math>1.005 \pm 0.005</math></b>	$1.529 \pm 0.025$	$1.772 \pm 0.032$
Random Forest	<b><math>1.566 \pm 0.024</math></b>	$1.797 \pm 0.033$	$2.001 \pm 0.039$
Ridge	<b><math>1.005 \pm 0.005</math></b>	$1.529 \pm 0.025$	$1.772 \pm 0.032$
SVR (RBF)	<b><math>1.418 \pm 0.020</math></b>	$1.651 \pm 0.027$	$1.879 \pm 0.035$
average	<b>1.562</b>	1.731	1.945

Table 36: Nonlinear dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.469 \pm 0.028$	<b><math>2.069 \pm 0.023</math></b>	$2.284 \pm 0.031$
Gradient Boosting	<b><math>1.430 \pm 0.016</math></b>	$1.835 \pm 0.030$	$2.087 \pm 0.036$
Lasso	$2.755 \pm 0.032$	<b><math>2.458 \pm 0.030</math></b>	$2.613 \pm 0.035$
Linear Regression	<b><math>1.496 \pm 0.008</math></b>	$1.896 \pm 0.023$	$2.139 \pm 0.031$
Random Forest	<b><math>1.735 \pm 0.025</math></b>	$2.004 \pm 0.033$	$2.230 \pm 0.039$
Ridge	<b><math>1.496 \pm 0.008</math></b>	$1.896 \pm 0.023$	$2.139 \pm 0.031$
SVR (RBF)	<b><math>1.622 \pm 0.020</math></b>	$1.870 \pm 0.028$	$2.118 \pm 0.035$
average	<b>1.858</b>	2.004	2.230

Table 37: Heteroscedastic noise dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.379 \pm 0.027$	<b><math>2.014 \pm 0.023</math></b>	$2.221 \pm 0.030$
Gradient Boosting	<b><math>1.697 \pm 0.014</math></b>	$1.977 \pm 0.026$	$2.192 \pm 0.033$
Lasso	$2.683 \pm 0.030$	<b><math>2.422 \pm 0.029</math></b>	$2.564 \pm 0.033$
Linear Regression	<b><math>1.427 \pm 0.010</math></b>	$1.826 \pm 0.023$	$2.062 \pm 0.030$
Random Forest	<b><math>1.875 \pm 0.022</math></b>	$2.069 \pm 0.030$	$2.267 \pm 0.036$
Ridge	<b><math>1.427 \pm 0.010</math></b>	$1.826 \pm 0.023$	$2.062 \pm 0.030$
SVR (RBF)	<b><math>1.777 \pm 0.019</math></b>	$1.945 \pm 0.026$	$2.163 \pm 0.033$
average	<b>1.895</b>	2.011	2.219

Table 38: Heavy-tailed noise dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.492 \pm 0.034$	<b><math>2.140 \pm 0.033</math></b>	$2.335 \pm 0.037$
Gradient Boosting	<b><math>1.907 \pm 0.032</math></b>	$2.122 \pm 0.036$	$2.321 \pm 0.040$
Lasso	$2.783 \pm 0.037$	<b><math>2.523 \pm 0.036</math></b>	$2.657 \pm 0.039$
Linear Regression	<b><math>1.650 \pm 0.031</math></b>	$1.970 \pm 0.035$	$2.191 \pm 0.038$
Random Forest	<b><math>2.043 \pm 0.034</math></b>	$2.196 \pm 0.039$	$2.384 \pm 0.042$
Ridge	<b><math>1.650 \pm 0.031</math></b>	$1.970 \pm 0.035$	$2.191 \pm 0.038$
SVR (RBF)	<b><math>1.932 \pm 0.034</math></b>	$2.078 \pm 0.036$	$2.279 \pm 0.039$
average	<b>2.065</b>	2.143	2.337

Table 39: Outlier contamination dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.329 \pm 0.032$	<b><math>1.952 \pm 0.032</math></b>	$2.159 \pm 0.034$
Gradient Boosting	<b><math>1.684 \pm 0.034</math></b>	$1.929 \pm 0.036$	$2.138 \pm 0.037$
Lasso	$2.634 \pm 0.034$	<b><math>2.375 \pm 0.035</math></b>	$2.513 \pm 0.036$
Linear Regression	<b><math>1.390 \pm 0.036</math></b>	$1.768 \pm 0.035$	$2.004 \pm 0.036$
Random Forest	<b><math>1.852 \pm 0.035</math></b>	$2.020 \pm 0.037$	$2.217 \pm 0.039$
Ridge	<b><math>1.390 \pm 0.036</math></b>	$1.768 \pm 0.035$	$2.004 \pm 0.036$
SVR (RBF)	<b><math>1.707 \pm 0.034</math></b>	$1.882 \pm 0.036$	$2.099 \pm 0.037$
average	<b>1.855</b>	1.956	2.162

Table 40: Sparse high-dimensional signals dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$1.153 \pm 0.016$	<b>1.062</b> $\pm 0.011$	$1.105 \pm 0.016$
Gradient Boosting	<b>1.039</b> $\pm 0.006$	$1.072 \pm 0.011$	$1.112 \pm 0.017$
Lasso	$1.204 \pm 0.018$	<b>1.134</b> $\pm 0.013$	$1.167 \pm 0.017$
Linear Regression	<b>0.993</b> $\pm 0.005$	$1.049 \pm 0.011$	$1.093 \pm 0.016$
Random Forest	<b>1.030</b> $\pm 0.006$	$1.070 \pm 0.012$	$1.110 \pm 0.017$
Ridge	<b>0.993</b> $\pm 0.005$	$1.049 \pm 0.011$	$1.093 \pm 0.016$
SVR (RBF)	<b>1.073</b> $\pm 0.009$	$1.091 \pm 0.014$	$1.128 \pm 0.019$
average	<b>1.069</b>	1.075	1.115

Table 41: Covariate shift dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$2.706 \pm 0.080$	<b>2.373</b> $\pm 0.088$	$2.688 \pm 0.105$
Gradient Boosting	<b>1.717</b> $\pm 0.037$	$2.395 \pm 0.090$	$2.705 \pm 0.106$
Lasso	$3.149 \pm 0.095$	<b>2.918</b> $\pm 0.089$	$3.136 \pm 0.104$
Linear Regression	<b>1.013</b> $\pm 0.006$	$2.128 \pm 0.091$	$2.491 \pm 0.108$
Random Forest	<b>2.152</b> $\pm 0.060$	$2.582 \pm 0.092$	$2.859 \pm 0.107$
Ridge	<b>1.013</b> $\pm 0.006$	$2.128 \pm 0.091$	$2.491 \pm 0.108$
SVR (RBF)	<b>2.740</b> $\pm 0.101$	$2.936 \pm 0.110$	$3.126 \pm 0.120$
average	<b>2.070</b>	2.494	2.785

Table 42: Friedman 1 dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$3.330 \pm 0.015$	<b>2.926</b> $\pm 0.014$	$3.175 \pm 0.015$
Gradient Boosting	<b>1.551</b> $\pm 0.009$	$2.243 \pm 0.012$	$2.583 \pm 0.015$
Lasso	$3.301 \pm 0.016$	<b>3.051</b> $\pm 0.016$	$3.289 \pm 0.016$
Linear Regression	<b>2.634</b> $\pm 0.014$	$2.907 \pm 0.014$	$3.155 \pm 0.015$
Random Forest	<b>2.039</b> $\pm 0.011$	$2.525 \pm 0.012$	$2.830 \pm 0.014$
Ridge	<b>2.633</b> $\pm 0.014$	$2.907 \pm 0.014$	$3.155 \pm 0.015$
SVR (RBF)	<b>2.172</b> $\pm 0.013$	$2.530 \pm 0.014$	$2.831 \pm 0.015$
average	<b>2.523</b>	2.727	3.003

Table 43: Friedman 2 dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$181.524 \pm 1.246$	<b><math>151.614 \pm 1.044</math></b>	$180.970 \pm 1.603$
Gradient Boosting	<b><math>23.714 \pm 0.181</math></b>	$142.014 \pm 0.949$	$171.767 \pm 1.660$
Lasso	<b><math>137.715 \pm 0.782</math></b>	$151.681 \pm 1.031$	$181.033 \pm 1.599$
Linear Regression	<b><math>137.720 \pm 0.782</math></b>	$151.707 \pm 1.032$	$181.048 \pm 1.599$
Random Forest	<b><math>19.889 \pm 0.206</math></b>	$141.642 \pm 0.918$	$171.495 \pm 1.637$
Ridge	<b><math>137.723 \pm 0.782</math></b>	$151.706 \pm 1.032$	$181.049 \pm 1.599$
SVR (RBF)	$345.751 \pm 2.424$	<b><math>170.534 \pm 1.669</math></b>	$193.720 \pm 2.034$
average	<b>140.576</b>	151.557	180.155

Table 44: Friedman 3 dataset ( $n = 1000$  and  $\sigma = 1$ )

	none	CVAR1	CVAR10
Elastic Net	$1.048 \pm 0.005$	$1.048 \pm 0.005$	<b><math>1.048 \pm 0.005</math></b>
Gradient Boosting	$1.045 \pm 0.005$	$1.029 \pm 0.005$	<b><math>1.028 \pm 0.005</math></b>
Lasso	$1.048 \pm 0.005$	$1.048 \pm 0.005$	<b><math>1.048 \pm 0.005</math></b>
Linear Regression	$1.024 \pm 0.005$	<b><math>1.021 \pm 0.005</math></b>	$1.022 \pm 0.005$
Random Forest	$1.054 \pm 0.005$	$1.031 \pm 0.005$	<b><math>1.030 \pm 0.005</math></b>
Ridge	$1.024 \pm 0.005$	<b><math>1.021 \pm 0.005</math></b>	$1.022 \pm 0.005$
SVR (RBF)	$1.035 \pm 0.005$	<b><math>1.026 \pm 0.005</math></b>	$1.026 \pm 0.005$
average	1.040	1.032	<b>1.032</b>

Table 45: Bias correction dataset

	none	CVAR1	CVAR10
Elastic Net	$1.836 \pm 0.003$	<b><math>1.590 \pm 0.003</math></b>	$1.605 \pm 0.003$
Gradient Boosting	<b><math>1.208 \pm 0.002</math></b>	$1.272 \pm 0.003$	$1.292 \pm 0.003$
Lasso	$1.977 \pm 0.004$	<b><math>1.738 \pm 0.003</math></b>	$1.751 \pm 0.003$
Linear Regression	<b><math>1.463 \pm 0.003</math></b>	$1.507 \pm 0.003$	$1.523 \pm 0.003$
Random Forest	<b><math>0.979 \pm 0.002</math></b>	$1.078 \pm 0.003$	$1.102 \pm 0.003$
Ridge	<b><math>1.463 \pm 0.003</math></b>	$1.507 \pm 0.003$	$1.523 \pm 0.003$
SVR (RBF)	<b><math>1.152 \pm 0.003</math></b>	$1.235 \pm 0.003$	$1.256 \pm 0.003$
average	1.440	<b>1.418</b>	1.436

Table 46: Wine quality dataset

	none	CVAR1	CVAR10
Elastic Net	<b><math>0.870 \pm 0.001</math></b>	$0.870 \pm 0.001$	$0.870 \pm 0.001$
Gradient Boosting	<b><math>0.681 \pm 0.001</math></b>	$0.684 \pm 0.001$	$0.684 \pm 0.001$
Lasso	<b><math>0.870 \pm 0.001</math></b>	$0.870 \pm 0.001$	$0.870 \pm 0.001$
Linear Regression	$0.733 \pm 0.001$	<b><math>0.732 \pm 0.001</math></b>	$0.732 \pm 0.001$
Random Forest	<b><math>0.604 \pm 0.002</math></b>	$0.618 \pm 0.002$	$0.619 \pm 0.002$
Ridge	$0.733 \pm 0.001$	<b><math>0.732 \pm 0.001</math></b>	$0.732 \pm 0.001$
SVR (RBF)	<b><math>0.676 \pm 0.001</math></b>	$0.677 \pm 0.001$	$0.678 \pm 0.001$
average	<b>0.738</b>	0.740	0.741

colour, red or white), Airfoil self-noise (Brooks et al., 1989) with 1503 examples and 5 features, and Student performance (Cortez, 2008) with 649 examples and 30 features. The results are mixed; in two cases, our methods slightly improve the performance of the base algorithms on average. (One of the datasets where our methods do not improve, and even slightly lower, the quality of predictions is the Wine quality dataset in Table 46; for this data set the label variable is bounded but with bounds that are far from typical values of the labels.)

## C Code Availability

Python code for reproducing our experiments in Sect. 8 and Appendix B is available at the GitHub repository <https://github.com/ip200/ivar-experiments> (Petej, 2026).

Table 47: Airfoil self-noise dataset

	none	CVAR1	CVAR10
Elastic Net	$5.616 \pm 0.018$	<b><math>4.865 \pm 0.020</math></b>	$4.992 \pm 0.019$
Gradient Boosting	<b><math>2.658 \pm 0.016</math></b>	$3.146 \pm 0.017$	$3.423 \pm 0.017$
Lasso	$5.541 \pm 0.019$	<b><math>5.072 \pm 0.022</math></b>	$5.182 \pm 0.021$
Linear Regression	$4.820 \pm 0.021$	<b><math>4.670 \pm 0.018</math></b>	$4.814 \pm 0.018$
Random Forest	<b><math>1.767 \pm 0.014</math></b>	$2.626 \pm 0.016$	$2.957 \pm 0.017$
Ridge	$4.820 \pm 0.020$	<b><math>4.670 \pm 0.018</math></b>	$4.814 \pm 0.018$
SVR (RBF)	<b><math>3.797 \pm 0.019</math></b>	$4.097 \pm 0.018$	$4.299 \pm 0.018$
average	<b>4.145</b>	4.164	4.354

Table 48: Student performance dataset

	none	CVAR1	CVAR10
Elastic Net	<b><math>241.009 \pm 0.626</math></b>	$241.617 \pm 0.632$	$241.694 \pm 0.640$
Gradient Boosting	<b><math>231.824 \pm 0.613</math></b>	$232.643 \pm 0.638$	$232.996 \pm 0.646$
Lasso	<b><math>240.262 \pm 0.633</math></b>	$240.960 \pm 0.641$	$241.154 \pm 0.648$
Linear Regression	<b><math>241.279 \pm 0.649</math></b>	$241.408 \pm 0.635$	$241.560 \pm 0.643$
Random Forest	$235.259 \pm 0.691$	<b><math>234.641 \pm 0.704</math></b>	$234.992 \pm 0.702$
Ridge	<b><math>240.821 \pm 0.636</math></b>	$241.245 \pm 0.634$	$241.422 \pm 0.644$
SVR (RBF)	$256.059 \pm 0.745$	$245.013 \pm 0.653$	<b><math>244.957 \pm 0.655</math></b>
average	240.931	<b>239.647</b>	239.825